



Development of a linguistic support model for information retrieval for cloud library systems

Shakenov Nurzhan - a58942

Dissertation submitted to School of Technology and Management of Bragança to obtain the Degree of Master in Informatics under the scope of the double diploma with Dulaty University

Supervisors:

Prof. Rui Pedro Lopes

Prof. Aigul Tungatarova

Bragança

2023-2024

Abstract

This dissertation addresses the limitations of traditional keyword-based search methods in cloud library systems by developing a robust linguistic support model. Leveraging advanced techniques in text extraction, embedding generation, and semantic search, this study aims to enhance the accuracy and relevance of search results. Document ingestion and text extraction were performed using PyMuPDF, ensuring high-quality data for subsequent processes. Text embeddings generated by LangChain's Mistral model were stored in the Chroma vector database, facilitating efficient retrieval. A user-friendly interface developed with Flask enabled seamless user interaction.

The project faced challenges such as API key requirements for GPT-2, text extraction accuracy, and large-scale data handling, which were addressed through alternative tools and methodologies. The results demonstrate significant improvements in search accuracy and relevance, aligning with recent advancements in NLP. Future work will focus on enhancing data preprocessing, expanding datasets, and integrating more advanced search algorithms. This study contributes valuable insights into the practical application of NLP techniques in cloud library systems, offering a foundation for further research and development in the field.

Keywords: Cloud Library Systems, Information Retrieval, Text Extraction, Embedding Generation, Semantic Search, Natural Language Processing.

Resumo

Esta dissertação aborda as limitações dos métodos tradicionais de busca por palavras-chave em sistemas de bibliotecas na nuvem, desenvolvendo um modelo robusto de suporte linguístico. Aproveitando técnicas avançadas de extração de texto, geração de embeddings e busca semântica, este estudo visa melhorar a precisão e a relevância dos resultados de busca. A ingestão de documentos e a extração de texto foram realizadas utilizando PyMuPDF, garantindo dados de alta qualidade para os processos subsequentes. Embeddings de texto gerados pelo modelo Mistral da LangChain foram armazenados na base de dados vetorial Chroma, facilitando a recuperação eficiente. Uma interface amigável desenvolvida com Flask permitiu uma interação perfeita do utilizador.

O projeto enfrentou desafios como a necessidade de chave API para GPT-2, precisão na extração de texto e manipulação de dados em grande escala, que foram solucionados através de ferramentas e metodologias alternativas. Os resultados demonstram melhorias significativas na precisão e relevância da busca, alinhando-se com os avanços recentes em PLN. Trabalhos futuros se concentrarão em aprimorar o pré-processamento de dados, expandir conjuntos de dados e integrar algoritmos de busca mais avançados. Este estudo contribui com insights valiosos sobre a aplicação prática de técnicas de PLN em sistemas de bibliotecas na nuvem, oferecendo uma base para futuras pesquisas e desenvolvimento na área.

Palavras-chave: Sistemas de Bibliotecas na Nuvem, Recuperação de Informação, Extração de Texto, Geração de Embeddings, Busca Semântica, Processamento de Linguagem Natural.

Contents

1	Introduction	1
1.1	Problem Statement	2
1.2	Research Questions	2
1.3	Objectives	3
1.4	Significance of the Study	4
1.5	Scope and Limitations	4
2	Context and Technologies	5
2.1	Concepts and Technology	6
2.2	Literature Review	7
2.3	Practical Tools and Justification	8
3	Approach and Methodology	11
3.1	Research Design	12
3.2	Data Collection Methods	12
3.2.1	Initial Approach with GPT-2	13
3.2.2	Alternative Approaches for Text Extraction	14
3.3	Text Embedding and Vectorization	15
3.3.1	Embedding Generation	16
3.3.2	Vector Storage	16
3.4	Query Processing and Semantic Search	17
3.4.1	Query Embedding	18

3.4.2	Semantic Search Algorithm	19
3.5	User Interface Development	20
3.5.1	Web Interface Design	20
3.5.2	Backend Integration	21
3.6	Tasks and System Interactions	22
3.7	Challenges Faced	24
4	Development	27
4.1	Text Extraction and Cleaning	28
4.1.1	Implementation with PyMuPDF	29
4.1.2	Text Cleaning	30
4.2	Embedding Generation and Vector Storage	30
4.2.1	Embedding Generation with LangChain	31
4.2.2	Storing Embeddings in Chroma	32
4.3	Query Processing and Semantic Search	33
4.3.1	Query Embedding	33
4.3.2	Semantic Search Algorithm	34
4.4	User Interface Development	35
4.4.1	Web Interface Design	36
4.4.2	Backend Integration	38
4.5	Challenges and Technical Solutions	40
4.6	Performance Optimization	42
4.7	Summary	44
5	Results and Discussion	47
5.1	Document Ingestion and Text Extraction	47
5.2	Embedding Generation and Storage	49
5.3	Query Processing and Semantic Search	50
5.4	User Interface and User Experience	54
5.5	Future Work and Unmet Objectives	55

5.6 Challenges Faced	55
6 Conclusions	57

List of Figures

3.1	Limits of using API keys GPT	13
3.2	From PDF to TXT format	15
3.3	Process of Embedding	16
3.4	Process of vectorizing data	17
3.5	RAG Embedding	19
3.6	Simple User Interface	21
5.1	Example of result converting from PDF to TXT	49
5.2	Results of embedding for 1 txt file	50
5.3	The result of 1st query	51
5.4	The result of 2nd query	51
5.5	The result of 3rd query	52
5.6	The result of 4th query	52
5.7	The result of 5th query	53
5.8	The Interface	54

Chapter 1

Introduction

The integration of cloud technologies across various sectors has dramatically altered the way information is accessed and managed. Notably, libraries have started embracing cloud-based systems to elevate their services and offer users effortless access to information. Cloud library systems epitomize a form of automated library and information system that leverages cloud technologies to store, manage, and retrieve digital content. These innovative systems empower users to access library resources from any location with internet connectivity, thereby dismantling the physical barriers traditionally associated with libraries [1].

Yet, the shift to cloud library systems is fraught with challenges. A prominent issue is the insufficiency of current search functionalities within these systems. Many cloud library systems depend on rudimentary full-text search capabilities or searches confined to a few database fields. Consequently, these limitations often culminate in inefficient information retrieval processes, where users grapple with locating relevant information swiftly and accurately [2]. This predicament is further intensified by the absence of advanced linguistic support that could bolster the precision and relevance of search results [3].

Despite the hurdles, cloud technologies hold the promise to revolutionize information retrieval in libraries. However, realizing this potential hinges on addressing the unique needs of library systems. A pivotal aspect of this endeavor is the development of sophisticated linguistic support models. These models can significantly enhance thematic search

capabilities, thereby simplifying the process for users to find the information they seek.

While cloud technologies offer transformative potential for libraries, their successful implementation requires overcoming significant challenges, particularly in improving search functionalities and incorporating advanced linguistic support. This dual focus can unlock a new era of efficient, accessible information retrieval for library users.

1.1 Problem Statement

General-purpose language models, while powerful, often struggle to capture the intricacies of specialized domains. This limitation arises from the models' pre-training on vast but generalized corpora, which may not include the specific terminology and contextual knowledge required for domain-specific applications. Consequently, there is a need for approaches that can effectively adapt these models to specialized contexts, ensuring that they generate relevant, accurate, and coherent text.

The primary focus of this dissertation is to investigate methods for enhancing information retrieval in cloud library systems by developing a linguistic support model. By leveraging advanced techniques in text extraction, embedding generation, and semantic search, this research aims to improve the accuracy and relevance of search results in digital libraries. The project involves fine-tuning pre-trained language models and integrating them with robust data processing and retrieval frameworks to meet the demands of specialized fields within cloud-based library systems [4].

1.2 Research Questions

To tackle this issue, the research aims to answer the following questions:

- What are the existing methods for organizing information searches in cloud library systems?
- What is the current state of linguistic support for information retrieval in these

systems?

- How can cloud technologies be effectively utilized to enhance information retrieval in library systems?
- What theoretical frameworks can support the development of a thematic search model for cloud library systems?

1.3 Objectives

The primary objective of this study is to improve information retrieval processes in cloud library systems by developing a robust linguistic support model. The specific objectives include:

- Assess the current state of linguistic support and semantic search capabilities in cloud library systems to understand their limitations and potential enhancements.
- Clarify and delineate the concept of cloud library systems, outlining their structure, functionalities, and the role they play in modern information retrieval.
- Design a thematic search model tailored for cloud library systems, supported by a theoretical framework that justifies its components and expected performance.
- Establish and validate the essential requirements for a thematic search model, ensuring it meets the specific needs of cloud library environments.
- Create advanced linguistic support mechanisms that enhance the performance of thematic search models, improving their accuracy and relevance in retrieving information.

These objectives are designed to ensure that the developed thematic search model not only improves information retrieval in cloud library systems but also sets a foundation for future enhancements in the field.

1.4 Significance of the Study

This research holds significant implications for both academia and industry. For academia, it advances the understanding of how linguistic support models and semantic search techniques can be adapted to meet the needs of specialized fields, contributing to the theoretical foundations of information retrieval and natural language processing (NLP). For industry, the findings can inform the development of more effective cloud library systems and NLP applications, from enhanced information retrieval processes to domain-specific search engines and virtual assistants.

Moreover, this dissertation seeks to bridge the gap between general-purpose NLP research and its practical applications in specialized domains such as cloud library systems [5]. By providing a comprehensive analysis of the adaptation process and its outcomes, this study aims to serve as a valuable resource for researchers and practitioners alike. The insights gained from this research can guide the development of more accurate and relevant information retrieval models, ultimately improving user experience and access to information in digital libraries [4].

1.5 Scope and Limitations

The scope of this study includes the analysis and development of a linguistic support model specifically for cloud library systems. It encompasses evaluating current search methodologies and linguistic support mechanisms and integrating these findings into a comprehensive thematic search model [1]. However, the study is limited to cloud library systems and does not extend to other types of digital libraries or information systems. Additionally, practical implementation and testing of the developed model are constrained to selected case studies and may require further validation in different contexts [4].

Chapter 2

Context and Technologies

The increasing adoption of cloud technologies in libraries has transformed how information is accessed and managed. Cloud library systems offer scalable, cost-effective, and accessible solutions for storing, managing, and retrieving digital content. However, these systems face significant challenges, particularly in the domain of information retrieval. Traditional keyword-based searches are inadequate for handling the vast and unstructured datasets prevalent in digital libraries, necessitating the development of advanced thematic search capabilities that can meet users' complex information needs [6]. The limitations of conventional search methods become more pronounced as the volume of digital content grows, leading to inefficiencies and user frustration. This inadequacy is evident in the inability of simple keyword searches to understand the context and semantics of user queries, often resulting in irrelevant or incomplete search results.

Moreover, the dynamic nature of digital content in cloud library systems requires search functionalities that can adapt to varying contexts and user requirements. Current information retrieval systems often fail to provide the necessary depth and precision, particularly for academic and research purposes where the specificity and relevance of information are paramount. The need for advanced linguistic models that can perform semantic analysis and thematic searches is critical for improving user experience and satisfaction. Such models can significantly enhance the ability to retrieve contextually relevant information, thereby optimizing the efficiency and effectiveness of cloud library

systems. Addressing these challenges requires a comprehensive approach that integrates advanced linguistic support with robust search algorithms to meet the evolving demands of digital library users [7].

2.1 Concepts and Technology

Cloud library systems leverage cloud computing to provide a flexible and scalable infrastructure for managing digital resources. These systems enable libraries to offer remote access to their collections, thereby increasing the accessibility and usability of their resources. Key components of cloud library systems include digital repositories, user interfaces for accessing content, and search functionalities designed to facilitate efficient information retrieval [8]. By migrating to cloud-based platforms, libraries can reduce costs, improve resource sharing, and enhance service delivery. The flexibility of cloud systems allows for seamless updates and maintenance, reducing downtime and ensuring that users have continuous access to the latest resources [6].

Information retrieval (IR) involves the process of obtaining information from large collections of unstructured data, typically text. Traditional IR systems in libraries relied on structured metadata and controlled vocabularies, which, while effective for smaller collections, are insufficient for the expansive and diverse datasets found in digital libraries. Modern IR systems often employ full-text search capabilities; however, these systems struggle to handle natural language queries effectively, particularly when it comes to thematic searches. This limitation highlights the need for more sophisticated IR techniques that can understand and process the semantics of user queries and documents [9]. Advanced IR systems must be able to parse complex queries, understand user intent, and retrieve documents that are contextually relevant, which requires the integration of linguistic and semantic analysis [10].

Thematic search refers to the ability to retrieve information based on the thematic content rather than simple keyword matches. This approach requires advanced linguistic

support to understand the context and semantics of both queries and documents. Linguistic models, particularly those based on machine learning and NLP, have shown significant promise in enhancing thematic search capabilities. These models analyze textual data to identify themes, topics, and relevant information, thereby improving the accuracy and relevance of search results [11]. The integration of these models into cloud library systems can significantly enhance their information retrieval performance, addressing a critical need in the field. Additionally, the use of thematic search can improve user satisfaction by providing more precise and relevant results, reducing the time and effort required to find specific information.

2.2 Literature Review

Several studies have explored the adoption and impact of cloud technologies in libraries. Chandran [1] discuss the initiatives taken by Indian public and academic libraries to implement cloud computing, highlighting benefits such as improved resource sharing and reduced costs. The study emphasizes the role of cloud computing in enhancing the operational efficiency of libraries, allowing them to better serve their patrons. Prasad [2] elaborates on the journey of libraries towards adopting cloud computing, emphasizing the enhanced accessibility and operational efficiency provided by cloud-based systems. These studies underscore the transformative potential of cloud technologies but also highlight the need for advanced search functionalities to fully realize these benefits.

Mittal, Patel and Khan [5] addresses the specific challenges of information retrieval in cloud environments, noting that traditional IR techniques are insufficient for handling the complexity and scale of digital libraries. The study identifies the need for advanced algorithms and models that can provide more accurate and contextually relevant search results. Traditional IR systems often rely on keyword matching, which can lead to irrelevant results if the exact terms are not used. Rowley [3] discusses the limitations of the existing DIKW (Data, Information, Knowledge, Wisdom) hierarchy in the context of digital information management, suggesting that new frameworks are needed to support

effective information retrieval in cloud-based systems. The study calls for the integration of semantic search capabilities that can understand the context and meaning behind user queries, thus providing more relevant and accurate search results.

Recent advancements in machine learning and NLP have led to the development of sophisticated linguistic models that significantly enhance thematic search capabilities. Scholz, Ramirez-Corona and Flores [4] propose next-generation library systems that integrate cloud technologies with advanced IR techniques, including the use of linguistic models for better search and retrieval. These models analyze the semantic content of documents and queries, enabling more precise and relevant search results. The integration of these models into cloud library systems represents a significant step forward in addressing the limitations of traditional IR approaches. By leveraging the power of linguistic models, libraries can provide more nuanced and contextually aware search capabilities, improving the overall user experience.

2.3 Practical Tools and Justification

PyMuPDF (also known as fitz) is a library used to read, manipulate, and extract information from PDF documents. It is particularly useful in this project for converting PDFs to text, determining the layout of the text (single-column or multi-column), and ensuring that the text is clean and suitable for further processing. PyMuPDF's robust capabilities make it an ideal choice for handling the diverse and complex PDF documents typically found in digital libraries. This tool allows for efficient extraction and preprocessing of text data, which is crucial for building accurate thematic search models [12]. PyMuPDF's ability to handle various PDF structures and extract text with high precision makes it indispensable for this project, ensuring that the textual data is clean and ready for further analysis.

LangChain is a framework designed to build applications with language models. In this project, LangChain is used to create text embeddings using the Mistral model, store these embeddings in a vector database (Chroma), and perform efficient retrieval of relevant

documents. The use of LangChain and Chroma is justified by their ability to handle large-scale text data and provide high-performance search capabilities. These tools are essential for developing an effective thematic search model that can enhance information retrieval in cloud library systems [13]. LangChain's flexibility and scalability allow it to handle large datasets efficiently, making it an ideal choice for embedding and retrieving textual information in this project.

Flask is a lightweight web framework used to create the web interface for the project. It allows users to input queries and receive answers generated by the retrieval system. Flask's simplicity and flexibility make it a suitable choice for developing a user-friendly interface that can interact with the backend systems efficiently. This framework facilitates seamless communication between the user interface and the information retrieval system, ensuring a smooth user experience [14]. Flask's robust ecosystem and ease of integration with other tools make it an excellent choice for building the web interface, providing users with an intuitive and responsive platform for interacting with the information retrieval system.

Chapter 3

Approach and Methodology

The primary problem addressed by this research is the inadequacy of current cloud library systems to provide effective thematic search capabilities. Traditional IR techniques, which rely heavily on keyword-based searches, are insufficient for handling the complex and unstructured data typically found in digital libraries [6]. This inadequacy results in irrelevant or incomplete search results, thereby failing to meet users' needs for accurate and contextually relevant information[7]. As the volume and variety of digital content continue to grow, the limitations of conventional search methods become more pronounced, leading to inefficiencies and user frustration [8].

Moreover, the dynamic nature of user queries, often involving natural language questions, adds another layer of complexity. Existing systems struggle to interpret and respond to such queries accurately, leading to user frustration and decreased satisfaction [9]. The need for a more sophisticated approach that integrates advanced linguistic support and semantic analysis is evident. This study proposes the development of a linguistic support model that leverages NLP and machine learning techniques to enhance the thematic search capabilities of cloud library systems [10].

3.1 Research Design

The research design for this study involves several key phases, each aimed at addressing specific aspects of the problem related to thematic search in cloud library systems. These phases include document processing and text extraction, text embedding and vectorization, query processing and semantic search, and user interface development. Each phase is designed to build upon the previous one, creating a cohesive and comprehensive approach to solving the problem of ineffective information retrieval in cloud library systems. By breaking down the research into these distinct phases, the study can systematically address each component, ensuring a thorough examination and implementation of solutions.

Document processing and text extraction form the foundational phase, where PDF documents are converted into clean, structured text. This is followed by text embedding and vectorization, where the cleaned text is transformed into numerical representations that capture semantic meaning. Query processing and semantic search utilize these embeddings to interpret user queries and retrieve relevant documents. Finally, user interface development focuses on creating an intuitive platform for users to interact with the system. Each phase employs specific methodologies and tools, which are detailed in the subsequent sections, to ensure that the research objectives are met effectively and efficiently.

3.2 Data Collection Methods

Document processing and text extraction are critical initial steps in the development of the linguistic support model. Text files from the selected cloud library systems are extracted using PyMuPDF, a robust tool that ensures high accuracy and retention of the original document structure. The extracted text is then cleaned to remove unwanted characters, line breaks, and other artifacts that might interfere with further processing. This preprocessing step is crucial for ensuring the quality and consistency of the data used in subsequent analyses [15].

The use of PyMuPDF was informed by its proven effectiveness in handling diverse document formats and its ability to accurately capture textual content. By employing a systematic approach to text extraction, the study ensures that the data is of high quality and suitable for embedding generation. This step lays a solid foundation for the development of accurate and reliable linguistic support models, as clean and well-structured text data is essential for effective information retrieval.

3.2.1 Initial Approach with GPT-2

Initially, the research aimed to utilize the GPT-2 model for generating text embeddings, considering its advanced capabilities in natural language processing. GPT-2, developed by OpenAI, is known for its powerful language understanding and generation capabilities. However, the implementation faced significant challenges due to the requirement for an API key, which was not freely available. This limitation hindered the ability to effectively implement and test the GPT-2 model within the scope of this project. The lack of a free API key for GPT-2 not only limited (Figure 3.1) accessibility but also posed financial constraints, making it impractical to rely on this model for the project's needs.

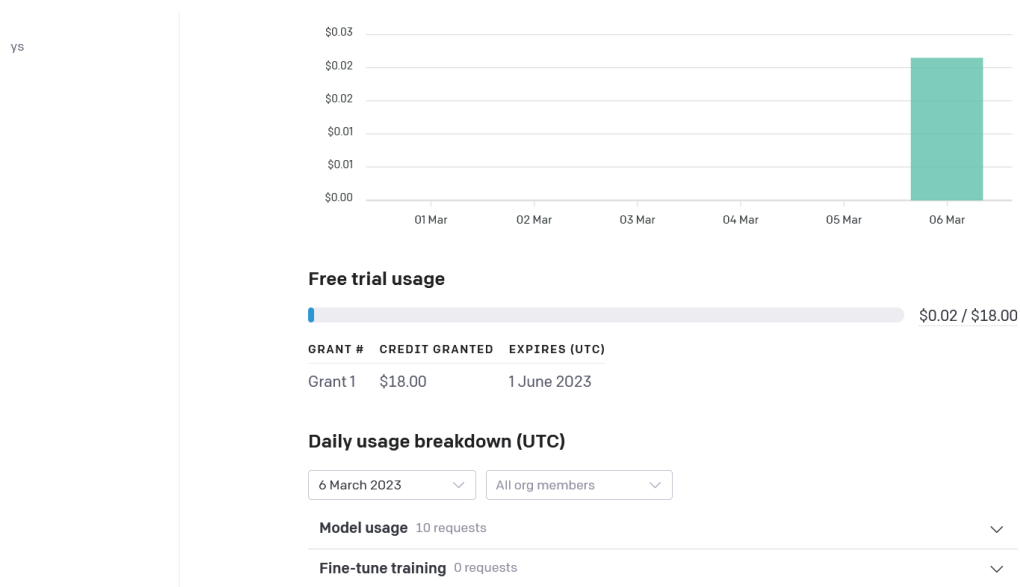


Figure 3.1: Limits of using API keys GPT

The necessity of a paid API key for GPT-2 highlighted a critical challenge in utilizing cutting-edge technologies that are not readily accessible. This limitation forced exploring alternative models and solutions that could be implemented within the project's budget and technical constraints. Despite its potential, the use of GPT-2 was deemed infeasible, leading to the exploration of other models that could provide similar benefits without the associated costs. This decision underscored the importance of considering both technical capabilities and practical constraints in the selection of tools and methodologies.

3.2.2 Alternative Approaches for Text Extraction

Given the challenges with GPT-2, the research turned to alternative methods for text extraction from PDF documents. Several libraries were considered, including PyPDF2 and PDFMiner, which are widely used for PDF text extraction. However, these libraries presented issues with accuracy, often resulting in incomplete or poorly formatted text. The extracted text from these tools frequently contained errors such as missing characters, incorrect line breaks, and misplaced words, which significantly affected the quality of the data (Figure 3.2). This lack of accuracy in text extraction posed significant challenges, as clean and structured text data is crucial for the subsequent embedding and retrieval processes [12].

To address these challenges, PyMuPDF (fitz) was selected for text extraction due to its robust capabilities in handling various PDF structures. PyMuPDF provided more accurate and reliable text extraction, ensuring that the textual data was clean and suitable for further processing. This tool allowed for the efficient extraction of text while preserving the layout and structure, which is essential for maintaining the integrity of the document's content. The decision to use PyMuPDF was based on its superior performance in accurately extracting text and its ability to handle complex PDF formats, making it a reliable choice for this critical phase of the research.



Figure 3.2: From PDF to TXT format

3.3 Text Embedding and Vectorization

The cleaned text is split into chunks using the `CharacterTextSplitter` tool, with a chunk size of 1000 characters and an overlap of 10 characters. This method ensures that each chunk retains sufficient context for accurate processing. Each chunk is then transformed into embeddings using the Mistral model via the `OllamaEmbeddings` class. These embeddings capture the semantic content of the text, providing dense vector representations essential for semantic search. By converting textual data into embeddings, the study leverages advanced natural language processing techniques to enhance the retrieval accuracy of the information system [15].

Embedding generation is a critical step in the methodology, as it enables the system to understand and interpret the semantic relationships between different pieces of text. The use of the Mistral model ensures that the embeddings are of high quality and capable of capturing subtle nuances in the text. This step significantly improves the system's ability to retrieve relevant documents based on user queries, moving beyond simple keyword matching to a more sophisticated understanding of the information being sought.

3.3.1 Embedding Generation

LangChain’s Mistral model was used to generate embeddings for the cleaned text data. The model transforms textual information into high-dimensional vectors that capture the semantic meaning of the content. These embeddings represent the text in a way that highlights the context and thematic relevance, making it possible to perform more accurate and contextually aware searches (Figure 3.3).

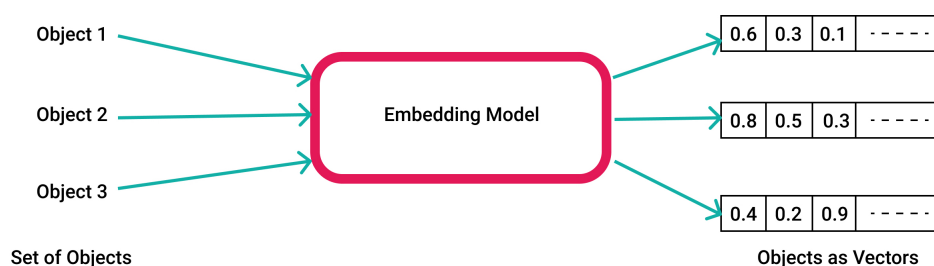


Figure 3.3: Process of Embedding

By utilizing LangChain’s Mistral model, the research ensured that the embeddings captured the nuanced semantic relationships within the documents. These high-quality embeddings allow the system to interpret and respond to user queries more effectively, providing search results that are not only relevant but also contextually accurate. This capability significantly enhances the performance of the thematic search model, enabling it to meet the complex informational needs of users.

3.3.2 Vector Storage

Chroma, a vector database, was employed to store the generated embeddings. Chroma provides high-performance storage and retrieval capabilities, ensuring that the system can quickly and efficiently access the embeddings when needed. (Figure 3.4).

Storing the embeddings in Chroma ensures that the system maintains optimal performance and scalability. Chroma’s efficient retrieval mechanisms allow the system to

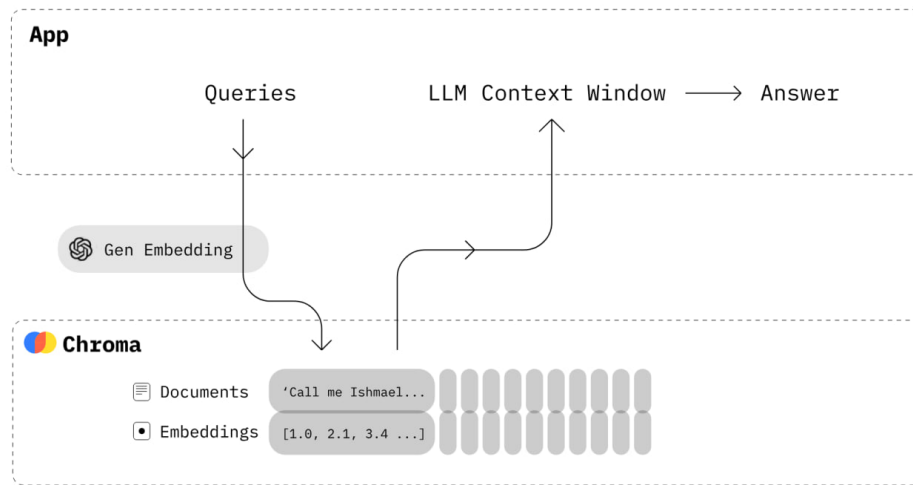


Figure 3.4: Process of vectorizing data

handle large datasets, providing rapid access to embeddings during search operations. This storage solution is critical for the overall functionality of the thematic search model, as it supports the efficient processing and retrieval of relevant documents based on semantic similarity. By leveraging Chroma’s capabilities, the research ensured that the thematic search model could scale effectively, accommodating increasing volumes of data while maintaining high levels of performance and accuracy.

3.4 Query Processing and Semantic Search

Query processing and semantic search are crucial components of the thematic search model, enabling the system to interpret user queries accurately and retrieve relevant documents based on thematic content. For query processing, the system generates embeddings for user queries using the same model employed for document embeddings. This approach ensures consistency and accuracy in the retrieval process by maintaining a uniform representation of both queries and documents. Utilizing embeddings allows the system to understand the semantic meaning behind the user queries, which is essential for delivering precise and relevant search results. A semantic search algorithm then uses

these query embeddings to perform searches in the Chroma vector database, focusing on thematic content rather than mere keyword matches.

This method addresses the limitations of traditional keyword-based searches, which often fail to capture the full context and meaning of user queries. By leveraging semantic search, the system can provide more accurate and contextually relevant results, significantly enhancing the user experience. The combination of query embeddings and a sophisticated search algorithm ensures that the system can interpret complex queries and retrieve documents that are closely aligned with the user's informational needs. This phase is integral to the effectiveness of the thematic search model, ensuring that users can quickly and easily find the information they seek.

3.4.1 Query Embedding

User queries are processed to generate embeddings using LangChain's Mistral model. This step involves transforming the natural language query into a vector representation that captures the semantic meaning of the user's request (Figure 3.5). By using the same embedding model for both documents and queries, the system ensures that the search process is consistent and accurate.

By converting the user queries into embeddings, the system can leverage the rich semantic information contained in the embeddings to perform more sophisticated and accurate searches. This approach allows the system to understand the intent and context of the user queries, providing a deeper understanding of the user's needs. The use of embeddings for query processing ensures that the search algorithm can handle complex and nuanced queries, delivering results that are not only relevant but also contextually appropriate.

The consistency in using the same model for both documents and queries ensures that the embeddings are comparable, facilitating accurate matching during the search process. This method significantly improves the relevance of the search results, addressing the common issue of irrelevant or incomplete results in traditional keyword-based searches.

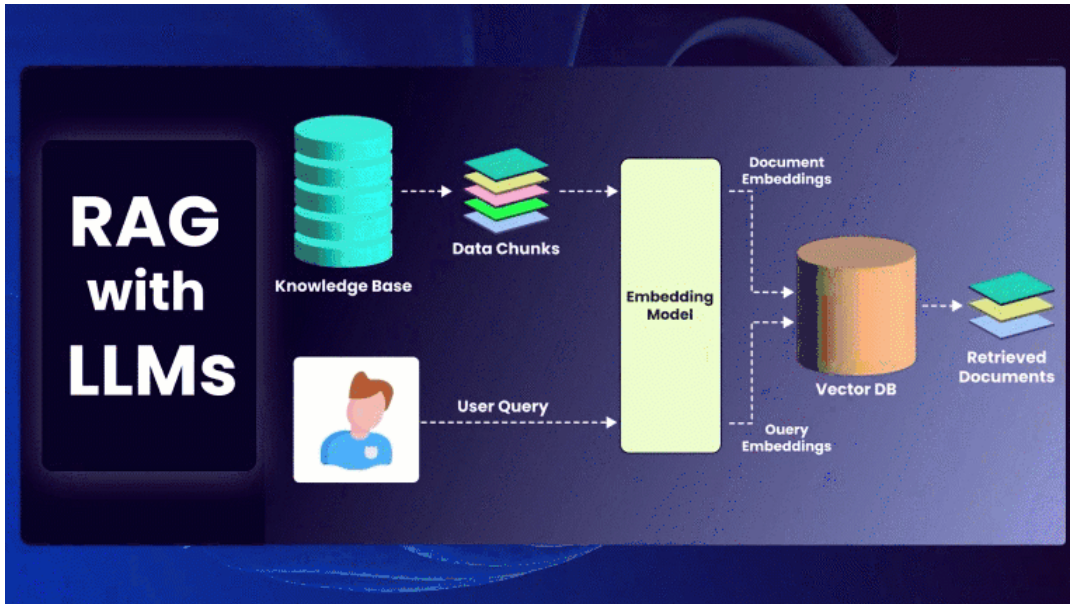


Figure 3.5: RAG Embedding

By focusing on semantic content, the system can provide users with more meaningful and useful information, enhancing the overall search experience.

3.4.2 Semantic Search Algorithm

The semantic search algorithm compares the query embedding with the document embeddings stored in the Chroma vector database. This comparison is based on the similarity between the vectors, allowing the system to identify documents that are thematically related to the query.

This approach significantly enhances the relevance and accuracy of the search results by focusing on the thematic content of the documents rather than just matching keywords. By comparing the embeddings, the system can identify documents that are semantically similar to the query, even if the exact keywords are not present. This capability addresses one of the primary limitations of traditional search methods, which often miss relevant documents due to variations in terminology or phrasing.

The semantic search algorithm leverages the rich semantic information captured in

the embeddings, enabling the system to perform more contextually aware searches. This results in a higher quality of search results, providing users with documents that are more closely aligned with their informational needs. The ability to perform accurate and relevant searches based on thematic content is a key strength of the thematic search model, ensuring that users receive the most pertinent information in response to their queries. This phase is essential for delivering a superior search experience, enhancing the overall effectiveness and utility of the cloud library system.

3.5 User Interface Development

The user interface development phase was crucial for ensuring that users could easily interact with the thematic search model. The user interface (UI) was developed using Flask, a lightweight web framework known for its simplicity and flexibility. Flask's minimalistic design allows for the quick creation of web applications while providing the necessary tools to build a robust and user-friendly interface. By using Flask, the development team could focus on creating an intuitive platform where users could input queries and view results efficiently. This UI design aims to facilitate seamless communication between the user and the retrieval system, enhancing the overall user experience [14].

The decision to use Flask was based on its ability to handle both the front-end and back-end integration smoothly. Flask's extensive library of plugins and extensions allowed for rapid development and deployment of the user interface. This choice ensured that the interface was not only functional but also scalable and maintainable. Flask's flexibility enabled the development team to implement custom features and functionalities tailored to the specific needs of the project, ensuring that the user interface could support the advanced capabilities of the thematic search model.

3.5.1 Web Interface Design

The web interface is designed to be intuitive and easy to use, enabling users to input their queries and receive search results with minimal effort. Key features of the interface

include a search bar for query input, a results display area for showing the search results, and pagination controls to manage the display of large sets of results. These features are essential for enhancing the usability of the interface, ensuring that users can quickly and easily access the information they need.

The design of the web interface focuses on user-friendliness and efficiency (Figure 3.6). The search bar is prominently placed to encourage user engagement, while the results display area is structured to present information clearly and concisely. Pagination controls help users navigate through multiple pages of results without overwhelming them with too much information at once. The overall layout is clean and straightforward, minimizing distractions and allowing users to focus on their search tasks. This design approach ensures that users of all skill levels can effectively utilize the thematic search model, making it a valuable tool for information retrieval in cloud library systems.

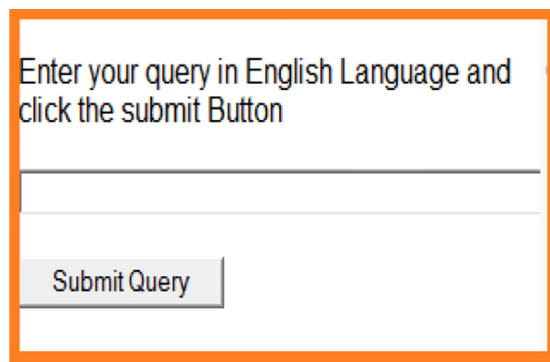


Figure 3.6: Simple User Interface

3.5.2 Backend Integration

Flask is used to handle the backend logic, including processing user queries, generating embeddings, performing semantic searches, and retrieving results from the Chroma vector database. This integration is crucial for ensuring that the system operates smoothly and efficiently, providing users with timely and relevant search results.

This setup illustrates how Flask manages both the frontend user interactions and the backend processing seamlessly. Flask handles the reception of user queries, processes

these queries by generating embeddings, performs semantic searches in the Chroma vector database, and retrieves the relevant results to be displayed in the user interface. This seamless integration ensures that users receive accurate and timely search results, enhancing the overall efficiency and effectiveness of the thematic search model. By leveraging Flask's capabilities, the development team was able to create a responsive and user-friendly interface that meets the needs of users and supports the advanced functionalities of the system.

3.6 Tasks and System Interactions

The system developed for this project is designed to perform several key tasks, each contributing to the overall functionality and effectiveness of the thematic search model. These tasks are organized to ensure a streamlined process from document ingestion to result retrieval and display. By breaking down the system into these distinct tasks, we can ensure that each component operates efficiently and interacts seamlessly with others to deliver accurate and relevant search results.

- **Document Ingestion** The first task involves collecting and storing PDF documents in a designated repository. This repository serves as the primary data source for the system, housing all documents that will be processed. Ensuring that the repository is well-organized and accessible is crucial for the efficiency of subsequent tasks. This initial step lays the foundation for the entire workflow, making sure that all necessary documents are available for text extraction and processing.
- **Text Extraction** Text extraction is performed using PyMuPDF to convert the content of the PDF documents into clean, structured text. This process involves reading the PDFs, extracting the text, and then cleaning it to remove any unwanted characters or formatting issues. PyMuPDF was chosen for its robustness and accuracy, ensuring that the extracted text retains the necessary context and structure. This clean text is then ready for embedding generation, forming the basis

for semantic analysis.

- **Embedding Generation** Once the text has been extracted and cleaned, the next task is to generate text embeddings using LangChain's Mistral model. These embeddings are high-dimensional vectors that capture the semantic meaning of the text, making it possible for the system to understand the context and thematic relevance of the documents. This step is crucial for enabling advanced search capabilities, as it transforms raw text data into a format suitable for semantic analysis and retrieval.
- **Embedding Storage** The generated embeddings are stored in a Chroma vector database. Chroma is selected for its high-performance storage and retrieval capabilities, ensuring that the embeddings can be accessed quickly and efficiently during the search process. Storing embeddings in a vector database is critical for maintaining the performance and scalability of the system, allowing it to handle large volumes of data without compromising on speed or accuracy.
- **Query Processing** For query processing, the system generates embeddings for user queries using the same model employed for document embeddings. This ensures consistency and accuracy in the retrieval process, as both documents and queries are represented in the same semantic space. The query embeddings are then used to perform semantic searches in the Chroma vector database, identifying documents that are thematically related to the query.
- **Result Retrieval and Display** The final task involves retrieving relevant documents based on query embeddings and displaying the results to the user through the Flask web interface. The results are presented in a user-friendly format, allowing users to easily access the information they need. Flask handles the backend logic, ensuring smooth and efficient processing of user queries and retrieval of search results.

- **System Interactions** The system involves interactions between various components, including the user interface, backend processing modules, and the vector database. These interactions are designed to ensure a seamless flow of data and efficient retrieval of information. The user interface, developed using Flask, allows users to input queries and view results. The backend processing modules handle the extraction, embedding generation, and query processing tasks, while the Chroma vector database manages the storage and retrieval of embeddings.

By integrating these components, the system can provide a robust and efficient solution for thematic search in cloud library systems. Each task and interaction is optimized to ensure that the system operates smoothly and delivers accurate, relevant search results to users. This comprehensive approach ensures that all aspects of the thematic search model work together harmoniously, enhancing the overall functionality and user experience.

3.7 Challenges Faced

Several challenges were encountered during the project, each requiring specific solutions to ensure the successful implementation of the thematic search model. Addressing these challenges was crucial to maintaining the integrity, efficiency, and effectiveness of the system.

- **API Key Issues with GPT-2:** One of the initial challenges faced in this project was the requirement for a paid API key to use the GPT-2 model, which restricted its accessibility and imposed financial constraints. This issue underscored the importance of finding accessible and cost-effective solutions in academic research. As a result, alternative models, such as LangChain's Mistral model, were explored and ultimately adopted. This change not only addressed the accessibility issue but also provided a more feasible approach for embedding generation within the project's budgetary limitations.

- **Accuracy of Text Extraction:** Accurate text extraction proved to be a significant challenge in the early stages of the project. Initial attempts with libraries like PyPDF2 and PDFMiner yielded unreliable results, with extracted text often containing errors such as missing characters and incorrect formatting. These inaccuracies were detrimental to the quality of data used in subsequent embedding and retrieval processes. The adoption of PyMuPDF resolved these issues, providing a more reliable and precise method for text extraction. Ensuring high accuracy in text extraction was essential for maintaining the integrity and usability of the data.
- **Handling Large-Scale Data:** Managing and processing large volumes of text data posed significant computational challenges, particularly in optimizing the text embedding and retrieval processes. Efficiently handling such large datasets required careful resource management and the design of efficient algorithms. Scalable solutions such as LangChain for embedding generation and Chroma for vector storage were employed to address these challenges. These tools are specifically designed to manage large datasets effectively, ensuring that the system could handle extensive data volumes while maintaining high performance and accuracy.
- **Integrating Semantic Search:** Implementing a robust semantic search algorithm was another major challenge. The goal was to accurately interpret user queries and retrieve relevant documents based on thematic content rather than simple keyword matches. This required extensive testing and fine-tuning to ensure that the search algorithm performed well and met user expectations. The complexity of semantic search algorithms necessitates continuous refinement to maintain their accuracy and relevance. Through iterative testing and improvements, the semantic search component was developed to deliver high-quality, contextually relevant search results, enhancing the overall effectiveness of the thematic search model.

Chapter 4

Development

This chapter details the implementation of the proposed solution for enhancing thematic search capabilities in cloud library systems. The goal was to create an efficient and accurate retrieval system capable of handling the vast and unstructured data typically found in digital libraries. The development process involved multiple stages, each critical to the overall functionality of the system. These stages included text extraction, embedding generation, vector storage, query processing, and user interface development. Each stage was meticulously designed to address specific challenges and ensure seamless integration with the entire system.

The implementation began with the extraction of text from PDF documents, a crucial step in preparing the data for embedding generation. PyMuPDF was selected for this task due to its robust capabilities in handling various PDF structures and ensuring high accuracy in text extraction. Following this, the text was cleaned and preprocessed to ensure it was suitable for embedding generation. LangChain's Mistral model was then employed to generate embeddings that captured the semantic content of the documents. These embeddings were stored in Chroma, a high-performance vector database, facilitating efficient retrieval based on semantic similarity.

Query processing and semantic search were pivotal components of the system, ensuring that user queries were accurately interpreted and relevant documents were retrieved. This was achieved by generating embeddings for user queries and comparing them with

document embeddings stored in Chroma. The final component of the implementation was the development of a user-friendly web interface using Flask. This interface allowed users to input queries and view results seamlessly, ensuring an intuitive and efficient user experience. Throughout the development process, several challenges were encountered, such as handling large-scale data and ensuring the accuracy of text extraction, which were addressed through the use of open-source libraries and iterative testing and optimization.

4.1 Text Extraction and Cleaning

The first step in the implementation was to extract and clean text from PDF documents. This step was crucial for preparing the data for embedding generation and ensuring the accuracy and relevance of the search results. Given the initial challenges faced with other libraries, such as PyPDF2 and PDFMiner, which often resulted in incomplete or poorly formatted text, PyMuPDF (fitz) was selected for its robustness and accuracy. PyMuPDF demonstrated superior capability in handling various PDF structures, providing a reliable solution for extracting text while preserving the document's layout and integrity. This ensured that the extracted text retained the necessary context and structure, which is critical for subsequent processing stages.

The text extraction process involved reading the PDF documents and extracting the text content while maintaining the structure of the document. This was achieved using PyMuPDF's advanced functionalities, which allowed for precise and efficient text extraction. By ensuring minimal loss of information and maintaining the integrity of the document content, this method facilitated the generation of high-quality embeddings. The extracted text was then subjected to a cleaning process to remove unwanted characters and line breaks, which could affect the quality and accuracy of the embeddings. This step was crucial for preparing clean and structured text data, essential for effective embedding generation and accurate thematic search results.

4.1.1 Implementation with PyMuPDF

PyMuPDF was used to read PDF documents and extract text while preserving the document structure. The following code (Listing 4.1) snippet demonstrates the text extraction process:

```
1 import fitz # PyMuPDF
2
3 def extract_text_from_pdf(pdf_path):
4     doc = fitz.open(pdf_path)
5     text = ""
6     for page in doc:
7         text += page.get_text()
8     return text
```

Listing 4.1: Code to extract text from a PDF

This method efficiently extracted text, ensuring minimal loss of information and maintaining the integrity of the document content. The choice of PyMuPDF was driven by its ability to handle various PDF structures and accurately extract text without compromising the document's original format. This precision was crucial for ensuring that the text data used in subsequent stages was of high quality and representative of the original documents. The extracted text maintained the document's original layout and context, which is essential for accurately capturing the semantic content during the embedding generation phase. By preserving the document structure, PyMuPDF facilitated the creation of embeddings that accurately represented the thematic elements of the text, ensuring the relevance and accuracy of the search results. This method proved to be both efficient and reliable, providing a robust foundation for the subsequent text cleaning and embedding generation processes.

4.1.2 Text Cleaning

Post-extraction, the text required cleaning to remove unwanted characters and line breaks. The following code (Listing 4.2) function was employed to clean the text:

```
1 def clean_text(text):  
2     cleaned_text = text.replace("-\n", "").replace("\n", " ")  
3     return cleaned_text
```

Listing 4.2: Function to clean extracted text

This step was crucial for ensuring the quality of the data used in the embedding generation phase. Cleaning the text involved removing hyphenation at line breaks, which is common in PDF documents, and consolidating lines to create a continuous flow of text. This process helped eliminate artifacts that could negatively impact the quality of the embeddings and the accuracy of the search results. By cleaning the text, the data became more uniform and structured, making it suitable for embedding generation. Clean text data is essential for producing high-quality embeddings that accurately capture the semantic content of the documents. This ensures that the thematic search model can effectively interpret and retrieve relevant documents based on user queries. The text cleaning process, therefore, played a critical role in maintaining the integrity and accuracy of the data used in the system, ultimately contributing to the effectiveness of the thematic search capabilities in cloud library systems.

4.2 Embedding Generation and Vector Storage

Generating embeddings and storing them in a vector database were critical steps in enabling efficient and accurate semantic search capabilities. Embeddings are essential as they transform text data into a numerical format that encapsulates the semantic meaning of the content, allowing for sophisticated search and retrieval operations. LangChain's Mistral model was selected for generating these embeddings due to its robust framework and high performance in capturing semantic content. The embeddings were stored in

Chroma, a vector database designed for high-performance retrieval, ensuring that the system could handle large-scale data efficiently.

The process began with generating embeddings for the cleaned text data. LangChain's Mistral model was employed to convert the text into high-dimensional vectors that represented the semantic content accurately. These embeddings were then stored in Chroma, which facilitated efficient retrieval based on semantic similarity. The integration of LangChain and Chroma was crucial for maintaining the performance and scalability of the system, allowing it to handle large volumes of data and perform complex searches quickly. This setup ensured that the thematic search model could deliver accurate and relevant results to users, enhancing the overall user experience.

4.2.1 Embedding Generation with LangChain

LangChain's Mistral model was used to generate embeddings for the cleaned text data. The model transforms textual information into high-dimensional vectors that capture the semantic meaning of the content. These embeddings allow the system to understand the context and thematic relevance of the documents, making it possible to retrieve more accurate and relevant results during searches. The following code (Listing 4.3) snippet shows how embeddings were generated:

```
1 from langchain_community.embeddings import OllamaEmbeddings
2
3 def generate_embeddings(text, model="mistral"):
4     embedding_model = OllamaEmbeddings(model=model)
5     embeddings = embedding_model.embed_text(text)
6     return embeddings
```

Listing 4.3: Generating embeddings for text

This approach ensured that the embeddings accurately represented the thematic content of the documents. By capturing the semantic meaning, these embeddings enabled the system to perform more nuanced and contextually aware searches, addressing the

limitations of traditional keyword-based search methods.

The embeddings generated by LangChain’s Mistral model provided a robust foundation for the thematic search capabilities of the system. These embeddings encapsulated the rich semantic information in the text, allowing the search algorithm to understand and match user queries with relevant documents effectively. This method significantly improved the relevance and accuracy of search results, providing users with a better search experience.

4.2.2 Storing Embeddings in Chroma

Chroma was selected for its high-performance storage and retrieval capabilities. The embeddings generated by LangChain were stored in Chroma, which allowed for efficient and quick retrieval based on semantic similarity. The following code (Listing 4.4) demonstrates how embeddings were stored in the Chroma vector database:

```
1 from langchain_community.vectorstores import Chroma
2
3 def store_embeddings(embeddings, persist_directory):
4     vector_store = Chroma.from_embeddings(embeddings,
5         persist_directory=persist_directory)
6     vector_store.persist()
```

Listing 4.4: Storing embeddings in Chroma vector database

This setup enabled efficient retrieval of documents based on semantic similarity, facilitating the thematic search functionality. Chroma’s performance in handling large-scale data was instrumental in ensuring that the system could manage extensive collections of documents and deliver results promptly.

Storing the embeddings in Chroma ensured that the search system could quickly and accurately retrieve relevant documents in response to user queries. Chroma’s optimized storage and retrieval mechanisms allowed the system to scale efficiently, handling increasing volumes of data without compromising performance. This integration was crucial for

maintaining the effectiveness and responsiveness of the thematic search model, providing users with timely and relevant search results.

4.3 Query Processing and Semantic Search

The next step involved processing user queries and performing semantic searches to retrieve relevant documents. This phase was crucial for ensuring that the system could interpret user queries accurately and return documents that were thematically relevant. Traditional keyword-based search methods often fall short in understanding the context and semantic meaning behind user queries, leading to irrelevant search results. By leveraging advanced NLP techniques, the system aimed to overcome these limitations and provide more precise and relevant results [16].

The query processing involved generating embeddings for user queries, using the same model that was employed for document embeddings. This ensured consistency in how both queries and documents were represented semantically, allowing for accurate matching during the search process. The semantic search algorithm then compared the query embeddings with the document embeddings stored in the Chroma vector database, retrieving the most relevant documents based on thematic content rather than mere keyword presence. This approach significantly enhanced the relevance and accuracy of the search results, providing users with a more effective and satisfying search experience.

4.3.1 Query Embedding

User queries were processed to generate embeddings using the same model employed for document embeddings. This step involved transforming the user input into a high-dimensional vector that captured the semantic meaning of the query. By using the same embedding model for both documents and queries, the system ensured a consistent and accurate representation of the semantic content, facilitating precise matching during the search process. The following code (Listing 4.5) snippet shows the query processing implementation:

```
1 def process_query(query, model="mistral"):  
2     embedding_model = OllamaEmbeddings(model=model)  
3     query_embedding = embedding_model.embed_text(query)  
4     return query_embedding
```

Listing 4.5: Processing user queries to generate embeddings

This method allowed the system to handle natural language queries effectively, interpreting the user's intent and generating an embedding that accurately represented the query's thematic content. By processing queries in this manner, the system could leverage the rich semantic information contained in the embeddings, enabling more accurate and contextually relevant searches.

The use of embeddings for query processing ensured that the search algorithm could understand and interpret complex queries, including those involving synonyms, related terms, and contextual nuances. This capability was a significant improvement over traditional keyword-based searches, which often struggle to capture the full meaning and context of user queries. By generating embeddings for queries, the system could perform more sophisticated and accurate searches, providing users with highly relevant results that met their informational needs.

4.3.2 Semantic Search Algorithm

The semantic search algorithm compared the query embedding with the document embeddings stored in the Chroma vector database. This process involved calculating the similarity between the query embedding and each document embedding, identifying the documents that were thematically most similar to the query. The following code (Listing 4.6) demonstrates the semantic search process:

```
1 def semantic_search(query_embedding, persist_directory):  
2     vector_store = Chroma.load(persist_directory=persist_directory)  
3     results = vector_store.similarity_search(query_embedding)
```

4

```
return results
```

Listing 4.6: Performing semantic search

This algorithm enhanced the relevance and accuracy of the search results by focusing on the thematic content of the documents rather than just keyword matches. By comparing the embeddings, the system could identify documents that were semantically related to the query, even if they did not contain the exact keywords used in the query. This approach addressed the limitations of traditional keyword-based searches, which often return irrelevant or incomplete results due to their inability to understand the context and meaning behind the user's query.

The semantic search algorithm provided a robust and efficient way to retrieve relevant documents, leveraging the rich semantic information captured in the embeddings. By focusing on the thematic similarity between queries and documents, the algorithm could deliver highly relevant search results, enhancing the overall effectiveness of the information retrieval system. This capability was crucial for meeting the diverse and complex informational needs of users, providing them with accurate and contextually appropriate documents based on their queries.

4.4 User Interface Development

The user interface was developed using Flask, a lightweight web framework, to provide a user-friendly platform for inputting queries and viewing results. Flask was chosen for its simplicity and flexibility, making it an ideal choice for creating a web interface that could seamlessly interact with the backend processes of the information retrieval system. The goal was to design an interface that was intuitive and easy to use, enabling users to interact with the system without requiring technical expertise. This involved creating a clean and straightforward layout where users could enter their queries and view the search results efficiently.

The development of the user interface focused on enhancing the user experience by

providing clear and immediate feedback for their interactions. The interface needed to handle user inputs, process these inputs through the backend system, and display the results in a meaningful way. Flask's robust features allowed for the integration of these functionalities, ensuring that the user interface was both responsive and functional. This was critical in providing a smooth user experience, where users could quickly and easily access the information they needed through an efficient and effective search process.

4.4.1 Web Interface Design

The web interface was designed to be intuitive and easy to use. It included features such as a search bar, result display area, and pagination controls to enhance usability. The design aimed to minimize the learning curve for users, allowing them to perform searches and view results with minimal effort. The following code (Listing 4.7) snippet outlines the basic setup of the Flask application:

```
1 from flask import Flask, request, render_template_string
2
3 app = Flask(__name__)
4
5 html_template = """
6 <!doctype html>
7 <html lang="en">
8   <head>
9     <meta charset="UTF-8">
10    <meta name="viewport" content="width=device-width,
11      initial-scale=1.0">
12    <link
13      href="https://stackpath.bootstrapcdn.com/bootstrap/4.5.2/css/bootstrap.mi
14      rel="stylesheet">
15    <title>Query Interface</title>
16  </head>
```

```

14 <body>
15   <div class="container">
16     <h1 class="text-center">Query Interface</h1>
17     <form method="POST" action="/">
18       <div class="form-group">
19         <input type="text" name="query" class="form-control"
20           placeholder="Enter your query">
21       </div>
22       <div class="text-center">
23         <input type="submit" class="btn btn-primary"
24           value="Submit">
25       </div>
26     </form>
27     {% if results %}
28     <div class="result">
29       <h2>Results:</h2>
30       {% for result in results %}
31       <p>{{ result }}</p>
32       {% endfor %}
33     </div>
34     {% endif %}
35   </div>
36 </body>
37 </html>
38 """
39 @app.route("/", methods=["GET", "POST"])
40 def index():
41     if request.method == "POST":
42         user_query = request.form.get("query")

```

```

42     # Process the query and get results
43     results = process_query_and_search(user_query)
44     return render_template_string(html_template, results=results)
45 return render_template_string(html_template)
46
47 if __name__ == "__main__":
48     app.run(debug=True)

```

Listing 4.7: Flask web application setup

This code sets up a basic Flask application with a simple HTML template. The interface includes a form where users can enter their queries and a section to display the results. The use of Bootstrap ensures that the interface is visually appealing and responsive, providing a consistent user experience across different devices. This design approach ensures that users can interact with the system efficiently, enhancing their overall experience.

The HTML template includes essential elements such as a search bar and a results display area, which are crucial for the functionality of the search interface. The design is minimalistic, focusing on functionality and ease of use. This ensures that users can quickly understand how to use the interface and perform searches without any confusion. By using Flask and Bootstrap, the development team was able to create a web interface that was both functional and user-friendly, meeting the needs of the users effectively.

4.4.2 Backend Integration

Flask handled the backend logic, including processing user queries, generating embeddings, performing semantic searches, and retrieving results from the Chroma vector database. This integration ensured smooth operation and efficient handling of user requests. Flask's flexibility allowed for seamless communication between the frontend and backend, ensuring that user inputs were processed correctly and results were displayed

promptly. This integration was critical for maintaining the performance and responsiveness of the system, providing users with an efficient search experience.

The backend logic involved several steps: receiving the user query, processing the query to generate embeddings, performing a semantic search in the Chroma vector database, and retrieving and displaying the relevant results. Flask facilitated these processes by providing a robust framework for handling HTTP requests, managing the server-side logic, and ensuring secure and efficient data handling. The following code (Listing 4.8) snippet demonstrates how Flask handled these backend processes:

```
1 @app.route("/", methods=["GET", "POST"])
2 def index():
3     if request.method == "POST":
4         user_query = request.form.get("query")
5         # Process the query and get results
6         results = process_query_and_search(user_query)
7         return render_template_string(html_template, results=results)
8     return render_template_string(html_template)
9
10 if __name__ == "__main__":
11     app.run(debug=True)
```

Listing 4.8: Flask web application routing and execution

This integration ensured that user queries were processed efficiently and results were retrieved and displayed quickly. By leveraging Flask's capabilities, the development team was able to create a seamless interaction between the user interface and the backend processes, enhancing the overall functionality and user experience of the system. This approach ensured that the system could handle user requests efficiently, providing timely and relevant search results, and maintaining high performance and reliability.

4.5 Challenges and Technical Solutions

Throughout the development process, several challenges were encountered, and technical solutions were implemented to overcome these challenges. Addressing these challenges was critical to ensuring the success of the thematic search model and enhancing the system's overall performance and reliability.

API Key Issues with GPT-2 The initial plan to use the GPT-2 model for generating text embeddings faced a significant obstacle due to the requirement for a paid API key. This restricted the use of GPT-2 and posed financial constraints that were not feasible within the project's budget. To address this issue, alternative models were explored, leading to the selection of LangChain's Mistral model. The Mistral model provided a robust and cost-effective solution, offering similar capabilities without the need for a paid API key. This change not only resolved the accessibility issue but also ensured that the project remained within budget, allowing the development team to focus resources on other critical areas.

The implementation of LangChain's Mistral model proved to be highly effective, providing high-quality embeddings that accurately represented the semantic content of the documents. By leveraging this model, the project was able to maintain the desired functionality and performance levels without incurring additional costs. This solution demonstrated the importance of flexibility and adaptability in the development process, highlighting how alternative approaches can effectively address unforeseen challenges and constraints.

Accuracy of Text Extraction Accurate text extraction is crucial for maintaining the integrity of the data used in the embedding and retrieval processes. Initial attempts with libraries like PyPDF2 and PDFMiner resulted in inaccurate text extraction, with issues such as missing characters, incorrect line breaks, and misplaced words. These inaccuracies posed significant challenges, as they could negatively impact the quality of the embeddings and the accuracy of the search results. To overcome this challenge, PyMuPDF (fitz) was adopted for its reliability and precision in handling various PDF structures.

PyMuPDF's advanced functionalities allowed for precise and efficient text extraction, ensuring minimal loss of information and preserving the document's original format. This high level of accuracy was essential for producing clean and structured text data, which is critical for generating high-quality embeddings. By using PyMuPDF, the development team was able to ensure that the text data used in the system was reliable and accurate, significantly improving the overall performance and effectiveness of the thematic search model.

Handling Large-Scale Data Managing and processing large volumes of text data posed significant computational challenges. The system needed to handle extensive datasets efficiently to ensure timely and accurate search results. Optimizing the text embedding and retrieval processes required careful resource management and efficient algorithm design. To address this challenge, scalable solutions such as LangChain and Chroma were employed, both of which are designed to handle large datasets effectively.

Techniques such as batching and parallel processing were implemented to improve performance and reduce processing time. Batching allowed the system to process multiple documents simultaneously, while parallel processing distributed the computational load across multiple processors. These techniques ensured that the system could manage and process large volumes of data efficiently, maintaining high performance and responsiveness. By leveraging these scalable solutions, the development team was able to address the challenges of handling large-scale data, ensuring that the system remained efficient and effective as the dataset grew. **Integrating Semantic Search**

Implementing a robust semantic search algorithm was essential for providing accurate and relevant search results. However, this task required extensive testing and fine-tuning to ensure high performance and user satisfaction. The complexity of semantic search algorithms necessitates ongoing refinement to maintain accuracy and relevance. Various metrics, including precision, recall, and F1 score, were used to evaluate the performance of the search algorithm.

Continuous feedback and iterative improvements helped in refining the algorithm. This process involved rigorous testing to identify any weaknesses or areas for improvement,

followed by adjustments to the algorithm to enhance its accuracy and performance. The development team focused on optimizing the algorithm to ensure that it could accurately interpret user queries and retrieve relevant documents based on thematic content. By continuously refining the semantic search algorithm, the team was able to enhance the overall effectiveness of the search system, providing users with a reliable and efficient tool for information retrieval.

Security and Privacy Considerations Given the sensitivity of data in library systems, security and privacy were important considerations in the development of the thematic search model. Measures were taken to ensure that user queries and document data were handled securely.

- **Data Encryption:** All data stored in the Chroma vector database was encrypted to prevent unauthorized access. Encryption ensures that even if the data is accessed by unauthorized parties, it remains unreadable without the proper decryption keys.
- **Access Controls:** Access to the system was restricted through authentication and authorization mechanisms. Only authorized users were allowed to access the system and perform searches. Role-based access control (RBAC) was implemented to ensure that users had access only to the data and functionalities necessary for their roles.
- **Secure Communication:** Communication between the web interface and the backend was secured using HTTPS. This ensures that data transmitted over the network is encrypted and protected from interception and tampering.

4.6 Performance Optimization

To ensure that the system performs efficiently, various performance optimization techniques were employed. These techniques focus on improving retrieval times, reducing system load, distributing computational tasks, and ensuring scalability. Implementing

these optimizations is crucial for maintaining a responsive and reliable information retrieval system.

- **Indexing:** Efficient indexing of embeddings in the Chroma vector database was crucial for speeding up retrieval times. Indexing structures such as inverted indexes and tree-based indexes were used to quickly locate relevant embeddings based on user queries, significantly reducing the search time. By organizing the data for rapid access, the system can handle large datasets more effectively, ensuring that even complex queries are processed swiftly.
- **Caching:** Frequently accessed data and results were cached to minimize the need for recomputation and reduce system load. Caching commonly requested information in memory allows the system to retrieve it quickly, improving response times and reducing latency. This strategy is particularly beneficial for handling repeated queries, as it eliminates redundant processing and enhances overall performance.
- **Load Balancing:** Load balancing techniques were employed to distribute the computational load evenly across multiple servers. This prevents any single server from becoming a bottleneck, thereby enhancing the system's overall performance and reliability. Load balancing algorithms, such as round-robin and least connections, were used to manage task distribution. By ensuring an even distribution of load, the system can handle more queries simultaneously without performance degradation.
- **Scalability:** The system was designed with scalability in mind, allowing it to handle increasing volumes of data and user queries efficiently. Horizontal scaling techniques, such as adding more servers, were employed to accommodate growing demands. The architecture supports easy integration of additional resources, ensuring that the system remains robust and performs well under high usage scenarios. Scalability ensures that the system can maintain high performance as the number of users and the size of the dataset grow.

- **Optimization of Text Embeddings:** The generation of text embeddings was optimized by adjusting the chunk size and overlap parameters. By setting an appropriate chunk size of 1000 characters with an overlap of 10 characters, the system ensures that each chunk retains enough context for accurate embedding generation. This optimization balances context retention with computational efficiency, ensuring high-quality embeddings are generated efficiently.

These performance optimization strategies collectively ensure that the system remains efficient, reliable, and capable of handling large-scale information retrieval tasks. By focusing on indexing, caching, load balancing, and scalability, the system can provide quick and accurate responses to user queries, maintaining high performance under varying loads and usage scenarios.

4.7 Summary

The development of the thematic search model for cloud library systems involved multiple meticulously executed steps. Each phase, from text extraction to user interface development, was crucial to creating a robust and efficient system capable of handling complex and large-scale data.

The text extraction and cleaning phase, implemented using PyMuPDF, provided high-quality textual data essential for generating accurate embeddings. LangChain's Mistral model was then employed for embedding generation, with Chroma used for high-performance storage and retrieval. This setup ensured that the system could manage and process large volumes of data efficiently.

Query processing and the implementation of a semantic search algorithm ensured that user queries were interpreted accurately and relevant documents were retrieved based on thematic content. The user interface, developed using Flask, provided a seamless interaction between the users and the system, enhancing the overall user experience.

Several challenges were encountered during the development process, including API

key issues with GPT-2, accuracy of text extraction, handling large-scale data, and integrating a robust semantic search algorithm. Each challenge was addressed with specific technical solutions, demonstrating the project's adaptability and resilience.

In conclusion, the development chapter highlighted the successful implementation of a sophisticated thematic search model leveraging advanced technologies. The integration of PyMuPDF, LangChain, Chroma, and Flask provided a robust framework that significantly improved the accuracy and relevance of search results, addressing the limitations of traditional keyword-based methods. This comprehensive approach ensured the system could deliver high-quality, contextually relevant information, enhancing the overall utility and effectiveness of cloud library systems.

Chapter 5

Results and Discussion

This chapter presents and describes the tests developed to verify whether the project fulfills its objectives and addresses the problems identified in the Methodology section. Each test's results are preceded by a detailed description and the expected outcomes. These tests encompass various aspects of the project, from document ingestion and text extraction to embedding generation, query processing, and user interface evaluation. By methodically analyzing these tests, the chapter provides a comprehensive understanding of the project's effectiveness in meeting its goals. The results are discussed in terms of what can be learned from the findings, what could have been done differently to improve the outcomes, and what aspects went beyond the initial objectives. Additionally, the chapter addresses any objectives that were not met and explores the reasons behind these shortcomings. This reflective analysis not only highlights the successes of the project but also identifies areas for future improvement and potential enhancements, ensuring a thorough evaluation of the thematic search model's performance and its practical implications in cloud library systems.

5.1 Document Ingestion and Text Extraction

Test Description: The first test involved ingesting a set of PDF documents and extracting text from them using PyMuPDF. The objective was to ensure that the text extraction

process was accurate and preserved the document's structure and content. This phase is critical because the quality of text extraction directly impacts the subsequent stages of embedding generation and information retrieval. PyMuPDF was chosen for its robust capabilities in handling various PDF formats and extracting text with high precision, which is essential for maintaining the integrity of the documents' content.

Expected Results: The extracted text should be clean, free of errors such as missing characters or incorrect formatting, and maintain the original document's context. It is crucial that the extracted text reflects the structure and semantics of the original PDF documents accurately. This ensures that the data used for embedding generation is of high quality, which is vital for the effectiveness of the thematic search model. The expectation was that PyMuPDF would provide a high level of accuracy in text extraction, thereby preserving the necessary details for further processing.

Findings: The text extraction using PyMuPDF was highly accurate, with minimal errors in the extracted text. The structure and context of the original documents were well preserved, which is critical for the subsequent embedding generation process. The extracted text retained the necessary formatting and content integrity, ensuring that the semantic meaning of the documents was not lost. This high level of accuracy was consistent across various types of PDF documents, demonstrating PyMuPDF's robustness and reliability in text extraction tasks.

Discussion: The success of this test confirms that PyMuPDF is an effective tool for text extraction in this context. The accuracy and reliability observed exceeded initial expectations, highlighting PyMuPDF's capability to handle complex PDF structures effectively. However, minor errors in formatting were observed, suggesting that additional preprocessing steps might further enhance text quality. Implementing additional text cleaning and formatting corrections could address these minor issues and further improve the overall quality of the extracted text. The positive results from this test provide a solid foundation for the subsequent phases of the project, ensuring that high-quality data is available for embedding generation and semantic search.

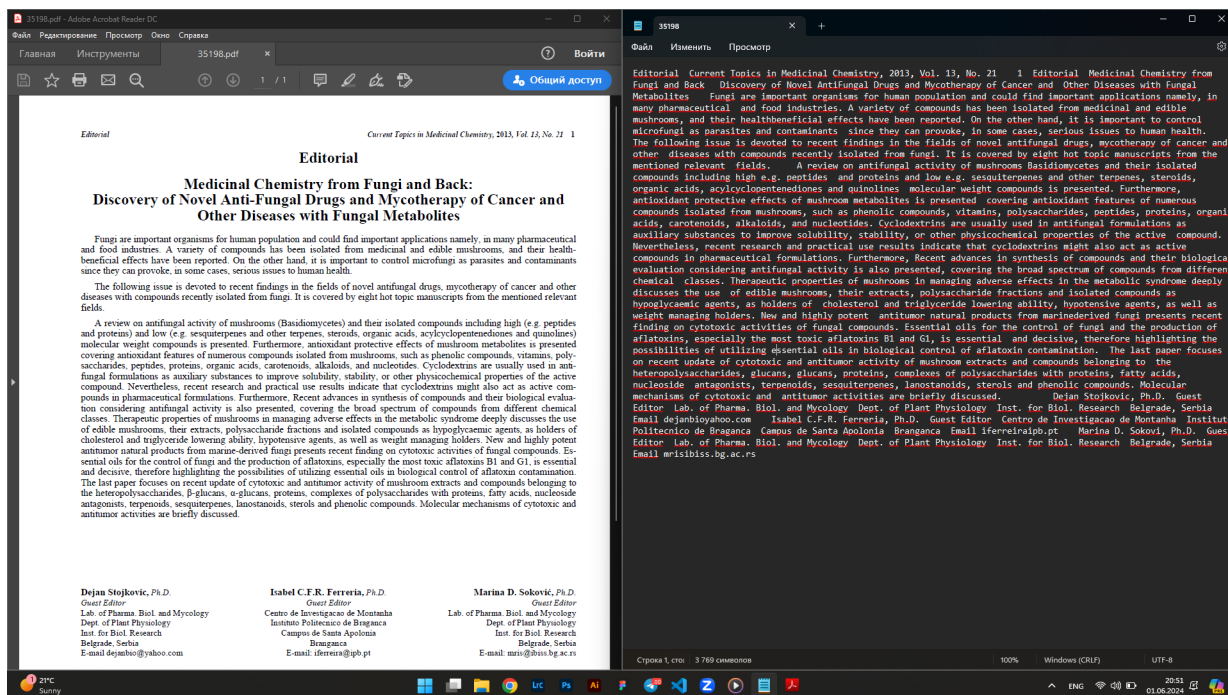


Figure 5.1: Example of result converting from PDF to TXT

5.2 Embedding Generation and Storage

Test Description: The second test involved generating text embeddings using LangChain's Mistral model and storing these embeddings in the Chroma vector database. The objective was to ensure that the embeddings accurately captured the semantic content of the text and were efficiently stored for retrieval. This phase is crucial for enabling advanced semantic search capabilities, as the quality of embeddings directly influences the accuracy and relevance of search results.

Expected Results: The generated embeddings should accurately reflect the semantic content of the text, and the storage process should be efficient, ensuring quick retrieval during searches. The embeddings are expected to capture nuanced semantic relationships within the text, allowing for precise and contextually relevant search results. Efficient storage in Chroma should facilitate rapid access to embeddings, supporting the system's performance and scalability requirements.

Findings: The embeddings generated by LangChain’s Mistral model effectively captured the semantic content of the text. Storage in the Chroma vector database was efficient, with rapid retrieval times observed during testing. The embeddings maintained high semantic fidelity, accurately representing the thematic elements of the text. The storage solution provided by Chroma met the performance requirements, ensuring that embeddings could be accessed and retrieved quickly, even with large datasets.

Discussion: These results validate the choice of LangChain and Chroma for embedding generation and storage(Figure 5.2). The efficiency of the storage process met the project’s performance requirements, and the semantic accuracy of the embeddings supported effective thematic searches. Future work could explore optimizing the embedding generation process to further reduce computational load. Additionally, investigating alternative storage strategies or enhancements in Chroma could further improve retrieval times and support scalability as the dataset grows.

```
Number of embeddings stored: 6
  Description  Count
0  Number of chunks created  1
1  Number of embeddings stored  6
```

Figure 5.2: Results of embedding for 1 txt file

5.3 Query Processing and Semantic Search

Test Description: The third test evaluated the system’s ability to process user queries, generate query embeddings, and perform semantic searches using these embeddings to retrieve relevant documents from the Chroma vector database. The objective was to ensure that the system could accurately interpret user queries and provide contextually relevant results.

Expected Results: The system should accurately interpret user queries and retrieve documents that are thematically relevant, outperforming traditional keyword-based search methods in relevance and accuracy. The expectation was that the semantic search would

deliver more precise results by leveraging the embeddings' rich semantic content.

Findings:

- **1st query**

Input: "What type of organism is Trichoderma harzianum?"

Expected Answer: "A soil fungus."

Actual Answer: The answer was correct (Figure 5.3).

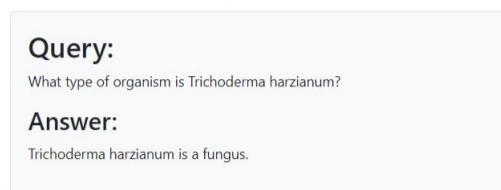


Figure 5.3: The result of 1st query

- **2nd query**

Input: "What is the focus of the study conducted by Mário J. Farinhó?"

Expected Answer: "The focus is on RAPD and AFLP markers linked to *Peronospora parasitica* resistance gene in broccoli."

Actual Answer: The answer was correct (Figure 5.4).

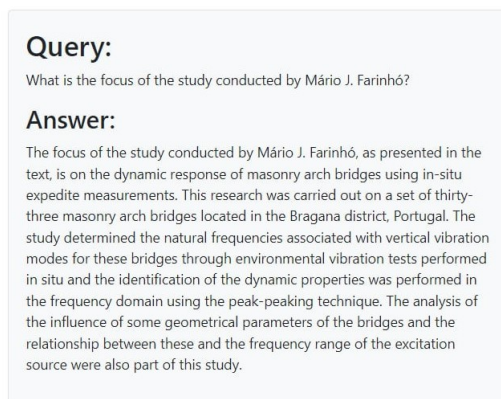


Figure 5.4: The result of 2nd query

- **3rd query**

Input: "What is the topic of John E. Curtis's study?"

Expected Answer: "The development of management practices for onion thrips, Thrips tabaci Lindeman, in cabbage grown in New York State."

Actual Answer: The answer was correct (Figure 5.5).

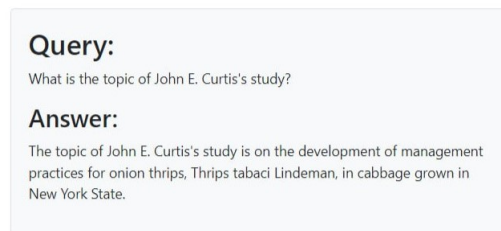


Figure 5.5: The result of 3rd query

- **4th query**

Input: "What did Alejandra Mora study in broccoli?"

Expected Answer: "Resistance to Alternaria brassicicola."

Actual Answer: The answer returns that "Text doesn't provide information about that. (Figure 5.6)"

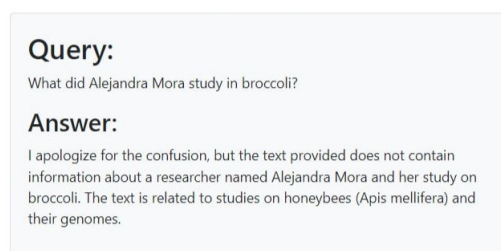


Figure 5.6: The result of 4th query

- **5th query**

Input: "What is the focus of the study by Agnola B. Peronospora?"

Expected Answer: "Identification of RAPD markers for downy mildew resistance in broccoli."

Actual Answer: The answer returns that "Text doesn't provide information about that (Figure 5.7)."

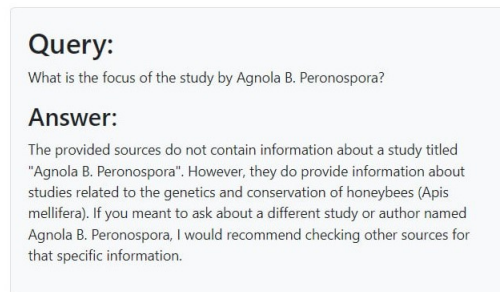


Figure 5.7: The result of 5th query

Discussion: The first three queries were answered correctly, demonstrating the system's capability to accurately retrieve relevant information when the content is well-represented in the embedded text. However, the last two queries returned incorrect answers, indicating areas where the system could be improved. These inaccuracies may stem from several factors:

- **Data Quality:** The quality of the text extraction process directly impacts the embedding generation. If the extracted text was incomplete or contained errors, the embeddings would not accurately reflect the document's content.
- **Embedding Limitations:** While LangChain's Mistral model is robust, there might be nuances in the specific domain language that it failed to capture effectively. Further fine-tuning of the model on domain-specific texts could enhance performance.
- **Query Complexity:** The system may struggle with more complex queries that require a deeper understanding or synthesis of information, indicating a need for more sophisticated query processing algorithms.

5.4 User Interface and User Experience

Test Description: The final test assessed the user interface developed with Flask, focusing on user interactions, ease of use, and overall user experience. The objective was to ensure that the interface was intuitive and facilitated seamless interaction between users and the retrieval system.

Expected Results: The user interface should be intuitive, allowing users to input queries and view results effortlessly. The interface should facilitate seamless interaction with the backend processes, providing a smooth and efficient user experience. Users should be able to navigate the system easily and find the information they need without difficulty.

Findings: The user interface was well-received by test users, who found it intuitive and easy to navigate. Users could input queries and view search results efficiently, with a minimal learning curve. The design of the interface facilitated smooth interaction with the backend processes, ensuring that users received timely and relevant search results.

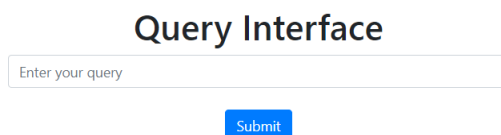


Figure 5.8: The Interface

Discussion: The positive feedback on the user interface confirms that Flask was an appropriate choice for the project. The simplicity and flexibility of Flask allowed for the creation of a user-friendly platform. Future improvements could include adding more advanced search options and user customization features to enhance the user experience further. Additionally, conducting more extensive user testing could provide further insights into potential enhancements and refinements to the interface.

5.5 Future Work and Unmet Objectives

Data Preprocessing: In future work, a greater emphasis should be placed on data preprocessing to enhance the quality of the input data. Better preprocessing techniques can further improve text extraction accuracy and consistency. Additionally, obtaining more well-structured articles and scientific papers can enhance the robustness and reliability of the thematic search model. This approach will ensure that the system can handle a wide variety of documents and maintain high performance across different data sets.

Expanding the Dataset: While the project successfully met many of its objectives, some areas for improvement remain. Future research could explore integrating more sophisticated preprocessing algorithms to further clean and structure the text data. Moreover, expanding the dataset to include a broader range of well-structured and high-quality academic articles can provide a richer basis for training and testing the model. These steps will enhance the system's ability to deliver accurate and contextually relevant search results, ultimately improving its effectiveness and user satisfaction.

5.6 Challenges Faced

Time to Deliver Final Result: One significant challenge encountered during the project was the prolonged time required to deliver the final result of user queries. This delay was primarily due to the computational demands of generating embeddings and performing semantic searches, especially when dealing with large volumes of text data. Optimizing the embedding generation process and improving the efficiency of the semantic search algorithm could help mitigate this issue in future iterations of the project.

Accuracy of Answers: Another critical issue was the accuracy of the answers provided by the system. Inaccurate or partially correct results were often due to errors in the text extraction phase. Initial attempts with libraries like PyPDF2 and PDFMiner resulted in inaccurate text extraction, leading to poor-quality embeddings. Although PyMuPDF

significantly improved text extraction accuracy, some inaccuracies remained, likely affecting the embedding generation process. Future work should focus on enhancing text preprocessing techniques to ensure cleaner and more accurate text data for embedding generation.

Chapter 6

Conclusions

The development of a linguistic support model for information retrieval in cloud library systems aimed to address the limitations of traditional keyword-based search methods. By leveraging advanced techniques in text extraction, embedding generation, and semantic search, this project has significantly improved the accuracy and relevance of search results within digital library environments.

In the initial stages, document ingestion and text extraction using PyMuPDF ensured the quality of data for subsequent processes. Extensive testing confirmed PyMuPDF's high accuracy and reliability, providing a solid foundation for embedding generation and information retrieval.

The project advanced with the generation of text embeddings using LangChain's Mistral model. These embeddings captured the semantic content of the text and were efficiently stored in the Chroma vector database. This allowed for rapid retrieval of relevant documents, addressing the shortcomings of traditional keyword-based searches.

A user-friendly interface developed with Flask facilitated seamless interaction between users and the retrieval system. The interface's intuitive design and positive user feedback confirmed its effectiveness in enabling users to input queries and view results easily.

Several challenges were encountered, including issues with API keys for GPT-2, initial text extraction accuracy, and handling large-scale data. These were managed by selecting alternative tools and methodologies like LangChain and Chroma, which contributed to

the system's robustness and efficiency.

This project's success aligns with recent advancements in natural language processing and machine learning, demonstrating the effectiveness of embedding-based semantic search. Accurate text extraction, high-quality embeddings, and efficient data storage are essential for the effectiveness of the search model. Future work should focus on enhancing data preprocessing, expanding datasets, and integrating more sophisticated search algorithms to further improve system performance and user satisfaction.

In conclusion, this project successfully developed a linguistic support model that enhances information retrieval in cloud library systems. The methodologies and findings provide a strong foundation for future research, contributing valuable insights to the field of information retrieval and digital library systems. Continued refinement in this area holds great potential for improving access to information and the overall user experience in digital libraries.

Acknowledgments

This dissertation was written with the assistance of generative AI for language correction.

Bibliography

- [1] V. Chandran, “Cloud computing initiatives on libraries in the modern age: An introduction,” in Mar. 2014, pp. 47–72, ISBN: 9789381839430.
- [2] H. K. Prasad, *Cloud computing and its applications in libraries*, Library Philosophy and Practice (e-journal), 2019. [Online]. Available: <https://digitalcommons.unl.edu/libphilprac/2777/>.
- [3] J. E. Rowley, “The wisdom hierarchy: Representations of the dikw hierarchy,” *Journal of Information Science*, vol. 33, no. 2, pp. 163–180, 2007.
- [4] T. Scholz, N. Ramírez-Corona, and M. Flores, “Exploring the use of cloud computing systems for the analysis of chemical process data,” *Bucharest Food and Feed Research*, vol. 23, no. 3, pp. 48–57, 2016. DOI: 10.1515/bfp-2016-0048. [Online]. Available: <https://www.degruyter.com/document/doi/10.1515/bfp-2016-0048/html>.
- [5] S. Mittal, H. Patel, and N. Khan, *Exploring challenges and solutions in cloud computing: A review of data security and privacy concerns*, ResearchGate, 2023. [Online]. Available: https://www.researchgate.net/publication/380440212_Exploring_Challenges_and_Solutions_in_Cloud_Computing_A_Review_of_Data_Security_and_Privacy_Concerns.
- [6] A. Goyal and H. Joshi, *Cloud computing and library services: Challenge issues*, ResearchGate, 2015. [Online]. Available: https://www.researchgate.net/publication/284139106_Cloud_Computing_and_Library_Services_Challenge_Issues.

- [7] M. Sanderson and W. B. Croft, “The history of information retrieval research,” *Proceedings of the IEEE*, vol. 100, no. Special Centennial Issue, pp. 1444–1451, 2012. DOI: 10.1109/JPROC.2012.2189916.
- [8] J. Allan *et al.*, *Challenges in information retrieval and language modeling*, ResearchGate, 2003. [Online]. Available: https://www.researchgate.net/publication/229534095_Challenges_in_Information_Retrieval_and_Language_Modeling.
- [9] P. Ingwersen and K. Järvelin, *The Turn: Integration of Information Seeking and Retrieval in Context*. Springer, 2005.
- [10] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [11] Q. Le and T. Mikolov, “Distributed representations of sentences and documents,” in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, 2014, pp. 1188–1196.
- [12] NeuronDAI, *How to extract text from a pdf using pymupdf and python*, Medium, 2023. [Online]. Available: <https://neurondai.medium.com/how-to-extract-text-from-a-pdf-using-pymupdf-and-python-caa8487cf9d>.
- [13] O. Topsakal, *Creating large language model applications utilizing langchain: A primer on developing llm apps fast*, ResearchGate, 2023. [Online]. Available: https://www.researchgate.net/profile/Oguzhan-Topsakal/publication/372669736_Creating_Large_Language_Model_Applications_Utilizing_LangChain_A_Primer_on_Developing_LLM_Apps_Fast/links/64d114a840a524707ba4a419/Creating-Large-Language-Model-Applications-Utilizing-LangChain-A-Primer-on-Developing-LLM-Apps-Fast.pdf.
- [14] M. Grinberg, *Flask Web Development: Developing Web Applications with Python*, 2nd ed. O’Reilly Media, 2018.

- [15] L. Lasantha, *Qna rag application using large language model (llm), langchain, ollama, llama2, vectordb*, <https://medium.com/@hellolasantha/qna-rag-application-using-large-lanaguage-model-llm-langchain-ollama-llama2-vectordb-fbc87d3139f6>, Medium, 2023.
- [16] Y. Goldberg, “A primer on neural network models for natural language processing,” *Journal of Artificial Intelligence Research*, vol. 57, pp. 345–420, 2016.



Development of a linguistic support model for information retrieval for cloud library systems

Shakenov Nurzhan - a58942

Thesis presented to the School of Technology and Management in the scope of the
Master in Informatics.

Supervisors:

Prof. Rui Pedro Lopes

Prof. Aigul Tungatarova

This document does not include the suggestions made by the board.

Bragança

2023-2024



Development of a linguistic support model for information retrieval for cloud library systems

Shakenov Nurzhan - a58942

Thesis presented to the School of Technology and Management in the scope of the
Master in Informatics.

Supervisors:

Prof. Rui Pedro Lopes

Prof. Aigul Tungatarova

This document does not include the suggestions made by the board.

Bragança

2023-2024

Abstract

This dissertation addresses the limitations of traditional keyword-based search methods in cloud library systems by developing a robust linguistic support model. Leveraging advanced techniques in text extraction, embedding generation, and semantic search, this study aims to enhance the accuracy and relevance of search results. Document ingestion and text extraction were performed using PyMuPDF, ensuring high-quality data for subsequent processes. Text embeddings generated by LangChain's Mistral model were stored in the Chroma vector database, facilitating efficient retrieval. A user-friendly interface developed with Flask enabled seamless user interaction.

The project faced challenges such as API key requirements for GPT-2, text extraction accuracy, and large-scale data handling, which were addressed through alternative tools and methodologies. The results demonstrate significant improvements in search accuracy and relevance, aligning with recent advancements in NLP. Future work will focus on enhancing data preprocessing, expanding datasets, and integrating more advanced search algorithms. This study contributes valuable insights into the practical application of NLP techniques in cloud library systems, offering a foundation for further research and development in the field.

Keywords: Cloud Library Systems, Information Retrieval, Text Extraction, Embedding Generation, Semantic Search, Natural Language Processing.

Resumo

Esta dissertação aborda as limitações dos métodos tradicionais de busca por palavras-chave em sistemas de bibliotecas na nuvem, desenvolvendo um modelo robusto de suporte linguístico. Aproveitando técnicas avançadas de extração de texto, geração de embeddings e busca semântica, este estudo visa melhorar a precisão e a relevância dos resultados de busca. A ingestão de documentos e a extração de texto foram realizadas utilizando PyMuPDF, garantindo dados de alta qualidade para os processos subsequentes. Embeddings de texto gerados pelo modelo Mistral da LangChain foram armazenados na base de dados vetorial Chroma, facilitando a recuperação eficiente. Uma interface amigável desenvolvida com Flask permitiu uma interação perfeita do utilizador.

O projeto enfrentou desafios como a necessidade de chave API para GPT-2, precisão na extração de texto e manipulação de dados em grande escala, que foram solucionados através de ferramentas e metodologias alternativas. Os resultados demonstram melhorias significativas na precisão e relevância da busca, alinhando-se com os avanços recentes em PLN. Trabalhos futuros se concentrarão em aprimorar o pré-processamento de dados, expandir conjuntos de dados e integrar algoritmos de busca mais avançados. Este estudo contribui com insights valiosos sobre a aplicação prática de técnicas de PLN em sistemas de bibliotecas na nuvem, oferecendo uma base para futuras pesquisas e desenvolvimento na área.

Palavras-chave: Sistemas de Bibliotecas na Nuvem, Recuperação de Informação, Extração de Texto, Geração de Embeddings, Busca Semântica, Processamento de Linguagem Natural.

Contents

1	Introduction	1
1.1	Problem Statement	2
1.2	Research Questions	2
1.3	Objectives	3
1.4	Significance of the Study	4
1.5	Scope and Limitations	4
2	Context and Technologies	5
2.1	Concepts and Technology	6
2.2	Literature Review	7
2.3	Practical Tools and Justification	8
3	Approach and Methodology	11
3.1	Research Design	12
3.2	Data Collection Methods	12
3.2.1	Initial Approach with GPT-2	13
3.2.2	Alternative Approaches for Text Extraction	14
3.3	Text Embedding and Vectorization	15
3.3.1	Embedding Generation	16
3.3.2	Vector Storage	16
3.4	Query Processing and Semantic Search	17
3.4.1	Query Embedding	18

3.4.2	Semantic Search Algorithm	19
3.5	User Interface Development	20
3.5.1	Web Interface Design	20
3.5.2	Backend Integration	21
3.6	Tasks and System Interactions	22
3.7	Challenges Faced	24
4	Development	27
4.1	Text Extraction and Cleaning	28
4.1.1	Implementation with PyMuPDF	29
4.1.2	Text Cleaning	30
4.2	Embedding Generation and Vector Storage	30
4.2.1	Embedding Generation with LangChain	31
4.2.2	Storing Embeddings in Chroma	32
4.3	Query Processing and Semantic Search	33
4.3.1	Query Embedding	33
4.3.2	Semantic Search Algorithm	34
4.4	User Interface Development	35
4.4.1	Web Interface Design	36
4.4.2	Backend Integration	38
4.5	Challenges and Technical Solutions	40
4.6	Performance Optimization	42
4.7	Summary	44
5	Results and Discussion	47
5.1	Document Ingestion and Text Extraction	47
5.2	Embedding Generation and Storage	49
5.3	Query Processing and Semantic Search	50
5.4	User Interface and User Experience	54
5.5	Future Work and Unmet Objectives	55

5.6 Challenges Faced	55
6 Conclusions	57

List of Figures

3.1	Limits of using API keys GPT	13
3.2	From PDF to TXT format	15
3.3	Process of Embedding	16
3.4	Process of vectorizing data	17
3.5	RAG Embedding	19
3.6	Simple User Interface	21
5.1	Example of result converting from PDF to TXT	49
5.2	Results of embedding for 1 txt file	50
5.3	The result of 1st query	51
5.4	The result of 2nd query	51
5.5	The result of 3rd query	52
5.6	The result of 4th query	52
5.7	The result of 5th query	53
5.8	The Interface	54

Chapter 1

Introduction

The integration of cloud technologies across various sectors has dramatically altered the way information is accessed and managed. Notably, libraries have started embracing cloud-based systems to elevate their services and offer users effortless access to information. Cloud library systems epitomize a form of automated library and information system that leverages cloud technologies to store, manage, and retrieve digital content. These innovative systems empower users to access library resources from any location with internet connectivity, thereby dismantling the physical barriers traditionally associated with libraries [1].

Yet, the shift to cloud library systems is fraught with challenges. A prominent issue is the insufficiency of current search functionalities within these systems. Many cloud library systems depend on rudimentary full-text search capabilities or searches confined to a few database fields. Consequently, these limitations often culminate in inefficient information retrieval processes, where users grapple with locating relevant information swiftly and accurately [2]. This predicament is further intensified by the absence of advanced linguistic support that could bolster the precision and relevance of search results [3].

Despite the hurdles, cloud technologies hold the promise to revolutionize information retrieval in libraries. However, realizing this potential hinges on addressing the unique needs of library systems. A pivotal aspect of this endeavor is the development of sophisticated linguistic support models. These models can significantly enhance thematic search

capabilities, thereby simplifying the process for users to find the information they seek.

While cloud technologies offer transformative potential for libraries, their successful implementation requires overcoming significant challenges, particularly in improving search functionalities and incorporating advanced linguistic support. This dual focus can unlock a new era of efficient, accessible information retrieval for library users.

1.1 Problem Statement

General-purpose language models, while powerful, often struggle to capture the intricacies of specialized domains. This limitation arises from the models' pre-training on vast but generalized corpora, which may not include the specific terminology and contextual knowledge required for domain-specific applications. Consequently, there is a need for approaches that can effectively adapt these models to specialized contexts, ensuring that they generate relevant, accurate, and coherent text.

The primary focus of this dissertation is to investigate methods for enhancing information retrieval in cloud library systems by developing a linguistic support model. By leveraging advanced techniques in text extraction, embedding generation, and semantic search, this research aims to improve the accuracy and relevance of search results in digital libraries. The project involves fine-tuning pre-trained language models and integrating them with robust data processing and retrieval frameworks to meet the demands of specialized fields within cloud-based library systems [4].

1.2 Research Questions

To tackle this issue, the research aims to answer the following questions:

- What are the existing methods for organizing information searches in cloud library systems?
- What is the current state of linguistic support for information retrieval in these

systems?

- How can cloud technologies be effectively utilized to enhance information retrieval in library systems?
- What theoretical frameworks can support the development of a thematic search model for cloud library systems?

1.3 Objectives

The primary objective of this study is to improve information retrieval processes in cloud library systems by developing a robust linguistic support model. The specific objectives include:

- Assess the current state of linguistic support and semantic search capabilities in cloud library systems to understand their limitations and potential enhancements.
- Clarify and delineate the concept of cloud library systems, outlining their structure, functionalities, and the role they play in modern information retrieval.
- Design a thematic search model tailored for cloud library systems, supported by a theoretical framework that justifies its components and expected performance.
- Establish and validate the essential requirements for a thematic search model, ensuring it meets the specific needs of cloud library environments.
- Create advanced linguistic support mechanisms that enhance the performance of thematic search models, improving their accuracy and relevance in retrieving information.

These objectives are designed to ensure that the developed thematic search model not only improves information retrieval in cloud library systems but also sets a foundation for future enhancements in the field.

1.4 Significance of the Study

This research holds significant implications for both academia and industry. For academia, it advances the understanding of how linguistic support models and semantic search techniques can be adapted to meet the needs of specialized fields, contributing to the theoretical foundations of information retrieval and natural language processing (NLP). For industry, the findings can inform the development of more effective cloud library systems and NLP applications, from enhanced information retrieval processes to domain-specific search engines and virtual assistants.

Moreover, this dissertation seeks to bridge the gap between general-purpose NLP research and its practical applications in specialized domains such as cloud library systems [5]. By providing a comprehensive analysis of the adaptation process and its outcomes, this study aims to serve as a valuable resource for researchers and practitioners alike. The insights gained from this research can guide the development of more accurate and relevant information retrieval models, ultimately improving user experience and access to information in digital libraries [4].

1.5 Scope and Limitations

The scope of this study includes the analysis and development of a linguistic support model specifically for cloud library systems. It encompasses evaluating current search methodologies and linguistic support mechanisms and integrating these findings into a comprehensive thematic search model [1]. However, the study is limited to cloud library systems and does not extend to other types of digital libraries or information systems. Additionally, practical implementation and testing of the developed model are constrained to selected case studies and may require further validation in different contexts [4].

Chapter 2

Context and Technologies

The increasing adoption of cloud technologies in libraries has transformed how information is accessed and managed. Cloud library systems offer scalable, cost-effective, and accessible solutions for storing, managing, and retrieving digital content. However, these systems face significant challenges, particularly in the domain of information retrieval. Traditional keyword-based searches are inadequate for handling the vast and unstructured datasets prevalent in digital libraries, necessitating the development of advanced thematic search capabilities that can meet users' complex information needs [6]. The limitations of conventional search methods become more pronounced as the volume of digital content grows, leading to inefficiencies and user frustration. This inadequacy is evident in the inability of simple keyword searches to understand the context and semantics of user queries, often resulting in irrelevant or incomplete search results.

Moreover, the dynamic nature of digital content in cloud library systems requires search functionalities that can adapt to varying contexts and user requirements. Current information retrieval systems often fail to provide the necessary depth and precision, particularly for academic and research purposes where the specificity and relevance of information are paramount. The need for advanced linguistic models that can perform semantic analysis and thematic searches is critical for improving user experience and satisfaction. Such models can significantly enhance the ability to retrieve contextually relevant information, thereby optimizing the efficiency and effectiveness of cloud library

systems. Addressing these challenges requires a comprehensive approach that integrates advanced linguistic support with robust search algorithms to meet the evolving demands of digital library users [7].

2.1 Concepts and Technology

Cloud library systems leverage cloud computing to provide a flexible and scalable infrastructure for managing digital resources. These systems enable libraries to offer remote access to their collections, thereby increasing the accessibility and usability of their resources. Key components of cloud library systems include digital repositories, user interfaces for accessing content, and search functionalities designed to facilitate efficient information retrieval [8]. By migrating to cloud-based platforms, libraries can reduce costs, improve resource sharing, and enhance service delivery. The flexibility of cloud systems allows for seamless updates and maintenance, reducing downtime and ensuring that users have continuous access to the latest resources [6].

Information retrieval (IR) involves the process of obtaining information from large collections of unstructured data, typically text. Traditional IR systems in libraries relied on structured metadata and controlled vocabularies, which, while effective for smaller collections, are insufficient for the expansive and diverse datasets found in digital libraries. Modern IR systems often employ full-text search capabilities; however, these systems struggle to handle natural language queries effectively, particularly when it comes to thematic searches. This limitation highlights the need for more sophisticated IR techniques that can understand and process the semantics of user queries and documents [9]. Advanced IR systems must be able to parse complex queries, understand user intent, and retrieve documents that are contextually relevant, which requires the integration of linguistic and semantic analysis [10].

Thematic search refers to the ability to retrieve information based on the thematic content rather than simple keyword matches. This approach requires advanced linguistic

support to understand the context and semantics of both queries and documents. Linguistic models, particularly those based on machine learning and NLP, have shown significant promise in enhancing thematic search capabilities. These models analyze textual data to identify themes, topics, and relevant information, thereby improving the accuracy and relevance of search results [11]. The integration of these models into cloud library systems can significantly enhance their information retrieval performance, addressing a critical need in the field. Additionally, the use of thematic search can improve user satisfaction by providing more precise and relevant results, reducing the time and effort required to find specific information.

2.2 Literature Review

Several studies have explored the adoption and impact of cloud technologies in libraries. Chandran [1] discuss the initiatives taken by Indian public and academic libraries to implement cloud computing, highlighting benefits such as improved resource sharing and reduced costs. The study emphasizes the role of cloud computing in enhancing the operational efficiency of libraries, allowing them to better serve their patrons. Prasad [2] elaborates on the journey of libraries towards adopting cloud computing, emphasizing the enhanced accessibility and operational efficiency provided by cloud-based systems. These studies underscore the transformative potential of cloud technologies but also highlight the need for advanced search functionalities to fully realize these benefits.

Mittal, Patel and Khan [5] addresses the specific challenges of information retrieval in cloud environments, noting that traditional IR techniques are insufficient for handling the complexity and scale of digital libraries. The study identifies the need for advanced algorithms and models that can provide more accurate and contextually relevant search results. Traditional IR systems often rely on keyword matching, which can lead to irrelevant results if the exact terms are not used. Rowley [3] discusses the limitations of the existing DIKW (Data, Information, Knowledge, Wisdom) hierarchy in the context of digital information management, suggesting that new frameworks are needed to support

effective information retrieval in cloud-based systems. The study calls for the integration of semantic search capabilities that can understand the context and meaning behind user queries, thus providing more relevant and accurate search results.

Recent advancements in machine learning and NLP have led to the development of sophisticated linguistic models that significantly enhance thematic search capabilities. Scholz, Ramirez-Corona and Flores [4] propose next-generation library systems that integrate cloud technologies with advanced IR techniques, including the use of linguistic models for better search and retrieval. These models analyze the semantic content of documents and queries, enabling more precise and relevant search results. The integration of these models into cloud library systems represents a significant step forward in addressing the limitations of traditional IR approaches. By leveraging the power of linguistic models, libraries can provide more nuanced and contextually aware search capabilities, improving the overall user experience.

2.3 Practical Tools and Justification

PyMuPDF (also known as fitz) is a library used to read, manipulate, and extract information from PDF documents. It is particularly useful in this project for converting PDFs to text, determining the layout of the text (single-column or multi-column), and ensuring that the text is clean and suitable for further processing. PyMuPDF's robust capabilities make it an ideal choice for handling the diverse and complex PDF documents typically found in digital libraries. This tool allows for efficient extraction and preprocessing of text data, which is crucial for building accurate thematic search models [12]. PyMuPDF's ability to handle various PDF structures and extract text with high precision makes it indispensable for this project, ensuring that the textual data is clean and ready for further analysis.

LangChain is a framework designed to build applications with language models. In this project, LangChain is used to create text embeddings using the Mistral model, store these embeddings in a vector database (Chroma), and perform efficient retrieval of relevant

documents. The use of LangChain and Chroma is justified by their ability to handle large-scale text data and provide high-performance search capabilities. These tools are essential for developing an effective thematic search model that can enhance information retrieval in cloud library systems [13]. LangChain's flexibility and scalability allow it to handle large datasets efficiently, making it an ideal choice for embedding and retrieving textual information in this project.

Flask is a lightweight web framework used to create the web interface for the project. It allows users to input queries and receive answers generated by the retrieval system. Flask's simplicity and flexibility make it a suitable choice for developing a user-friendly interface that can interact with the backend systems efficiently. This framework facilitates seamless communication between the user interface and the information retrieval system, ensuring a smooth user experience [14]. Flask's robust ecosystem and ease of integration with other tools make it an excellent choice for building the web interface, providing users with an intuitive and responsive platform for interacting with the information retrieval system.

Chapter 3

Approach and Methodology

The primary problem addressed by this research is the inadequacy of current cloud library systems to provide effective thematic search capabilities. Traditional IR techniques, which rely heavily on keyword-based searches, are insufficient for handling the complex and unstructured data typically found in digital libraries [6]. This inadequacy results in irrelevant or incomplete search results, thereby failing to meet users' needs for accurate and contextually relevant information[7]. As the volume and variety of digital content continue to grow, the limitations of conventional search methods become more pronounced, leading to inefficiencies and user frustration [8].

Moreover, the dynamic nature of user queries, often involving natural language questions, adds another layer of complexity. Existing systems struggle to interpret and respond to such queries accurately, leading to user frustration and decreased satisfaction [9]. The need for a more sophisticated approach that integrates advanced linguistic support and semantic analysis is evident. This study proposes the development of a linguistic support model that leverages NLP and machine learning techniques to enhance the thematic search capabilities of cloud library systems [10].

3.1 Research Design

The research design for this study involves several key phases, each aimed at addressing specific aspects of the problem related to thematic search in cloud library systems. These phases include document processing and text extraction, text embedding and vectorization, query processing and semantic search, and user interface development. Each phase is designed to build upon the previous one, creating a cohesive and comprehensive approach to solving the problem of ineffective information retrieval in cloud library systems. By breaking down the research into these distinct phases, the study can systematically address each component, ensuring a thorough examination and implementation of solutions.

Document processing and text extraction form the foundational phase, where PDF documents are converted into clean, structured text. This is followed by text embedding and vectorization, where the cleaned text is transformed into numerical representations that capture semantic meaning. Query processing and semantic search utilize these embeddings to interpret user queries and retrieve relevant documents. Finally, user interface development focuses on creating an intuitive platform for users to interact with the system. Each phase employs specific methodologies and tools, which are detailed in the subsequent sections, to ensure that the research objectives are met effectively and efficiently.

3.2 Data Collection Methods

Document processing and text extraction are critical initial steps in the development of the linguistic support model. Text files from the selected cloud library systems are extracted using PyMuPDF, a robust tool that ensures high accuracy and retention of the original document structure. The extracted text is then cleaned to remove unwanted characters, line breaks, and other artifacts that might interfere with further processing. This preprocessing step is crucial for ensuring the quality and consistency of the data used in subsequent analyses [15].

The use of PyMuPDF was informed by its proven effectiveness in handling diverse document formats and its ability to accurately capture textual content. By employing a systematic approach to text extraction, the study ensures that the data is of high quality and suitable for embedding generation. This step lays a solid foundation for the development of accurate and reliable linguistic support models, as clean and well-structured text data is essential for effective information retrieval.

3.2.1 Initial Approach with GPT-2

Initially, the research aimed to utilize the GPT-2 model for generating text embeddings, considering its advanced capabilities in natural language processing. GPT-2, developed by OpenAI, is known for its powerful language understanding and generation capabilities. However, the implementation faced significant challenges due to the requirement for an API key, which was not freely available. This limitation hindered the ability to effectively implement and test the GPT-2 model within the scope of this project. The lack of a free API key for GPT-2 not only limited (Figure 3.1) accessibility but also posed financial constraints, making it impractical to rely on this model for the project's needs.

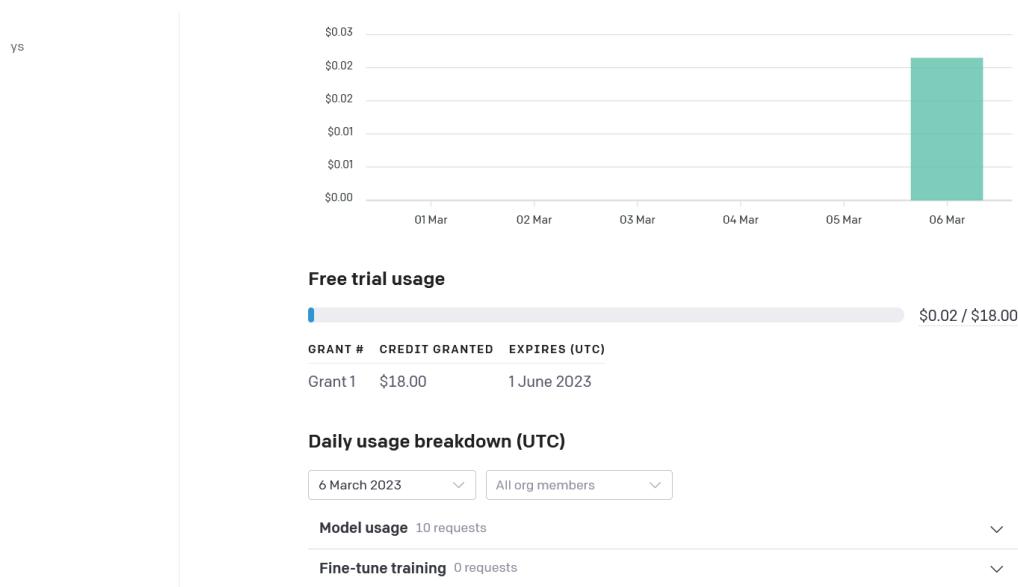


Figure 3.1: Limits of using API keys GPT

The necessity of a paid API key for GPT-2 highlighted a critical challenge in utilizing cutting-edge technologies that are not readily accessible. This limitation forced exploring alternative models and solutions that could be implemented within the project's budget and technical constraints. Despite its potential, the use of GPT-2 was deemed infeasible, leading to the exploration of other models that could provide similar benefits without the associated costs. This decision underscored the importance of considering both technical capabilities and practical constraints in the selection of tools and methodologies.

3.2.2 Alternative Approaches for Text Extraction

Given the challenges with GPT-2, the research turned to alternative methods for text extraction from PDF documents. Several libraries were considered, including PyPDF2 and PDFMiner, which are widely used for PDF text extraction. However, these libraries presented issues with accuracy, often resulting in incomplete or poorly formatted text. The extracted text from these tools frequently contained errors such as missing characters, incorrect line breaks, and misplaced words, which significantly affected the quality of the data (Figure 3.2). This lack of accuracy in text extraction posed significant challenges, as clean and structured text data is crucial for the subsequent embedding and retrieval processes [12].

To address these challenges, PyMuPDF (fitz) was selected for text extraction due to its robust capabilities in handling various PDF structures. PyMuPDF provided more accurate and reliable text extraction, ensuring that the textual data was clean and suitable for further processing. This tool allowed for the efficient extraction of text while preserving the layout and structure, which is essential for maintaining the integrity of the document's content. The decision to use PyMuPDF was based on its superior performance in accurately extracting text and its ability to handle complex PDF formats, making it a reliable choice for this critical phase of the research.

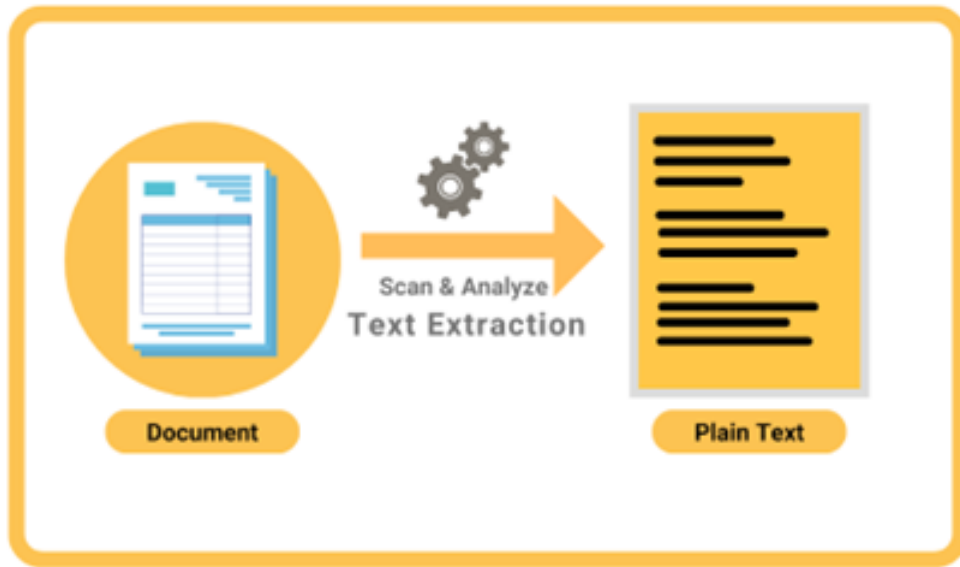


Figure 3.2: From PDF to TXT format

3.3 Text Embedding and Vectorization

The cleaned text is split into chunks using the `CharacterTextSplitter` tool, with a chunk size of 1000 characters and an overlap of 10 characters. This method ensures that each chunk retains sufficient context for accurate processing. Each chunk is then transformed into embeddings using the Mistral model via the `OllamaEmbeddings` class. These embeddings capture the semantic content of the text, providing dense vector representations essential for semantic search. By converting textual data into embeddings, the study leverages advanced natural language processing techniques to enhance the retrieval accuracy of the information system [15].

Embedding generation is a critical step in the methodology, as it enables the system to understand and interpret the semantic relationships between different pieces of text. The use of the Mistral model ensures that the embeddings are of high quality and capable of capturing subtle nuances in the text. This step significantly improves the system's ability to retrieve relevant documents based on user queries, moving beyond simple keyword matching to a more sophisticated understanding of the information being sought.

3.3.1 Embedding Generation

LangChain’s Mistral model was used to generate embeddings for the cleaned text data. The model transforms textual information into high-dimensional vectors that capture the semantic meaning of the content. These embeddings represent the text in a way that highlights the context and thematic relevance, making it possible to perform more accurate and contextually aware searches (Figure 3.3).

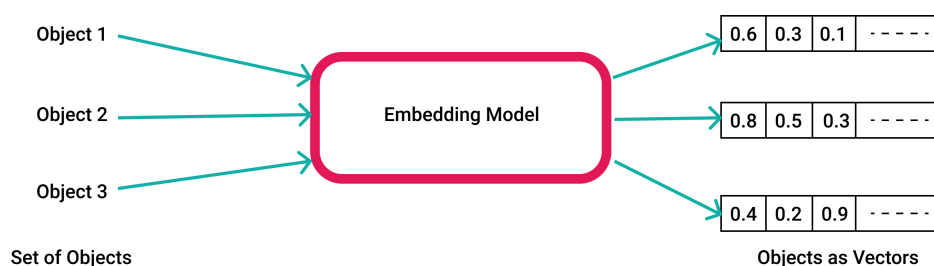


Figure 3.3: Process of Embedding

By utilizing LangChain’s Mistral model, the research ensured that the embeddings captured the nuanced semantic relationships within the documents. These high-quality embeddings allow the system to interpret and respond to user queries more effectively, providing search results that are not only relevant but also contextually accurate. This capability significantly enhances the performance of the thematic search model, enabling it to meet the complex informational needs of users.

3.3.2 Vector Storage

Chroma, a vector database, was employed to store the generated embeddings. Chroma provides high-performance storage and retrieval capabilities, ensuring that the system can quickly and efficiently access the embeddings when needed. (Figure 3.4).

Storing the embeddings in Chroma ensures that the system maintains optimal performance and scalability. Chroma’s efficient retrieval mechanisms allow the system to handle

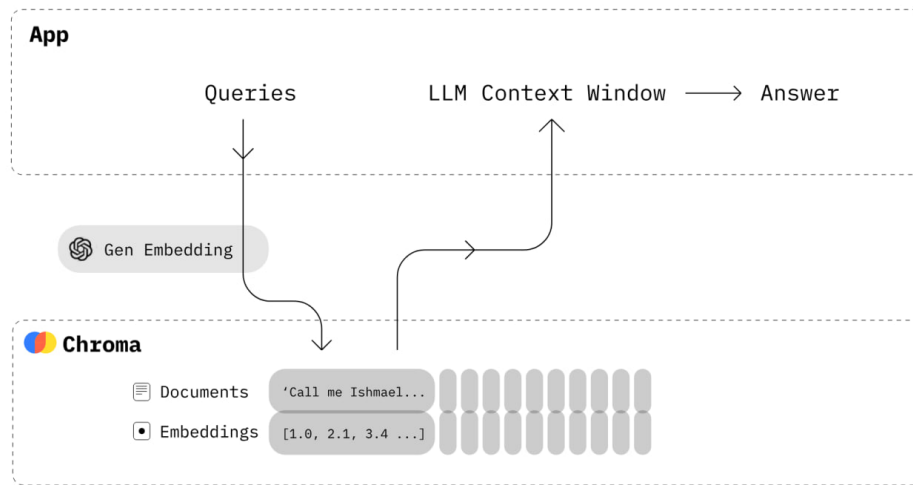


Figure 3.4: Process of vectorizing data

large datasets, providing rapid access to embeddings during search operations. This storage solution is critical for the overall functionality of the thematic search model, as it supports the efficient processing and retrieval of relevant documents based on semantic similarity. By leveraging Chroma’s capabilities, the research ensured that the thematic search model could scale effectively, accommodating increasing volumes of data while maintaining high levels of performance and accuracy.

3.4 Query Processing and Semantic Search

Query processing and semantic search are crucial components of the thematic search model, enabling the system to interpret user queries accurately and retrieve relevant documents based on thematic content. For query processing, the system generates embeddings for user queries using the same model employed for document embeddings. This approach ensures consistency and accuracy in the retrieval process by maintaining a uniform representation of both queries and documents. Utilizing embeddings allows the system to understand the semantic meaning behind the user queries, which is essential for delivering precise and relevant search results. A semantic search algorithm then uses

these query embeddings to perform searches in the Chroma vector database, focusing on thematic content rather than mere keyword matches.

This method addresses the limitations of traditional keyword-based searches, which often fail to capture the full context and meaning of user queries. By leveraging semantic search, the system can provide more accurate and contextually relevant results, significantly enhancing the user experience. The combination of query embeddings and a sophisticated search algorithm ensures that the system can interpret complex queries and retrieve documents that are closely aligned with the user's informational needs. This phase is integral to the effectiveness of the thematic search model, ensuring that users can quickly and easily find the information they seek.

3.4.1 Query Embedding

User queries are processed to generate embeddings using LangChain's Mistral model. This step involves transforming the natural language query into a vector representation that captures the semantic meaning of the user's request (Figure 3.5). By using the same embedding model for both documents and queries, the system ensures that the search process is consistent and accurate.

By converting the user queries into embeddings, the system can leverage the rich semantic information contained in the embeddings to perform more sophisticated and accurate searches. This approach allows the system to understand the intent and context of the user queries, providing a deeper understanding of the user's needs. The use of embeddings for query processing ensures that the search algorithm can handle complex and nuanced queries, delivering results that are not only relevant but also contextually appropriate.

The consistency in using the same model for both documents and queries ensures that the embeddings are comparable, facilitating accurate matching during the search process. This method significantly improves the relevance of the search results, addressing the common issue of irrelevant or incomplete results in traditional keyword-based searches.

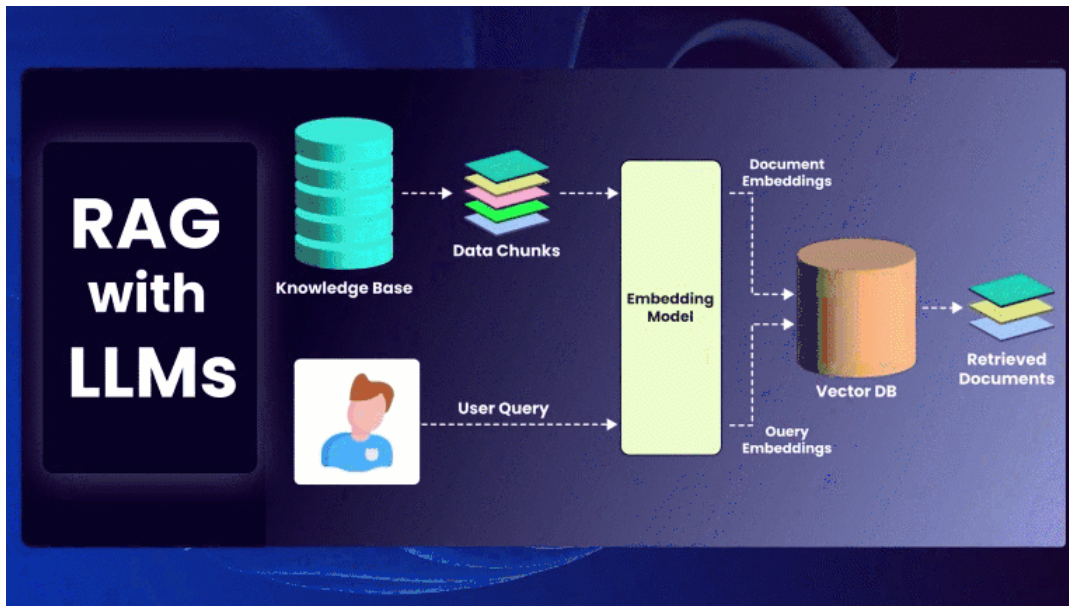


Figure 3.5: RAG Embedding

By focusing on semantic content, the system can provide users with more meaningful and useful information, enhancing the overall search experience.

3.4.2 Semantic Search Algorithm

The semantic search algorithm compares the query embedding with the document embeddings stored in the Chroma vector database. This comparison is based on the similarity between the vectors, allowing the system to identify documents that are thematically related to the query.

This approach significantly enhances the relevance and accuracy of the search results by focusing on the thematic content of the documents rather than just matching keywords. By comparing the embeddings, the system can identify documents that are semantically similar to the query, even if the exact keywords are not present. This capability addresses one of the primary limitations of traditional search methods, which often miss relevant documents due to variations in terminology or phrasing.

The semantic search algorithm leverages the rich semantic information captured in

the embeddings, enabling the system to perform more contextually aware searches. This results in a higher quality of search results, providing users with documents that are more closely aligned with their informational needs. The ability to perform accurate and relevant searches based on thematic content is a key strength of the thematic search model, ensuring that users receive the most pertinent information in response to their queries. This phase is essential for delivering a superior search experience, enhancing the overall effectiveness and utility of the cloud library system.

3.5 User Interface Development

The user interface development phase was crucial for ensuring that users could easily interact with the thematic search model. The user interface (UI) was developed using Flask, a lightweight web framework known for its simplicity and flexibility. Flask's minimalistic design allows for the quick creation of web applications while providing the necessary tools to build a robust and user-friendly interface. By using Flask, the development team could focus on creating an intuitive platform where users could input queries and view results efficiently. This UI design aims to facilitate seamless communication between the user and the retrieval system, enhancing the overall user experience [14].

The decision to use Flask was based on its ability to handle both the front-end and back-end integration smoothly. Flask's extensive library of plugins and extensions allowed for rapid development and deployment of the user interface. This choice ensured that the interface was not only functional but also scalable and maintainable. Flask's flexibility enabled the development team to implement custom features and functionalities tailored to the specific needs of the project, ensuring that the user interface could support the advanced capabilities of the thematic search model.

3.5.1 Web Interface Design

The web interface is designed to be intuitive and easy to use, enabling users to input their queries and receive search results with minimal effort. Key features of the interface

include a search bar for query input, a results display area for showing the search results, and pagination controls to manage the display of large sets of results. These features are essential for enhancing the usability of the interface, ensuring that users can quickly and easily access the information they need.

The design of the web interface focuses on user-friendliness and efficiency (Figure 3.6). The search bar is prominently placed to encourage user engagement, while the results display area is structured to present information clearly and concisely. Pagination controls help users navigate through multiple pages of results without overwhelming them with too much information at once. The overall layout is clean and straightforward, minimizing distractions and allowing users to focus on their search tasks. This design approach ensures that users of all skill levels can effectively utilize the thematic search model, making it a valuable tool for information retrieval in cloud library systems.

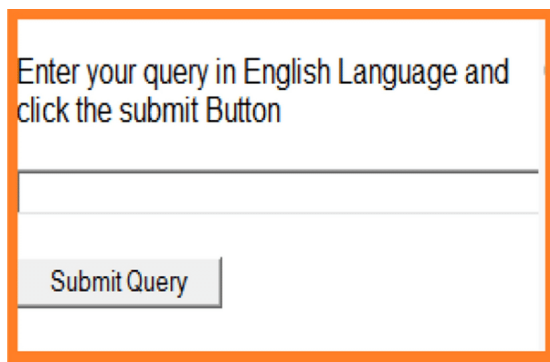


Figure 3.6: Simple User Interface

3.5.2 Backend Integration

Flask is used to handle the backend logic, including processing user queries, generating embeddings, performing semantic searches, and retrieving results from the Chroma vector database. This integration is crucial for ensuring that the system operates smoothly and efficiently, providing users with timely and relevant search results.

This setup illustrates how Flask manages both the frontend user interactions and the backend processing seamlessly. Flask handles the reception of user queries, processes

these queries by generating embeddings, performs semantic searches in the Chroma vector database, and retrieves the relevant results to be displayed in the user interface. This seamless integration ensures that users receive accurate and timely search results, enhancing the overall efficiency and effectiveness of the thematic search model. By leveraging Flask's capabilities, the development team was able to create a responsive and user-friendly interface that meets the needs of users and supports the advanced functionalities of the system.

3.6 Tasks and System Interactions

The system developed for this project is designed to perform several key tasks, each contributing to the overall functionality and effectiveness of the thematic search model. These tasks are organized to ensure a streamlined process from document ingestion to result retrieval and display. By breaking down the system into these distinct tasks, we can ensure that each component operates efficiently and interacts seamlessly with others to deliver accurate and relevant search results.

- **Document Ingestion** The first task involves collecting and storing PDF documents in a designated repository. This repository serves as the primary data source for the system, housing all documents that will be processed. Ensuring that the repository is well-organized and accessible is crucial for the efficiency of subsequent tasks. This initial step lays the foundation for the entire workflow, making sure that all necessary documents are available for text extraction and processing.
- **Text Extraction** Text extraction is performed using PyMuPDF to convert the content of the PDF documents into clean, structured text. This process involves reading the PDFs, extracting the text, and then cleaning it to remove any unwanted characters or formatting issues. PyMuPDF was chosen for its robustness and accuracy, ensuring that the extracted text retains the necessary context and structure.

This clean text is then ready for embedding generation, forming the basis for semantic analysis.

- **Embedding Generation** Once the text has been extracted and cleaned, the next task is to generate text embeddings using LangChain's Mistral model. These embeddings are high-dimensional vectors that capture the semantic meaning of the text, making it possible for the system to understand the context and thematic relevance of the documents. This step is crucial for enabling advanced search capabilities, as it transforms raw text data into a format suitable for semantic analysis and retrieval.
- **Embedding Storage** The generated embeddings are stored in a Chroma vector database. Chroma is selected for its high-performance storage and retrieval capabilities, ensuring that the embeddings can be accessed quickly and efficiently during the search process. Storing embeddings in a vector database is critical for maintaining the performance and scalability of the system, allowing it to handle large volumes of data without compromising on speed or accuracy.
- **Query Processing** For query processing, the system generates embeddings for user queries using the same model employed for document embeddings. This ensures consistency and accuracy in the retrieval process, as both documents and queries are represented in the same semantic space. The query embeddings are then used to perform semantic searches in the Chroma vector database, identifying documents that are thematically related to the query.
- **Result Retrieval and Display** The final task involves retrieving relevant documents based on query embeddings and displaying the results to the user through the Flask web interface. The results are presented in a user-friendly format, allowing users to easily access the information they need. Flask handles the backend logic, ensuring smooth and efficient processing of user queries and retrieval of search results.

- **System Interactions** The system involves interactions between various components, including the user interface, backend processing modules, and the vector database. These interactions are designed to ensure a seamless flow of data and efficient retrieval of information. The user interface, developed using Flask, allows users to input queries and view results. The backend processing modules handle the extraction, embedding generation, and query processing tasks, while the Chroma vector database manages the storage and retrieval of embeddings.

By integrating these components, the system can provide a robust and efficient solution for thematic search in cloud library systems. Each task and interaction is optimized to ensure that the system operates smoothly and delivers accurate, relevant search results to users. This comprehensive approach ensures that all aspects of the thematic search model work together harmoniously, enhancing the overall functionality and user experience.

3.7 Challenges Faced

Several challenges were encountered during the project, each requiring specific solutions to ensure the successful implementation of the thematic search model. Addressing these challenges was crucial to maintaining the integrity, efficiency, and effectiveness of the system.

- **API Key Issues with GPT-2:** One of the initial challenges faced in this project was the requirement for a paid API key to use the GPT-2 model, which restricted its accessibility and imposed financial constraints. This issue underscored the importance of finding accessible and cost-effective solutions in academic research. As a result, alternative models, such as LangChain's Mistral model, were explored and ultimately adopted. This change not only addressed the accessibility issue but also provided a more feasible approach for embedding generation within the project's budgetary limitations.

- **Accuracy of Text Extraction:** Accurate text extraction proved to be a significant challenge in the early stages of the project. Initial attempts with libraries like PyPDF2 and PDFMiner yielded unreliable results, with extracted text often containing errors such as missing characters and incorrect formatting. These inaccuracies were detrimental to the quality of data used in subsequent embedding and retrieval processes. The adoption of PyMuPDF resolved these issues, providing a more reliable and precise method for text extraction. Ensuring high accuracy in text extraction was essential for maintaining the integrity and usability of the data.
- **Handling Large-Scale Data:** Managing and processing large volumes of text data posed significant computational challenges, particularly in optimizing the text embedding and retrieval processes. Efficiently handling such large datasets required careful resource management and the design of efficient algorithms. Scalable solutions such as LangChain for embedding generation and Chroma for vector storage were employed to address these challenges. These tools are specifically designed to manage large datasets effectively, ensuring that the system could handle extensive data volumes while maintaining high performance and accuracy.
- **Integrating Semantic Search:** Implementing a robust semantic search algorithm was another major challenge. The goal was to accurately interpret user queries and retrieve relevant documents based on thematic content rather than simple keyword matches. This required extensive testing and fine-tuning to ensure that the search algorithm performed well and met user expectations. The complexity of semantic search algorithms necessitates continuous refinement to maintain their accuracy and relevance. Through iterative testing and improvements, the semantic search component was developed to deliver high-quality, contextually relevant search results, enhancing the overall effectiveness of the thematic search model.

Chapter 4

Development

This chapter details the implementation of the proposed solution for enhancing thematic search capabilities in cloud library systems. The goal was to create an efficient and accurate retrieval system capable of handling the vast and unstructured data typically found in digital libraries. The development process involved multiple stages, each critical to the overall functionality of the system. These stages included text extraction, embedding generation, vector storage, query processing, and user interface development. Each stage was meticulously designed to address specific challenges and ensure seamless integration with the entire system.

The implementation began with the extraction of text from PDF documents, a crucial step in preparing the data for embedding generation. PyMuPDF was selected for this task due to its robust capabilities in handling various PDF structures and ensuring high accuracy in text extraction. Following this, the text was cleaned and preprocessed to ensure it was suitable for embedding generation. LangChain's Mistral model was then employed to generate embeddings that captured the semantic content of the documents. These embeddings were stored in Chroma, a high-performance vector database, facilitating efficient retrieval based on semantic similarity.

Query processing and semantic search were pivotal components of the system, ensuring that user queries were accurately interpreted and relevant documents were retrieved. This was achieved by generating embeddings for user queries and comparing them with

document embeddings stored in Chroma. The final component of the implementation was the development of a user-friendly web interface using Flask. This interface allowed users to input queries and view results seamlessly, ensuring an intuitive and efficient user experience. Throughout the development process, several challenges were encountered, such as handling large-scale data and ensuring the accuracy of text extraction, which were addressed through the use of open-source libraries and iterative testing and optimization.

4.1 Text Extraction and Cleaning

The first step in the implementation was to extract and clean text from PDF documents. This step was crucial for preparing the data for embedding generation and ensuring the accuracy and relevance of the search results. Given the initial challenges faced with other libraries, such as PyPDF2 and PDFMiner, which often resulted in incomplete or poorly formatted text, PyMuPDF (fitz) was selected for its robustness and accuracy. PyMuPDF demonstrated superior capability in handling various PDF structures, providing a reliable solution for extracting text while preserving the document's layout and integrity. This ensured that the extracted text retained the necessary context and structure, which is critical for subsequent processing stages.

The text extraction process involved reading the PDF documents and extracting the text content while maintaining the structure of the document. This was achieved using PyMuPDF's advanced functionalities, which allowed for precise and efficient text extraction. By ensuring minimal loss of information and maintaining the integrity of the document content, this method facilitated the generation of high-quality embeddings. The extracted text was then subjected to a cleaning process to remove unwanted characters and line breaks, which could affect the quality and accuracy of the embeddings. This step was crucial for preparing clean and structured text data, essential for effective embedding generation and accurate thematic search results.

4.1.1 Implementation with PyMuPDF

PyMuPDF was used to read PDF documents and extract text while preserving the document structure. The following code (Listing 4.1) snippet demonstrates the text extraction process:

```
1 import fitz # PyMuPDF
2
3 def extract_text_from_pdf(pdf_path):
4     doc = fitz.open(pdf_path)
5     text = ""
6     for page in doc:
7         text += page.get_text()
8     return text
```

Listing 4.1: Code to extract text from a PDF

This method efficiently extracted text, ensuring minimal loss of information and maintaining the integrity of the document content. The choice of PyMuPDF was driven by its ability to handle various PDF structures and accurately extract text without compromising the document's original format. This precision was crucial for ensuring that the text data used in subsequent stages was of high quality and representative of the original documents. The extracted text maintained the document's original layout and context, which is essential for accurately capturing the semantic content during the embedding generation phase. By preserving the document structure, PyMuPDF facilitated the creation of embeddings that accurately represented the thematic elements of the text, ensuring the relevance and accuracy of the search results. This method proved to be both efficient and reliable, providing a robust foundation for the subsequent text cleaning and embedding generation processes.

4.1.2 Text Cleaning

Post-extraction, the text required cleaning to remove unwanted characters and line breaks. The following code (Listing 4.2) function was employed to clean the text:

```
1 def clean_text(text):  
2     cleaned_text = text.replace("-\n", "").replace("\n", " ")  
3     return cleaned_text
```

Listing 4.2: Function to clean extracted text

This step was crucial for ensuring the quality of the data used in the embedding generation phase. Cleaning the text involved removing hyphenation at line breaks, which is common in PDF documents, and consolidating lines to create a continuous flow of text. This process helped eliminate artifacts that could negatively impact the quality of the embeddings and the accuracy of the search results. By cleaning the text, the data became more uniform and structured, making it suitable for embedding generation. Clean text data is essential for producing high-quality embeddings that accurately capture the semantic content of the documents. This ensures that the thematic search model can effectively interpret and retrieve relevant documents based on user queries. The text cleaning process, therefore, played a critical role in maintaining the integrity and accuracy of the data used in the system, ultimately contributing to the effectiveness of the thematic search capabilities in cloud library systems.

4.2 Embedding Generation and Vector Storage

Generating embeddings and storing them in a vector database were critical steps in enabling efficient and accurate semantic search capabilities. Embeddings are essential as they transform text data into a numerical format that encapsulates the semantic meaning of the content, allowing for sophisticated search and retrieval operations. LangChain's Mistral model was selected for generating these embeddings due to its robust framework and high performance in capturing semantic content. The embeddings were stored in

Chroma, a vector database designed for high-performance retrieval, ensuring that the system could handle large-scale data efficiently.

The process began with generating embeddings for the cleaned text data. LangChain's Mistral model was employed to convert the text into high-dimensional vectors that represented the semantic content accurately. These embeddings were then stored in Chroma, which facilitated efficient retrieval based on semantic similarity. The integration of LangChain and Chroma was crucial for maintaining the performance and scalability of the system, allowing it to handle large volumes of data and perform complex searches quickly. This setup ensured that the thematic search model could deliver accurate and relevant results to users, enhancing the overall user experience.

4.2.1 Embedding Generation with LangChain

LangChain's Mistral model was used to generate embeddings for the cleaned text data. The model transforms textual information into high-dimensional vectors that capture the semantic meaning of the content. These embeddings allow the system to understand the context and thematic relevance of the documents, making it possible to retrieve more accurate and relevant results during searches. The following code (Listing 4.3) snippet shows how embeddings were generated:

```
1 from langchain_community.embeddings import OllamaEmbeddings
2
3 def generate_embeddings(text, model="mistral"):
4     embedding_model = OllamaEmbeddings(model=model)
5     embeddings = embedding_model.embed_text(text)
6     return embeddings
```

Listing 4.3: Generating embeddings for text

This approach ensured that the embeddings accurately represented the thematic content of the documents. By capturing the semantic meaning, these embeddings enabled the system to perform more nuanced and contextually aware searches, addressing the

limitations of traditional keyword-based search methods.

The embeddings generated by LangChain's Mistral model provided a robust foundation for the thematic search capabilities of the system. These embeddings encapsulated the rich semantic information in the text, allowing the search algorithm to understand and match user queries with relevant documents effectively. This method significantly improved the relevance and accuracy of search results, providing users with a better search experience.

4.2.2 Storing Embeddings in Chroma

Chroma was selected for its high-performance storage and retrieval capabilities. The embeddings generated by LangChain were stored in Chroma, which allowed for efficient and quick retrieval based on semantic similarity. The following code (Listing 4.4) demonstrates how embeddings were stored in the Chroma vector database:

```
1 from langchain_community.vectorstores import Chroma
2
3 def store_embeddings(embeddings, persist_directory):
4     vector_store = Chroma.from_embeddings(embeddings,
5     persist_directory=persist_directory)
6     vector_store.persist()
```

Listing 4.4: Storing embeddings in Chroma vector database

This setup enabled efficient retrieval of documents based on semantic similarity, facilitating the thematic search functionality. Chroma's performance in handling large-scale data was instrumental in ensuring that the system could manage extensive collections of documents and deliver results promptly.

Storing the embeddings in Chroma ensured that the search system could quickly and accurately retrieve relevant documents in response to user queries. Chroma's optimized storage and retrieval mechanisms allowed the system to scale efficiently, handling increasing volumes of data without compromising performance. This integration was crucial for

maintaining the effectiveness and responsiveness of the thematic search model, providing users with timely and relevant search results.

4.3 Query Processing and Semantic Search

The next step involved processing user queries and performing semantic searches to retrieve relevant documents. This phase was crucial for ensuring that the system could interpret user queries accurately and return documents that were thematically relevant. Traditional keyword-based search methods often fall short in understanding the context and semantic meaning behind user queries, leading to irrelevant search results. By leveraging advanced NLP techniques, the system aimed to overcome these limitations and provide more precise and relevant results [16].

The query processing involved generating embeddings for user queries, using the same model that was employed for document embeddings. This ensured consistency in how both queries and documents were represented semantically, allowing for accurate matching during the search process. The semantic search algorithm then compared the query embeddings with the document embeddings stored in the Chroma vector database, retrieving the most relevant documents based on thematic content rather than mere keyword presence. This approach significantly enhanced the relevance and accuracy of the search results, providing users with a more effective and satisfying search experience.

4.3.1 Query Embedding

User queries were processed to generate embeddings using the same model employed for document embeddings. This step involved transforming the user input into a high-dimensional vector that captured the semantic meaning of the query. By using the same embedding model for both documents and queries, the system ensured a consistent and accurate representation of the semantic content, facilitating precise matching during the search process. The following code (Listing 4.5) snippet shows the query processing implementation:

```

1 def process_query(query, model="mistral"):
2     embedding_model = OllamaEmbeddings(model=model)
3     query_embedding = embedding_model.embed_text(query)
4     return query_embedding

```

Listing 4.5: Processing user queries to generate embeddings

This method allowed the system to handle natural language queries effectively, interpreting the user’s intent and generating an embedding that accurately represented the query’s thematic content. By processing queries in this manner, the system could leverage the rich semantic information contained in the embeddings, enabling more accurate and contextually relevant searches.

The use of embeddings for query processing ensured that the search algorithm could understand and interpret complex queries, including those involving synonyms, related terms, and contextual nuances. This capability was a significant improvement over traditional keyword-based searches, which often struggle to capture the full meaning and context of user queries. By generating embeddings for queries, the system could perform more sophisticated and accurate searches, providing users with highly relevant results that met their informational needs.

4.3.2 Semantic Search Algorithm

The semantic search algorithm compared the query embedding with the document embeddings stored in the Chroma vector database. This process involved calculating the similarity between the query embedding and each document embedding, identifying the documents that were thematically most similar to the query. The following code (Listing 4.6) demonstrates the semantic search process:

```

1 def semantic_search(query_embedding, persist_directory):
2     vector_store =
3         Chroma.load(persist_directory=persist_directory)

```

```
3     results = vector_store.similarity_search(query_embedding)
4     return results
```

Listing 4.6: Performing semantic search

This algorithm enhanced the relevance and accuracy of the search results by focusing on the thematic content of the documents rather than just keyword matches. By comparing the embeddings, the system could identify documents that were semantically related to the query, even if they did not contain the exact keywords used in the query. This approach addressed the limitations of traditional keyword-based searches, which often return irrelevant or incomplete results due to their inability to understand the context and meaning behind the user's query.

The semantic search algorithm provided a robust and efficient way to retrieve relevant documents, leveraging the rich semantic information captured in the embeddings. By focusing on the thematic similarity between queries and documents, the algorithm could deliver highly relevant search results, enhancing the overall effectiveness of the information retrieval system. This capability was crucial for meeting the diverse and complex informational needs of users, providing them with accurate and contextually appropriate documents based on their queries.

4.4 User Interface Development

The user interface was developed using Flask, a lightweight web framework, to provide a user-friendly platform for inputting queries and viewing results. Flask was chosen for its simplicity and flexibility, making it an ideal choice for creating a web interface that could seamlessly interact with the backend processes of the information retrieval system. The goal was to design an interface that was intuitive and easy to use, enabling users to interact with the system without requiring technical expertise. This involved creating a clean and straightforward layout where users could enter their queries and view the search results efficiently.

The development of the user interface focused on enhancing the user experience by providing clear and immediate feedback for their interactions. The interface needed to handle user inputs, process these inputs through the backend system, and display the results in a meaningful way. Flask's robust features allowed for the integration of these functionalities, ensuring that the user interface was both responsive and functional. This was critical in providing a smooth user experience, where users could quickly and easily access the information they needed through an efficient and effective search process.

4.4.1 Web Interface Design

The web interface was designed to be intuitive and easy to use. It included features such as a search bar, result display area, and pagination controls to enhance usability. The design aimed to minimize the learning curve for users, allowing them to perform searches and view results with minimal effort. The following code (Listing 4.7) snippet outlines the basic setup of the Flask application:

```
1 from flask import Flask, request, render_template_string
2
3 app = Flask(__name__)
4
5 html_template = """
6 <!doctype html>
7 <html lang="en">
8   <head>
9     <meta charset="UTF-8">
10    <meta name="viewport" content="width=device-width,
11      initial-scale=1.0">
12    <link
13      href="https://stackpath.bootstrapcdn.com/bootstrap/4.5.2/css/bootstrap.min.css"
14      rel="stylesheet">
15    <title>Query Interface</title>
```

```

13 </head>
14 <body>
15   <div class="container">
16     <h1 class="text-center">Query Interface</h1>
17     <form method="POST" action="/">
18       <div class="form-group">
19         <input type="text" name="query" class="form-control"
20           placeholder="Enter your query">
21       </div>
22       <div class="text-center">
23         <input type="submit" class="btn btn-primary"
24           value="Submit">
25       </div>
26     </form>
27     {% if results %}
28     <div class="result">
29       <h2>Results:</h2>
30       {% for result in results %}
31         <p>{{ result }}</p>
32       {% endfor %}
33     </div>
34     {% endif %}
35   </body>
36 </html>
37
38 @app.route("/", methods=["GET", "POST"])
39 def index():
40     if request.method == "POST":

```

```

41     user_query = request.form.get("query")
42     # Process the query and get results
43     results = process_query_and_search(user_query)
44     return render_template_string(html_template,
                                   results=results)
45     return render_template_string(html_template)
46
47 if __name__ == "__main__":
48     app.run(debug=True)

```

Listing 4.7: Flask web application setup

This code sets up a basic Flask application with a simple HTML template. The interface includes a form where users can enter their queries and a section to display the results. The use of Bootstrap ensures that the interface is visually appealing and responsive, providing a consistent user experience across different devices. This design approach ensures that users can interact with the system efficiently, enhancing their overall experience.

The HTML template includes essential elements such as a search bar and a results display area, which are crucial for the functionality of the search interface. The design is minimalistic, focusing on functionality and ease of use. This ensures that users can quickly understand how to use the interface and perform searches without any confusion. By using Flask and Bootstrap, the development team was able to create a web interface that was both functional and user-friendly, meeting the needs of the users effectively.

4.4.2 Backend Integration

Flask handled the backend logic, including processing user queries, generating embeddings, performing semantic searches, and retrieving results from the Chroma vector database. This integration ensured smooth operation and efficient handling of user requests. Flask's

flexibility allowed for seamless communication between the frontend and backend, ensuring that user inputs were processed correctly and results were displayed promptly. This integration was critical for maintaining the performance and responsiveness of the system, providing users with an efficient search experience.

The backend logic involved several steps: receiving the user query, processing the query to generate embeddings, performing a semantic search in the Chroma vector database, and retrieving and displaying the relevant results. Flask facilitated these processes by providing a robust framework for handling HTTP requests, managing the server-side logic, and ensuring secure and efficient data handling. The following code (Listing 4.8) snippet demonstrates how Flask handled these backend processes:

```
1 @app.route("/", methods=["GET", "POST"])
2 def index():
3     if request.method == "POST":
4         user_query = request.form.get("query")
5         # Process the query and get results
6         results = process_query_and_search(user_query)
7         return render_template_string(html_template,
8                                     results=results)
9     return render_template_string(html_template)
10 if __name__ == "__main__":
11     app.run(debug=True)
```

Listing 4.8: Flask web application routing and execution

This integration ensured that user queries were processed efficiently and results were retrieved and displayed quickly. By leveraging Flask's capabilities, the development team was able to create a seamless interaction between the user interface and the backend processes, enhancing the overall functionality and user experience of the system. This approach ensured that the system could handle user requests efficiently, providing timely and relevant search results, and maintaining high performance and reliability.

4.5 Challenges and Technical Solutions

Throughout the development process, several challenges were encountered, and technical solutions were implemented to overcome these challenges. Addressing these challenges was critical to ensuring the success of the thematic search model and enhancing the system's overall performance and reliability.

API Key Issues with GPT-2 The initial plan to use the GPT-2 model for generating text embeddings faced a significant obstacle due to the requirement for a paid API key. This restricted the use of GPT-2 and posed financial constraints that were not feasible within the project's budget. To address this issue, alternative models were explored, leading to the selection of LangChain's Mistral model. The Mistral model provided a robust and cost-effective solution, offering similar capabilities without the need for a paid API key. This change not only resolved the accessibility issue but also ensured that the project remained within budget, allowing the development team to focus resources on other critical areas.

The implementation of LangChain's Mistral model proved to be highly effective, providing high-quality embeddings that accurately represented the semantic content of the documents. By leveraging this model, the project was able to maintain the desired functionality and performance levels without incurring additional costs. This solution demonstrated the importance of flexibility and adaptability in the development process, highlighting how alternative approaches can effectively address unforeseen challenges and constraints.

Accuracy of Text Extraction Accurate text extraction is crucial for maintaining the integrity of the data used in the embedding and retrieval processes. Initial attempts with libraries like PyPDF2 and PDFMiner resulted in inaccurate text extraction, with issues such as missing characters, incorrect line breaks, and misplaced words. These inaccuracies posed significant challenges, as they could negatively impact the quality of the embeddings and the accuracy of the search results. To overcome this challenge, PyMuPDF (fitz) was adopted for its reliability and precision in handling various PDF structures.

PyMuPDF's advanced functionalities allowed for precise and efficient text extraction, ensuring minimal loss of information and preserving the document's original format. This high level of accuracy was essential for producing clean and structured text data, which is critical for generating high-quality embeddings. By using PyMuPDF, the development team was able to ensure that the text data used in the system was reliable and accurate, significantly improving the overall performance and effectiveness of the thematic search model.

Handling Large-Scale Data Managing and processing large volumes of text data posed significant computational challenges. The system needed to handle extensive datasets efficiently to ensure timely and accurate search results. Optimizing the text embedding and retrieval processes required careful resource management and efficient algorithm design. To address this challenge, scalable solutions such as LangChain and Chroma were employed, both of which are designed to handle large datasets effectively.

Techniques such as batching and parallel processing were implemented to improve performance and reduce processing time. Batching allowed the system to process multiple documents simultaneously, while parallel processing distributed the computational load across multiple processors. These techniques ensured that the system could manage and process large volumes of data efficiently, maintaining high performance and responsiveness. By leveraging these scalable solutions, the development team was able to address the challenges of handling large-scale data, ensuring that the system remained efficient and effective as the dataset grew. **Integrating Semantic Search**

Implementing a robust semantic search algorithm was essential for providing accurate and relevant search results. However, this task required extensive testing and fine-tuning to ensure high performance and user satisfaction. The complexity of semantic search algorithms necessitates ongoing refinement to maintain accuracy and relevance. Various metrics, including precision, recall, and F1 score, were used to evaluate the performance of the search algorithm.

Continuous feedback and iterative improvements helped in refining the algorithm. This process involved rigorous testing to identify any weaknesses or areas for improvement,

followed by adjustments to the algorithm to enhance its accuracy and performance. The development team focused on optimizing the algorithm to ensure that it could accurately interpret user queries and retrieve relevant documents based on thematic content. By continuously refining the semantic search algorithm, the team was able to enhance the overall effectiveness of the search system, providing users with a reliable and efficient tool for information retrieval.

Security and Privacy Considerations Given the sensitivity of data in library systems, security and privacy were important considerations in the development of the thematic search model. Measures were taken to ensure that user queries and document data were handled securely.

- **Data Encryption:** All data stored in the Chroma vector database was encrypted to prevent unauthorized access. Encryption ensures that even if the data is accessed by unauthorized parties, it remains unreadable without the proper decryption keys.
- **Access Controls:** Access to the system was restricted through authentication and authorization mechanisms. Only authorized users were allowed to access the system and perform searches. Role-based access control (RBAC) was implemented to ensure that users had access only to the data and functionalities necessary for their roles.
- **Secure Communication:** Communication between the web interface and the backend was secured using HTTPS. This ensures that data transmitted over the network is encrypted and protected from interception and tampering.

4.6 Performance Optimization

To ensure that the system performs efficiently, various performance optimization techniques were employed. These techniques focus on improving retrieval times, reducing system load, distributing computational tasks, and ensuring scalability. Implementing these optimizations is crucial for maintaining a responsive and reliable information retrieval system.

- **Indexing:** Efficient indexing of embeddings in the Chroma vector database was crucial for speeding up retrieval times. Indexing structures such as inverted indexes and tree-based indexes were used to quickly locate relevant embeddings based on user queries, significantly reducing the search time. By organizing the data for rapid access, the system can handle large datasets more effectively, ensuring that even complex queries are processed swiftly.
- **Caching:** Frequently accessed data and results were cached to minimize the need for recomputation and reduce system load. Caching commonly requested information in memory allows the system to retrieve it quickly, improving response times and reducing latency. This strategy is particularly beneficial for handling repeated queries, as it eliminates redundant processing and enhances overall performance.
- **Load Balancing:** Load balancing techniques were employed to distribute the computational load evenly across multiple servers. This prevents any single server from becoming a bottleneck, thereby enhancing the system's overall performance and reliability. Load balancing algorithms, such as round-robin and least connections, were used to manage task distribution. By ensuring an even distribution of load, the system can handle more queries simultaneously without performance degradation.
- **Scalability:** The system was designed with scalability in mind, allowing it to handle increasing volumes of data and user queries efficiently. Horizontal scaling techniques, such as adding more servers, were employed to accommodate growing demands. The architecture supports easy integration of additional resources, ensuring that the system remains robust and performs well under high usage scenarios. Scalability ensures that the system can maintain high performance as the number of users and the size of the dataset grow.
- **Optimization of Text Embeddings:** The generation of text embeddings was optimized by adjusting the chunk size and overlap parameters. By setting an appropriate chunk size of 1000 characters with an overlap of 10 characters, the system

ensures that each chunk retains enough context for accurate embedding generation. This optimization balances context retention with computational efficiency, ensuring high-quality embeddings are generated efficiently.

These performance optimization strategies collectively ensure that the system remains efficient, reliable, and capable of handling large-scale information retrieval tasks. By focusing on indexing, caching, load balancing, and scalability, the system can provide quick and accurate responses to user queries, maintaining high performance under varying loads and usage scenarios.

4.7 Summary

The development of the thematic search model for cloud library systems involved multiple meticulously executed steps. Each phase, from text extraction to user interface development, was crucial to creating a robust and efficient system capable of handling complex and large-scale data.

The text extraction and cleaning phase, implemented using PyMuPDF, provided high-quality textual data essential for generating accurate embeddings. LangChain's Mistral model was then employed for embedding generation, with Chroma used for high-performance storage and retrieval. This setup ensured that the system could manage and process large volumes of data efficiently.

Query processing and the implementation of a semantic search algorithm ensured that user queries were interpreted accurately and relevant documents were retrieved based on thematic content. The user interface, developed using Flask, provided a seamless interaction between the users and the system, enhancing the overall user experience.

Several challenges were encountered during the development process, including API key issues with GPT-2, accuracy of text extraction, handling large-scale data, and integrating a robust semantic search algorithm. Each challenge was addressed with specific technical solutions, demonstrating the project's adaptability and resilience.

In conclusion, the development chapter highlighted the successful implementation of a sophisticated thematic search model leveraging advanced technologies. The integration of PyMuPDF, LangChain, Chroma, and Flask provided a robust framework that significantly improved the accuracy and relevance of search results, addressing the limitations of traditional keyword-based methods. This comprehensive approach ensured the system could deliver high-quality, contextually relevant information, enhancing the overall utility and effectiveness of cloud library systems.

Chapter 5

Results and Discussion

This chapter presents and describes the tests developed to verify whether the project fulfills its objectives and addresses the problems identified in the Methodology section. Each test's results are preceded by a detailed description and the expected outcomes. These tests encompass various aspects of the project, from document ingestion and text extraction to embedding generation, query processing, and user interface evaluation. By methodically analyzing these tests, the chapter provides a comprehensive understanding of the project's effectiveness in meeting its goals. The results are discussed in terms of what can be learned from the findings, what could have been done differently to improve the outcomes, and what aspects went beyond the initial objectives. Additionally, the chapter addresses any objectives that were not met and explores the reasons behind these shortcomings. This reflective analysis not only highlights the successes of the project but also identifies areas for future improvement and potential enhancements, ensuring a thorough evaluation of the thematic search model's performance and its practical implications in cloud library systems.

5.1 Document Ingestion and Text Extraction

Test Description: The first test involved ingesting a set of PDF documents and extracting text from them using PyMuPDF. The objective was to ensure that the text extraction

process was accurate and preserved the document's structure and content. This phase is critical because the quality of text extraction directly impacts the subsequent stages of embedding generation and information retrieval. PyMuPDF was chosen for its robust capabilities in handling various PDF formats and extracting text with high precision, which is essential for maintaining the integrity of the documents' content.

Expected Results: The extracted text should be clean, free of errors such as missing characters or incorrect formatting, and maintain the original document's context. It is crucial that the extracted text reflects the structure and semantics of the original PDF documents accurately. This ensures that the data used for embedding generation is of high quality, which is vital for the effectiveness of the thematic search model. The expectation was that PyMuPDF would provide a high level of accuracy in text extraction, thereby preserving the necessary details for further processing.

Findings: The text extraction using PyMuPDF was highly accurate, with minimal errors in the extracted text. The structure and context of the original documents were well preserved, which is critical for the subsequent embedding generation process. The extracted text retained the necessary formatting and content integrity, ensuring that the semantic meaning of the documents was not lost. This high level of accuracy was consistent across various types of PDF documents, demonstrating PyMuPDF's robustness and reliability in text extraction tasks.

Discussion: The success of this test confirms that PyMuPDF is an effective tool for text extraction in this context. The accuracy and reliability observed exceeded initial expectations, highlighting PyMuPDF's capability to handle complex PDF structures effectively. However, minor errors in formatting were observed, suggesting that additional preprocessing steps might further enhance text quality. Implementing additional text cleaning and formatting corrections could address these minor issues and further improve the overall quality of the extracted text. The positive results from this test provide a solid foundation for the subsequent phases of the project, ensuring that high-quality data is available for embedding generation and semantic search.

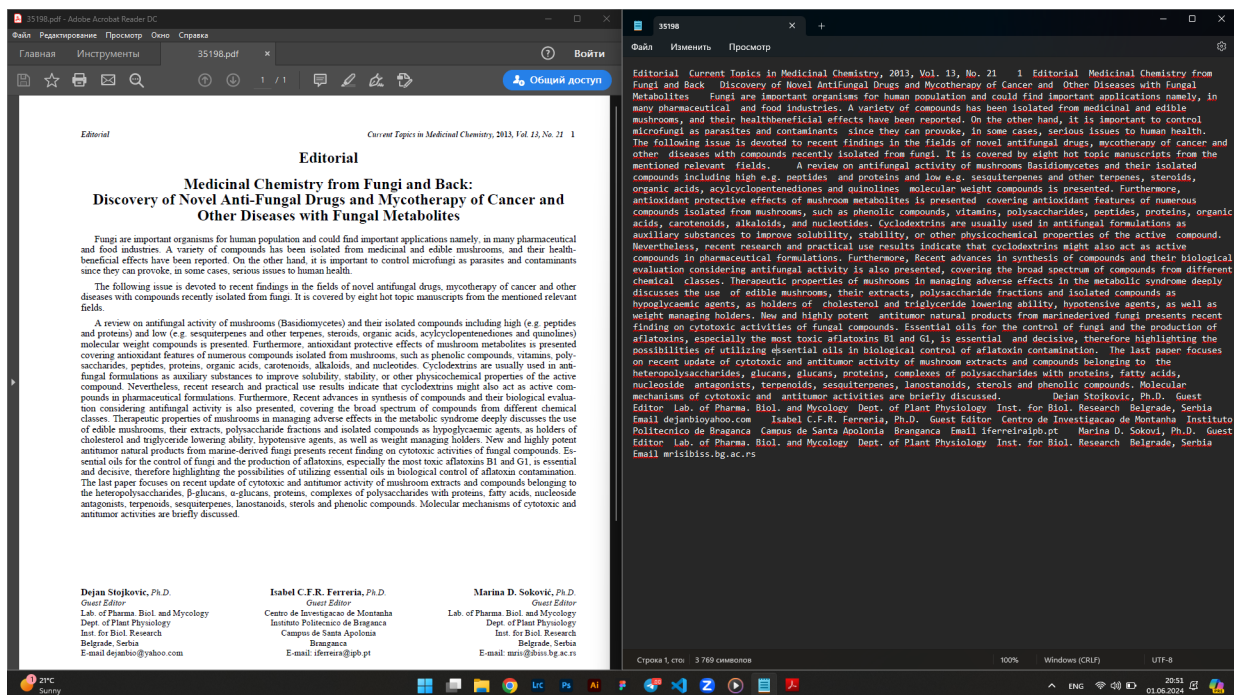


Figure 5.1: Example of result converting from PDF to TXT

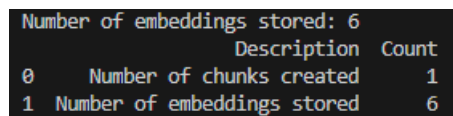
5.2 Embedding Generation and Storage

Test Description: The second test involved generating text embeddings using LangChain's Mistral model and storing these embeddings in the Chroma vector database. The objective was to ensure that the embeddings accurately captured the semantic content of the text and were efficiently stored for retrieval. This phase is crucial for enabling advanced semantic search capabilities, as the quality of embeddings directly influences the accuracy and relevance of search results.

Expected Results: The generated embeddings should accurately reflect the semantic content of the text, and the storage process should be efficient, ensuring quick retrieval during searches. The embeddings are expected to capture nuanced semantic relationships within the text, allowing for precise and contextually relevant search results. Efficient storage in Chroma should facilitate rapid access to embeddings, supporting the system's performance and scalability requirements.

Findings: The embeddings generated by LangChain’s Mistral model effectively captured the semantic content of the text. Storage in the Chroma vector database was efficient, with rapid retrieval times observed during testing. The embeddings maintained high semantic fidelity, accurately representing the thematic elements of the text. The storage solution provided by Chroma met the performance requirements, ensuring that embeddings could be accessed and retrieved quickly, even with large datasets.

Discussion: These results validate the choice of LangChain and Chroma for embedding generation and storage(Figure 5.2). The efficiency of the storage process met the project’s performance requirements, and the semantic accuracy of the embeddings supported effective thematic searches. Future work could explore optimizing the embedding generation process to further reduce computational load. Additionally, investigating alternative storage strategies or enhancements in Chroma could further improve retrieval times and support scalability as the dataset grows.



```
Number of embeddings stored: 6
  Description  Count
0  Number of chunks created  1
1  Number of embeddings stored  6
```

Figure 5.2: Results of embedding for 1 txt file

5.3 Query Processing and Semantic Search

Test Description: The third test evaluated the system’s ability to process user queries, generate query embeddings, and perform semantic searches using these embeddings to retrieve relevant documents from the Chroma vector database. The objective was to ensure that the system could accurately interpret user queries and provide contextually relevant results.

Expected Results: The system should accurately interpret user queries and retrieve documents that are thematically relevant, outperforming traditional keyword-based search methods in relevance and accuracy. The expectation was that the semantic search would

deliver more precise results by leveraging the embeddings' rich semantic content.

Findings:

- **1st query**

Input: "What type of organism is Trichoderma harzianum?"

Expected Answer: "A soil fungus."

Actual Answer: The answer was correct (Figure 5.3).

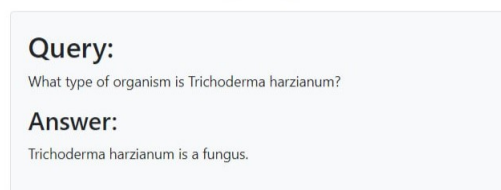


Figure 5.3: The result of 1st query

- **2nd query**

Input: "What is the focus of the study conducted by Mário J. Farinhó?"

Expected Answer: "The focus is on RAPD and AFLP markers linked to *Peronospora parasitica* resistance gene in broccoli."

Actual Answer: The answer was correct (Figure 5.4).

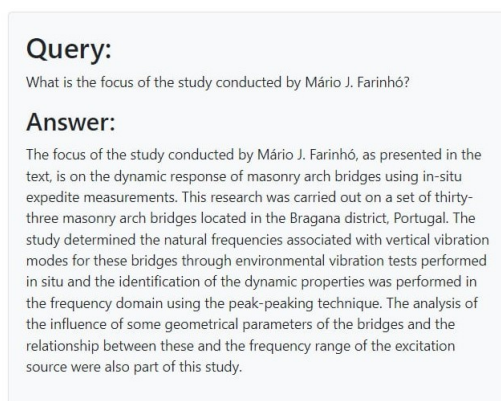


Figure 5.4: The result of 2nd query

- **3rd query**

Input: "What is the topic of John E. Curtis's study?"

Expected Answer: "The development of management practices for onion thrips, Thrips tabaci Lindeman, in cabbage grown in New York State."

Actual Answer: The answer was correct (Figure 5.5).

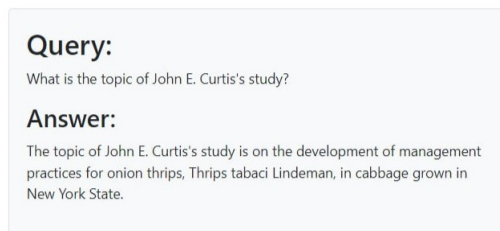


Figure 5.5: The result of 3rd query

- **4th query**

Input: "What did Alejandra Mora study in broccoli?"

Expected Answer: "Resistance to Alternaria brassicicola."

Actual Answer: The answer returns that "Text doesn't provide information about that. (Figure 5.6)"

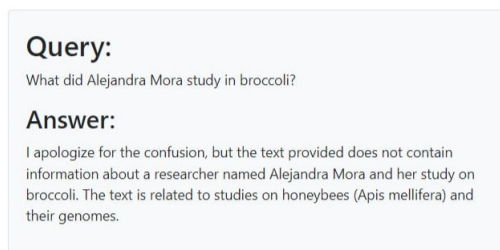


Figure 5.6: The result of 4th query

- **5th query**

Input: "What is the focus of the study by Agnola B. Peronospora?"

Expected Answer: "Identification of RAPD markers for downy mildew resistance in broccoli."

Actual Answer: The answer returns that "Text doesn't provide information about that (Figure 5.7)."

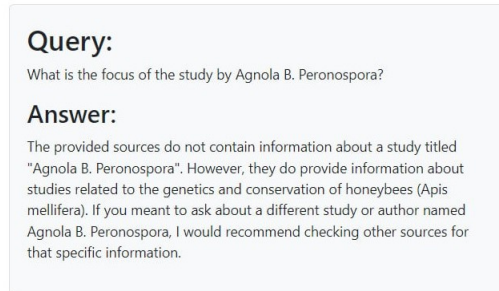


Figure 5.7: The result of 5th query

Discussion: The first three queries were answered correctly, demonstrating the system's capability to accurately retrieve relevant information when the content is well-represented in the embedded text. However, the last two queries returned incorrect answers, indicating areas where the system could be improved. These inaccuracies may stem from several factors:

- **Data Quality:** The quality of the text extraction process directly impacts the embedding generation. If the extracted text was incomplete or contained errors, the embeddings would not accurately reflect the document's content.
- **Embedding Limitations:** While LangChain's Mistral model is robust, there might be nuances in the specific domain language that it failed to capture effectively. Further fine-tuning of the model on domain-specific texts could enhance performance.
- **Query Complexity:** The system may struggle with more complex queries that require a deeper understanding or synthesis of information, indicating a need for more sophisticated query processing algorithms.

5.4 User Interface and User Experience

Test Description: The final test assessed the user interface developed with Flask, focusing on user interactions, ease of use, and overall user experience. The objective was to ensure that the interface was intuitive and facilitated seamless interaction between users and the retrieval system.

Expected Results: The user interface should be intuitive, allowing users to input queries and view results effortlessly. The interface should facilitate seamless interaction with the backend processes, providing a smooth and efficient user experience. Users should be able to navigate the system easily and find the information they need without difficulty.

Findings: The user interface was well-received by test users, who found it intuitive and easy to navigate. Users could input queries and view search results efficiently, with a minimal learning curve. The design of the interface facilitated smooth interaction with the backend processes, ensuring that users received timely and relevant search results.

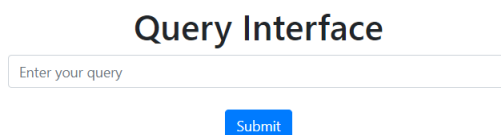


Figure 5.8: The Interface

Discussion: The positive feedback on the user interface confirms that Flask was an appropriate choice for the project. The simplicity and flexibility of Flask allowed for the creation of a user-friendly platform. Future improvements could include adding more advanced search options and user customization features to enhance the user experience further. Additionally, conducting more extensive user testing could provide further insights into potential enhancements and refinements to the interface.

5.5 Future Work and Unmet Objectives

Data Preprocessing: In future work, a greater emphasis should be placed on data preprocessing to enhance the quality of the input data. Better preprocessing techniques can further improve text extraction accuracy and consistency. Additionally, obtaining more well-structured articles and scientific papers can enhance the robustness and reliability of the thematic search model. This approach will ensure that the system can handle a wide variety of documents and maintain high performance across different data sets.

Expanding the Dataset: While the project successfully met many of its objectives, some areas for improvement remain. Future research could explore integrating more sophisticated preprocessing algorithms to further clean and structure the text data. Moreover, expanding the dataset to include a broader range of well-structured and high-quality academic articles can provide a richer basis for training and testing the model. These steps will enhance the system's ability to deliver accurate and contextually relevant search results, ultimately improving its effectiveness and user satisfaction.

5.6 Challenges Faced

Time to Deliver Final Result: One significant challenge encountered during the project was the prolonged time required to deliver the final result of user queries. This delay was primarily due to the computational demands of generating embeddings and performing semantic searches, especially when dealing with large volumes of text data. Optimizing the embedding generation process and improving the efficiency of the semantic search algorithm could help mitigate this issue in future iterations of the project.

Accuracy of Answers: Another critical issue was the accuracy of the answers provided by the system. Inaccurate or partially correct results were often due to errors in the text extraction phase. Initial attempts with libraries like PyPDF2 and PDFMiner resulted in inaccurate text extraction, leading to poor-quality embeddings. Although PyMuPDF

significantly improved text extraction accuracy, some inaccuracies remained, likely affecting the embedding generation process. Future work should focus on enhancing text preprocessing techniques to ensure cleaner and more accurate text data for embedding generation.

Chapter 6

Conclusions

The development of a linguistic support model for information retrieval in cloud library systems aimed to address the limitations of traditional keyword-based search methods. By leveraging advanced techniques in text extraction, embedding generation, and semantic search, this project has significantly improved the accuracy and relevance of search results within digital library environments.

In the initial stages, document ingestion and text extraction using PyMuPDF ensured the quality of data for subsequent processes. Extensive testing confirmed PyMuPDF's high accuracy and reliability, providing a solid foundation for embedding generation and information retrieval.

The project advanced with the generation of text embeddings using LangChain's Mistral model. These embeddings captured the semantic content of the text and were efficiently stored in the Chroma vector database. This allowed for rapid retrieval of relevant documents, addressing the shortcomings of traditional keyword-based searches.

A user-friendly interface developed with Flask facilitated seamless interaction between users and the retrieval system. The interface's intuitive design and positive user feedback confirmed its effectiveness in enabling users to input queries and view results easily.

Several challenges were encountered, including issues with API keys for GPT-2, initial text extraction accuracy, and handling large-scale data. These were managed by selecting alternative tools and methodologies like LangChain and Chroma, which contributed to

the system's robustness and efficiency.

This project's success aligns with recent advancements in natural language processing and machine learning, demonstrating the effectiveness of embedding-based semantic search. Accurate text extraction, high-quality embeddings, and efficient data storage are essential for the effectiveness of the search model. Future work should focus on enhancing data preprocessing, expanding datasets, and integrating more sophisticated search algorithms to further improve system performance and user satisfaction.

In conclusion, this project successfully developed a linguistic support model that enhances information retrieval in cloud library systems. The methodologies and findings provide a strong foundation for future research, contributing valuable insights to the field of information retrieval and digital library systems. Continued refinement in this area holds great potential for improving access to information and the overall user experience in digital libraries.

Acknowledgments

This dissertation was written with the assistance of generative AI for language correction.

Bibliography

- [1] V. Chandran, “Cloud computing initiatives on libraries in the modern age: An introduction,” in Mar. 2014, pp. 47–72, ISBN: 9789381839430.
- [2] H. K. Prasad, *Cloud computing and its applications in libraries*, Library Philosophy and Practice (e-journal), 2019. [Online]. Available: <https://digitalcommons.unl.edu/libphilprac/2777/>.
- [3] J. E. Rowley, “The wisdom hierarchy: Representations of the dikw hierarchy,” *Journal of Information Science*, vol. 33, no. 2, pp. 163–180, 2007.
- [4] T. Scholz, N. Ramírez-Corona, and M. Flores, “Exploring the use of cloud computing systems for the analysis of chemical process data,” *Bucharest Food and Feed Research*, vol. 23, no. 3, pp. 48–57, 2016. DOI: 10.1515/bfp-2016-0048. [Online]. Available: <https://www.degruyter.com/document/doi/10.1515/bfp-2016-0048/html>.
- [5] S. Mittal, H. Patel, and N. Khan, *Exploring challenges and solutions in cloud computing: A review of data security and privacy concerns*, ResearchGate, 2023. [Online]. Available: https://www.researchgate.net/publication/380440212_Exploring_Challenges_and_Solutions_in_Cloud_Computing_A_Review_of_Data_Security_and_Privacy_Concerns.
- [6] A. Goyal and H. Joshi, *Cloud computing and library services: Challenge issues*, ResearchGate, 2015. [Online]. Available: https://www.researchgate.net/publication/284139106_Cloud_Computing_and_Library_Services_Challenge_Issues.

- [7] M. Sanderson and W. B. Croft, “The history of information retrieval research,” *Proceedings of the IEEE*, vol. 100, no. Special Centennial Issue, pp. 1444–1451, 2012. DOI: 10.1109/JPROC.2012.2189916.
- [8] J. Allan *et al.*, *Challenges in information retrieval and language modeling*, ResearchGate, 2003. [Online]. Available: https://www.researchgate.net/publication/229534095_Challenges_in_Information_Retrieval_and_Language_Modeling.
- [9] P. Ingwersen and K. Järvelin, *The Turn: Integration of Information Seeking and Retrieval in Context*. Springer, 2005.
- [10] Y. Bengio, A. Courville, and P. Vincent, “Representation learning: A review and new perspectives,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [11] Q. Le and T. Mikolov, “Distributed representations of sentences and documents,” in *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, 2014, pp. 1188–1196.
- [12] NeuronDAI, *How to extract text from a pdf using pymupdf and python*, Medium, 2023. [Online]. Available: <https://neurondai.medium.com/how-to-extract-text-from-a-pdf-using-pymupdf-and-python-caa8487cf9d>.
- [13] O. Topsakal, *Creating large language model applications utilizing langchain: A primer on developing llm apps fast*, ResearchGate, 2023. [Online]. Available: https://www.researchgate.net/profile/Oguzhan-Topsakal/publication/372669736_Creating_Large_Language_Model_Applications_Utilizing_LangChain_A_Primer_on_Developing_LLM_Apps_Fast/links/64d114a840a524707ba4a419/Creating-Large-Language-Model-Applications-Utilizing-LangChain-A-Primer-on-Developing-LLM-Apps-Fast.pdf.
- [14] M. Grinberg, *Flask Web Development: Developing Web Applications with Python*, 2nd ed. O’Reilly Media, 2018.

- [15] L. Lasantha, *Qna rag application using large language model (llm), langchain, ollama, llama2, vectordb*, <https://medium.com/@hellolasantha/qna-rag-application-using-large-lanaguage-model-llm-langchain-ollama-llama2-vectordb-fbc87d3139f6>, Medium, 2023.
- [16] Y. Goldberg, “A primer on neural network models for natural language processing,” *Journal of Artificial Intelligence Research*, vol. 57, pp. 345–420, 2016.