

GenQuery

A REPORT

SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE COLLOQUIUM (BCCS-4105)

OF

BACHELOR OF TECHNOLOGY

IN

COMPUTER SCIENCE & ENGINEERING

Submitted by:

Harsh Jha (BCS2021-078)

Under the supervision of

Dr. Debanjan Sadhya



DEPT. OF COMPUTER SCIENCE & ENGINEERING

IIITM GWALIOR

(ABV-IIITM Gwalior)

Morena Link Road, Gwalior-474015

JULY 2024

DEPT. OF COMPUTER SCIENCE & ENGINEERING

A

(ABV-IIITM Gwalior)

Morena Link Road, Gwalior-474015

CANDIDATE DECLARATION

I hereby certify that the work, which is being presented in the report, entitled **Gen-Query Empowering Insightful Analysis and Discovery through Advanced LLM-Driven News Research** , in partial fulfillment of the requirement for summer project (BCCS-2999) for Bachelor of Technology in Computer Science and Engineering and submitted to the institution is an authentic record of our own work carried out during the period May 2024 to July 2024 under the supervision of Dr. Debanjan Sadhya.

(Harsh Jha)

Place : Gwalior

Date : 19/08/2024

Dr. Debanjan Sadhya

(SUPERVISOR)

Department of Computer Science

ABV-IIITM Gwalior

ABSTRACT

Keywords - LLM's, LangChain , OpenAI Embeddings, Vector DB

GenQuery is designed to streamline the process of analyzing and reacquiring information from news articles and websites. By allowing Users to input URLs or textbook lines, the tool fetches and processes composition content using LangChain's UnstructuredURLLoader. Embedding vectors are generated through OpenAI's embeddings and stored in the FAISS vector database to enable swift and effective similarity searches. The tool addresses the issue of high API costs due to token operation by employing Map Reduce Document Chaining, which breaks down large documents into manageable Chunks, reducing token operation per API call while maintaining essential information. Users can interact with Large Language Models(LLMs) to input queries and admit detailed answers with source references, easing translucency and verification. Users can interact with Large Language Models (LLMs), such as ChatGPT, by submitting queries and receiving comprehensive answers along with source URLs, ensuring transparency and verification. The GenQuery exemplifies a robust and cost-effective methodology for managing and querying extensive news datasets.

ACKNOWLEDGEMENT

I am highly indebted to Dr. Debanjan Sadhya and am obliged for giving me the autonomy of functioning and experimenting with ideas. I would like to take this opportunity to express my profound gratitude to them not only for their academic guidance but also for their personal interest in our project and constant support coupled with confidence-boosting and motivating sessions, which proved very fruitful and were instrumental in infusing self-assurance and trust within me. The nurturing and blossoming of the present work is mainly due to their valuable guidance, suggestions, astute judgment, constructive criticism, and an eye for perfection. My mentor always answered myriad of my doubts with smiling graciousness and prodigious patience, never letting me feel that I am a novice by always lending an ear to my views, appreciating and improving them, and by giving me a free hand in my project. It's only because of their overwhelming interest and helpful attitude that the present work has attained the stage it has.

Finally, I am grateful to my institution and colleagues whose constant encouragement served to renew my spirit, refocus my attention and energy, and helped me in carrying out this work.

Contents

CANDIDATE DECLARATION	2
ABSTRACT	3
ACKNOWLEDGEMENT	4
LIST OF FIGURES	6
CHAPTER 1: INTRODUCTION	1
1.1 Overview	1
1.2 Problem Statement	1
1.3 Objectives	2
1.4 Key Features	2
1.5 Technological Integration	3
1.5.1 Large Language Models (LLMs)	3
1.5.2 Transformer Model	3
1.5.3 FAISS Vector Database	4
1.5.4 LangChain	4
1.6 Literature Review	5
CHAPTER 2: SYSTEM DESIGN	6
2.1 Overall System Architecture	6
2.2 Map-Reduce Chaining for Query Handling	7
2.3 Stuff Method for Efficient Query Management	7
2.4 RAG (Retrieval-Augmented Generation) Working Mechanism	8
2.5 Vector Databases and Semantic Search	9
2.5.1 Embeddings: Transforming Data into Meaningful Vectors	9
2.5.2 Metrics for Evaluation and Fine-Tuning	10

2.6	Interface Design and Integration	11
2.6.1	Initial Interface View	13
2.6.2	Processed Interface View	14
CHAPTER 3: RESULTS		15
3.1	Dataset Loading and Input Preparation	15
3.1.1	Dataset Inputs and Filtering	15
3.1.2	T5 Model Summarization and Accuracy Evaluation	17
3.2	Model Performance	18
3.2.1	Summary of the ROUGE Metrics	19
3.3	Model Analysis	19
CHAPTER 4: CONCLUSION		22
Bibliography		

List of Figures

Figure 2.1 : High-Level System Architecture	7
Figure 2.2 : Map-Reduce Chaining Process	7
Figure 2.3 : Stuff Method for Query Management	8
Figure 2.4 : Illustration of Vector Database and Semantic Search	9
Figure 2.5 : Streamlit UI with Empty Prompt	13
Figure 2.6 : Streamlit UI with URL and Prompt Provided	14
Figure 3.1 : Code for Dataset Loading and Input Preparation	15
Figure 3.2 : Filtered dataset inputs from CNN articles	16
Figure 3.3 : 5 Model Summarization	17
Figure 3.4 : Accuracy evaluation shows limitations for summarization tasks	18
Figure 3.5 : T5-small: A 60M parameter transformer excelling in ROUGE metrics for summarization	19
Figure 3.6 : FLAN-T5 Base: A robust transformer with strong ROUGE performance, tailored for high-quality summarization	20
Figure 3.7 : GPT-3: A versatile model with high ROUGE-1 scores, excelling in content generation but less focused on summarization coherence	21

Chapter 1

INTRODUCTION

1.1 Overview

The report on the GenQuery begins with an introduction outlining its purpose and objectives. It addresses challenges in efficient article retrieval and query handling, emphasizing the tool's key features such as URL/text content fetching, embedding construction using OpenAI's models [1], and interaction with Large Language Models (LLMs) for query responses linked with source URLs. The technological integration involves leveraging LangChain [2] for data processing and FAISS for efficient similarity search, ensuring robust backend support. The design and implementation section detail the system architecture, data processing workflows, and user interaction mechanisms. It highlights the seamless integration of user interfaces with backend functionalities, including the handling of queries through LLMs and transparent result presentation. This report aims to provide a comprehensive view of the tool's development, implementation, and operational capabilities in enhancing research efficiency and information retrieval processes.

1.2 Problem Statement

GenQuery aims to address the challenge of efficient and cost-effective information retrieval from large documents or articles. The primary issue lies in the substantial API costs incurred due to extensive token usage when processing lengthy textual content. This cost not only poses financial constraints but also limits the scalability of the tool for widespread adoption. The goal is to minimize the number of tokens processed per

API call without compromising on the integrity and relevance of the retrieved information. By implementing strategies such as Map Reduce Document Chaining, the tool seeks to break down large documents into manageable segments, process each segment individually to generate embeddings efficiently, and then aggregate these embeddings. This approach not only reduces token consumption per API transaction but also enhances the overall efficiency of information retrieval, making the **GenQuery** a cost-effective solution for research and academic endeavors.

1.3 Objectives

Develop an AI-Driven Platform for Streamlining Research Queries

The "GenQuery" aims to develop a content processing platform that uses LangChain's UnstructuredURLLoader, GPT-3.5 Turbo, and the FAISS [3] vector database to enhance academic research. The tool will enable users to load URLs or text files, efficiently process large text data, and generate accurate embedding vectors for similarity searches. By employing Map Reduce Document Chaining [4], the project aims to minimize API costs while preserving the integrity of the information, ensuring relevant and concise results.

The **GenQuery** aims to streamline academic research by integrating LangChain for content extraction, OpenAI's embeddings for semantic accuracy, and FAISS for rapid similarity searches. The tool efficiently processes large volumes of text, generating precise embeddings to ensure relevant query responses. By leveraging transformer models, it enhances user interaction with contextually accurate answers, optimizing the research process and facilitating swift information retrieval.

1.4 Key Features

GenQuery efficiently retrieves and analyzes article content using LangChain's UnstructuredURLLoader for URL content extraction and RecursiveCharacterTextSplitter for text segmentation. It integrates with OpenAI's embeddings and FAISS vector database for fast information retrieval, enabling seamless interaction with LLMs for querying and transparent result verification through source URLs. This tool optimizes research workflows by

enhancing information exploration and analysis capabilities.

1.5 Technological Integration

The GenQuery harnesses advanced Transformer models [5], including GPT-3.5 Turbo from OpenAI, to power its natural language processing capabilities. These models enable the tool to perform sophisticated tasks such as text generation, question answering, and semantic understanding. Each of these technologies plays a crucial role in ensuring the tool's efficiency, accuracy, and scalability. Below is a detailed explanation of these technologies and their integration into the system:

1.5.1 Large Language Models (LLMs)

Large Language Models (LLMs), such as GPT-3.5 Turbo, are foundational to the GenQuery's capabilities. LLMs are designed to understand and generate human-like text based on vast amounts of data they have been trained on. In the context of this tool, LLMs are used for:

- **Natural Language Processing (NLP):** LLMs are capable of understanding user queries and generating accurate, context-aware responses.
- **Query Handling:** LLMs interact with user queries in natural language, linking responses to relevant sources and ensuring that information is presented coherently.
- **Text Summarization:** LLMs help in summarizing large volumes of text into concise summaries, reducing the amount of data that needs to be processed during searches.

1.5.2 Transformer Model

The Transformer model, upon which LLMs are built, is a neural network architecture that has revolutionized natural language processing. Key features include:

- **Self-Attention Mechanism:** This mechanism allows the model to weigh the importance of different words in a sentence relative to one another, leading to better understanding of context and meaning.

- **Parallel Processing:** Transformers can process multiple words simultaneously, making them more efficient than previous sequential models like RNNs.
- **Scalability:** The Transformer model's architecture allows it to be scaled to handle vast amounts of data and more complex tasks, contributing to the power of LLMs.

1.5.3 FAISS Vector Database

FAISS (Facebook AI Similarity Search) is an open-source library that is used to efficiently search for similar vectors, which is crucial for high-dimensional data like embeddings generated by LLMs. It is used in this tool for:

- **Efficient Similarity Search:** FAISS is optimized for searching through large datasets to find the most similar items quickly.
- **High-Dimensional Data Handling:** FAISS is capable of handling the high-dimensional vectors generated by LLMs, making it ideal for similarity searches in this context.
- **Indexing Mechanisms:** It employs various indexing methods, such as Inverted File Index (IVF) and HNSW (Hierarchical Navigable Small World graphs), to improve search speed and accuracy.

1.5.4 LangChain

LangChain is a framework that facilitates the integration of LLMs into various applications, enhancing their functionality. In this tool, LangChain is responsible for:

- **Data Processing:** LangChain helps in processing and organizing large amounts of data before feeding it into LLMs.
- **Task Orchestration:** It manages the flow of tasks, ensuring that the data is correctly pre-processed, fed into the model, and that the results are post-processed for output.
- **Extensibility:** LangChain's modular architecture allows for easy integration of new models, data sources, and processing techniques, making the system adaptable to new research needs.

1.6 Literature Review

The literature surrounding the development of advanced information retrieval systems has evolved significantly, with key milestones in natural language processing (NLP) and data management technologies playing pivotal roles. The introduction of Transformer models, such as those used in large language models (LLMs) like GPT-3.5 Turbo, has dramatically enhanced the capability of machines to understand and generate human-like text, leading to more accurate and context-aware information retrieval. These models, based on the Transformer architecture introduced by Vaswani et al. (2017), leverage self-attention mechanisms to process text in a way that captures long-range dependencies, making them exceptionally powerful for tasks such as document summarization, question answering, and semantic search. This evolution is complemented by the emergence of vector-based representations of text, where words and documents are encoded as dense vectors within high-dimensional spaces. Technologies like FAISS (Facebook AI Similarity Search), which facilitate efficient similarity searches within these vector spaces, have become indispensable tools for handling large-scale data retrieval tasks. By enabling rapid indexing and searching of document embeddings, FAISS enhances the ability of systems to retrieve relevant information from vast datasets, a capability that is crucial for the efficient functioning of modern research tools.

However, while these technological advancements have significantly improved the field, there are still notable gaps and challenges, particularly in integrating these sophisticated models into user-friendly, cost-effective research tools. Existing systems often struggle with balancing computational efficiency and accuracy, especially when dealing with large volumes of text data. The GenQuery seeks to address these challenges by combining the strengths of LLMs, FAISS, and frameworks like LangChain, which streamlines the processing of textual data through modular pipelines. The tool's innovative use of Map Reduce Document Chaining and embedding-based retrieval ensures that complex research queries can be handled efficiently, without the prohibitive costs typically associated with high API usage. By bridging these gaps, the GenQuery not only builds upon the foundation laid by previous research but also pushes the boundaries of what is possible in automated information retrieval, making it a valuable asset for researchers in various domains.

Chapter 2

SYSTEM DESIGN

2.1 Overall System Architecture

The system architecture is designed for efficient retrieval, processing, and querying of large web-sourced documents. The key components include:

- **Document Retrieval:** Using LangChain's `UnstructuredURLLoader`, the system fetches documents directly from the web, preparing the content for processing.
- **Text Segmentation:** Retrieved documents are split into smaller chunks or segments to facilitate efficient processing and embedding.
- **Embedding Construction:** Each text segment is converted into vector embeddings using OpenAI's GPT-3.5 Turbo, ensuring the data is prepared for similarity searches.
- **FAISS Vector Database:** The embeddings are stored in the FAISS database. When a user queries the system, FAISS efficiently searches for the most relevant embeddings.
- **Query Handling:** The system processes user queries by accessing the vector embeddings stored in FAISS, returning accurate and relevant information.

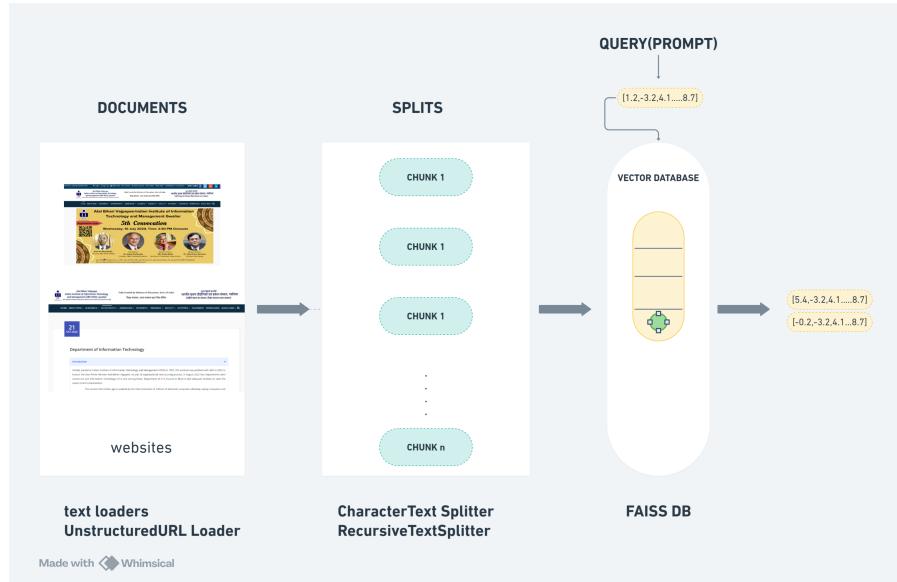


Figure 2.1: High-Level System Architecture

2.2 Map-Reduce Chaining for Query Handling

To handle large documents efficiently, the system uses a Map-Reduce Chaining approach. This method breaks down documents into smaller segments, processes each segment individually, and aggregates the results to form a comprehensive response. The following diagram explains the Map-Reduce Chaining process:

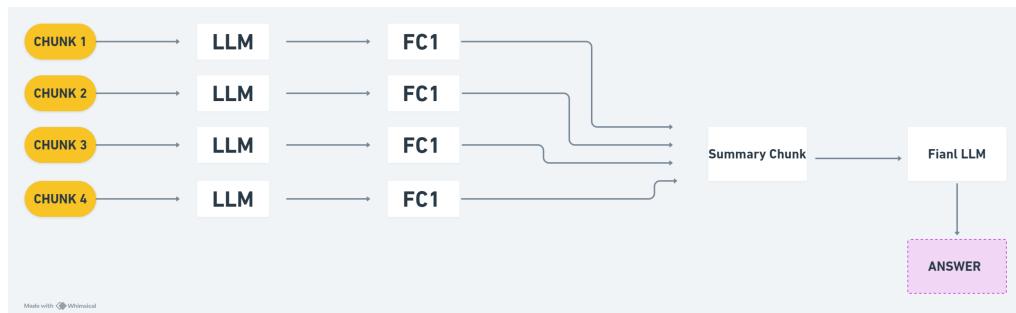


Figure 2.2: Map-Reduce Chaining Process

2.3 Stuff Method for Efficient Query Management

The Stuff Method is employed to manage queries more efficiently by reducing API costs and optimizing the way data is processed. The diagram below illustrates how the Stuff Method is integrated into the system:

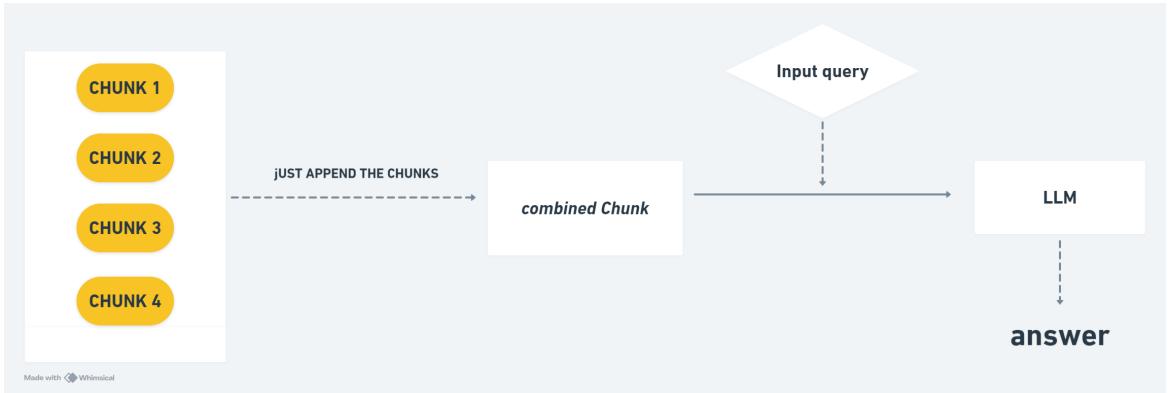


Figure 2.3: Stuff Method for Query Management

2.4 RAG (Retrieval-Augmented Generation) Working Mechanism

The system also incorporates a RAG (Retrieval-Augmented Generation) mechanism to improve the relevance of the responses generated by the LLM. Expanding the context window involves providing more contextual information to LLMs, theoretically allowing them to make more informed responses. Anthropic, for instance, introduced the Claude model with an impressive 100K token context window. OpenAI followed suit, unveiling a 32K token GPT-4 model and a 16K token GPT-3.5 model. While the idea of an extensive context window may seem like a panacea, it's important to acknowledge that the approach of "context stuffing" has its drawbacks. The following figure demonstrates how RAG works within the system:

The **GenQuery** exemplifies a successful application of advanced technologies to streamline academic research processes. By integrating LangChain's UnstructuredURL Loader, OpenAI's embeddings, and FAISS for similarity search, the tool offers a robust solution for managing and retrieving extensive textual data. The architecture leverages transformer models, particularly for generating high-quality embeddings that capture the semantic nuances of the content, thereby facilitating accurate and relevant information retrieval.

The successful implementation of this project highlights several key technological and architectural achievements:

2.5 Vector Databases and Semantic Search

In our earlier examination of Retrieval Augmented Generation (RAG), we highlighted its role in enriching LLMs with current and domain-specific information. Central to this enhancement is the use of vector databases, which are essential for effective data retrieval and semantic search.

Vector databases excel in handling multi-dimensional data points, or vectors, representing characteristics or features of data. Unlike traditional databases, which struggle with high-dimensional data, vector databases are optimized for storing and querying these vectors based on similarity rather than exact matches. This capability allows for more precise and contextually relevant search results.

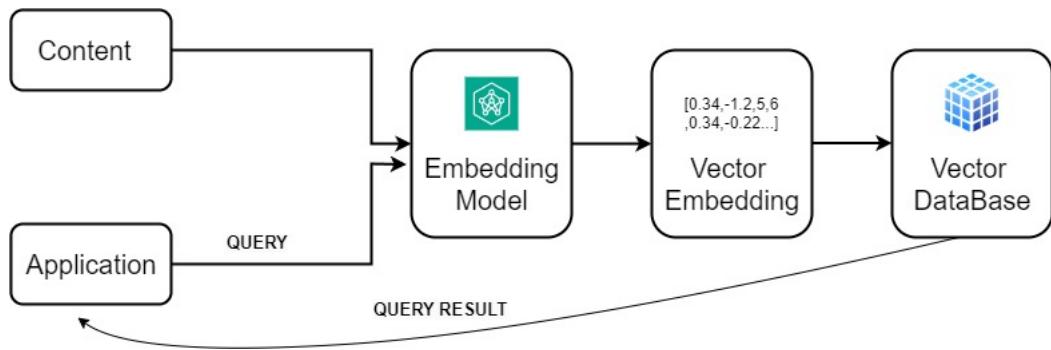


Figure 2.4: Illustration of Vector Database and Semantic Search

2.5.1 Embeddings: Transforming Data into Meaningful Vectors

Embeddings are crucial for converting unstructured data, such as text, images, and audio, into numerical representations. These embeddings allow machine learning models to process and understand data more effectively. By representing data points as vectors, embeddings capture the essence and meaning of the data, making it possible for algorithms to discern patterns and relationships.

Creating a vector database involves embedding data and incorporating metadata to enhance functionality. Metadata, such as source documents and custom tags, enriches the database, enabling more refined searches and improving the overall utility.

2.5.2 Metrics for Evaluation and Fine-Tuning

To ensure the effectiveness of our system, we utilize various metrics for evaluation and fine-tuning. Techniques such as Parameter-Efficient Fine-Tuning (PEFT) enhance model performance without extensive retraining. Metrics like ROUGE (for Summarization learning) and few-shot learning contexts help in assessing the quality and relevance of generated responses. These evaluations guide the optimization of our tool, ensuring accurate and efficient responses to user queries.

This structured approach to system design leverages the power of vector databases and embeddings to automate and streamline information retrieval, making the research process more efficient and user-friendly.

- **Efficient Data Retrieval:** The tool automates the extraction of content from URLs and text files, ensuring that users can quickly access large volumes of information without manual intervention. LangChain’s UnstructuredURL Loader plays a crucial role in processing raw web content into a structured format suitable for further analysis.
- **Accurate Embedding Generation:** Using OpenAI’s embeddings, the tool creates precise vector representations of the content. These embeddings capture the semantic essence of the text, which is crucial for effective similarity searches and query responses.
- **Swift Information Retrieval:** FAISS, a powerful similarity search library, is employed to enable fast and efficient retrieval of relevant information from the embedding vectors. This ensures that users receive prompt and accurate answers to their queries, enhancing the overall efficiency of the research process.
- **Seamless User Interaction:** The tool is designed to be user-friendly, providing an intuitive interface for querying and interacting with the embedded data. The integration of transformer models facilitates contextually relevant responses, making it easier for users to obtain the information they need without extensive manual effort.

The successful integration of these technologies demonstrates how modern AI solutions can transform the research experience. By automating content retrieval and embedding generation, the GenQuery not only simplifies the research process but also ensures

that users can efficiently access and analyze relevant information. This approach significantly reduces manual effort and enhances the accuracy and relevance of the information retrieved, showcasing the potential of AI-driven solutions in academic research.

Furthermore, the ****Retrieval-Augmented Generation (RAG)**** framework effectively addresses some limitations of Large Language Models (LLMs) by incorporating up-to-date information, domain-specific data, and organizational knowledge. Here's how RAG enhances LLM performance:

- **Retrieval Component:** RAG includes a retrieval mechanism that fetches context-specific data from external databases or documents. This ensures that the data retrieved is highly relevant to the query being processed.
- **Generation Component:** The retrieved information is combined with the original query, creating an enriched context for the LLM. This results in more accurate and relevant responses by improving the model's understanding of the query.

The result is that RAG allows LLMs to cite their sources, improving auditability and significantly enhancing the accuracy and relevance of their responses. This integration of retrieval and generation mechanisms highlights the potential for advanced AI techniques to improve information retrieval and user interaction in academic and other domains.

2.6 Interface Design and Integration

The development of the GenQuery prioritizes creating a user-friendly and intuitive interface that simplifies the research process for users. To achieve this, we will leverage Streamlit **[6]**, a powerful and efficient framework for building interactive web applications. Streamlit allows us to rapidly develop and deploy a dynamic interface, where users can easily input URLs or text files, process them, and retrieve relevant information. Streamlit's ability to integrate seamlessly with Python code ensures that the user interface is not only responsive but also tightly coupled with the backend functionalities of the tool. This integration enables the tool to interact smoothly with components like LangChain for data processing, GPT-3.5 Turbo for generating embeddings, and FAISS for efficient similarity searches. Users can perform complex queries and receive results in

real-time, enhancing their research experience.

Our objective is to ensure that the interface is both intuitive and functional, allowing users to focus on their research rather than the technical complexities of the tool. The simplicity of Streamlit's design capabilities, combined with its flexibility, ensures that the tool remains accessible to a wide range of users, from seasoned researchers to those less familiar with advanced computational tools.

In the future, we plan to integrate the GenQuery with existing academic and research platforms, allowing for real-world testing and validation. This step will ensure that the tool meets the practical needs of researchers in various fields. Additionally, we will explore the potential for commercializing the tool as a SaaS-based application, making it available to a broader academic audience and enhancing research capabilities across institutions.

2.6.1 Initial Interface View

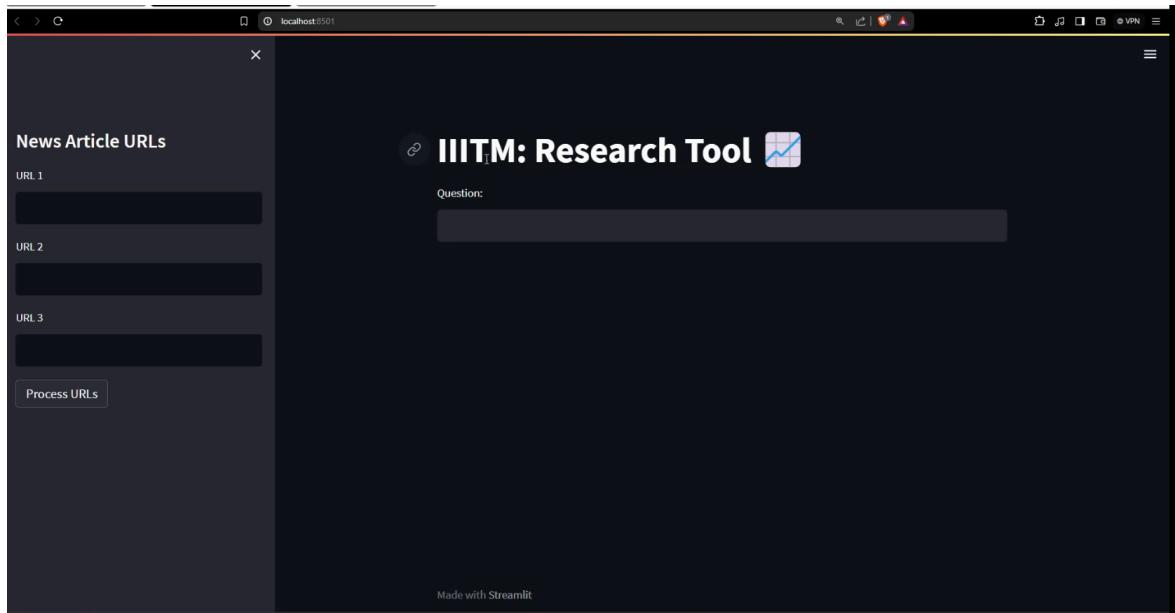


Figure 2.5: Streamlit UI with Empty Prompt

The screenshot above shows the Streamlit UI when no URL or prompt has been entered. This initial view provides a clean interface where users can input their queries and URLs.

2.6.2 Processed Interface View

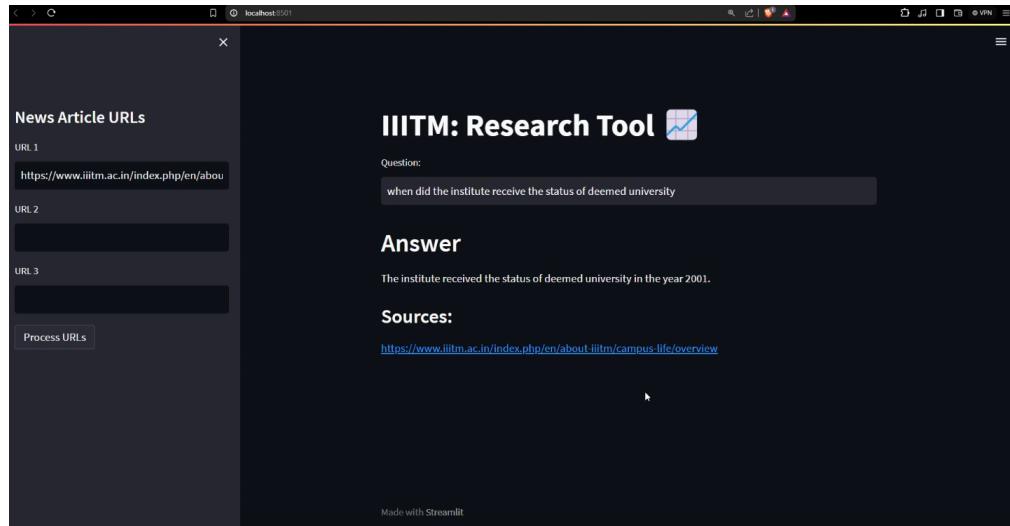


Figure 2.6: Streamlit UI with URL and Prompt Provided

Chapter 3

RESULTS

3.1 Dataset Loading and Input Preparation

For this evaluation, we utilized a subset of the CNN / DailyMail Dataset is an English-language dataset containing just over 300k unique news articles as written by journalists at CNN and the Daily Mail. This dataset provides a comprehensive collection of news articles paired with summaries in the "highlights" column, making it an ideal resource for assessing the performance of summarization models. The dataset includes a diverse array of topics and writing styles, which helps in evaluating how well the models perform across different contexts. To ensure the models were evaluated under consistent conditions, the dataset was preprocessed and loaded using the following code. The image below shows the code used and the corresponding input data:

```
import torch
from datasets import load_dataset

full_dataset = load_dataset(
    "cnn_dailymail",
    version="3.0.0",
    cache_dir=DA.paths.datasets) # Note: We specify cache_dir to use pre-cached data.

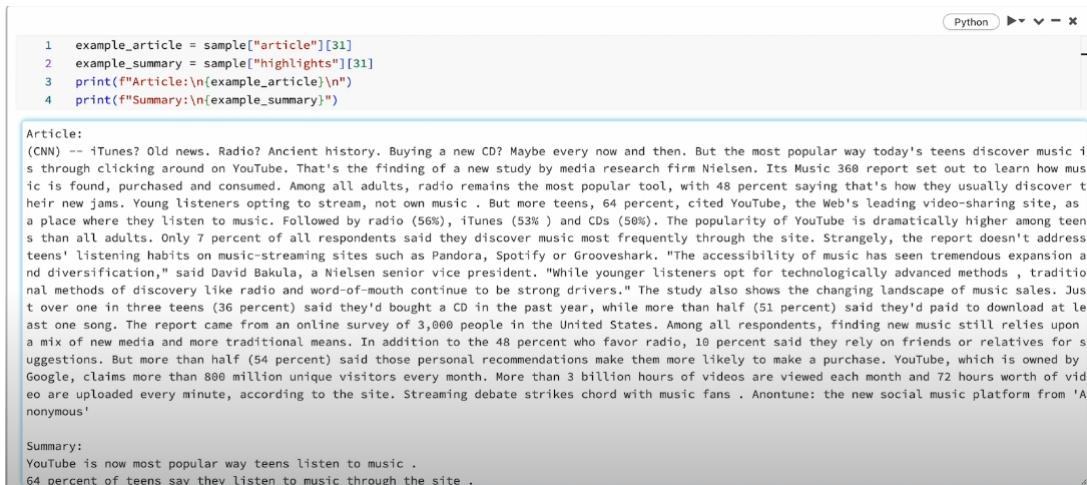
# Use a small sample of the data during this lab, for speed.
sample_size = 100
sample = full_dataset["train"].filter(lambda r: 'CNN' in r["article"][:25]).shuffle(seed=42).select(range(sample_size))
sample
```

Figure 3.1: Code for Dataset Loading and Input Preparation

3.1.1 Dataset Inputs and Filtering

In this subsection, we detail the process of importing and filtering the dataset used for model evaluation. The dataset, comprising CNN articles, was initially imported and pro-

cessed to isolate relevant content. We applied a lambda function to filter articles based on predefined criteria, ensuring that only high-quality and pertinent data were included. The following example demonstrates the inputs after filtering, showcasing how the dataset was refined for accurate model performance assessment.



```
1 example_article = sample["article"][31]
2 example_summary = sample["highlights"][31]
3 print(f"Article:{\n{example_article}\n}")
4 print(f"Summary:{\n{example_summary}\n}")

Article:
(CNN) -- iTunes? Old news. Radio? Ancient history. Buying a new CD? Maybe every now and then. But the most popular way today's teens discover music is through clicking around on YouTube. That's the finding of a new study by media research firm Nielsen. Its Music 360 report set out to learn how music is found, purchased and consumed. Among all adults, radio remains the most popular tool, with 48 percent saying that's how they usually discover their new jams. Young listeners opting to stream, not own music. But more teens, 64 percent, cited YouTube, the Web's leading video-sharing site, as a place where they listen to music. Followed by radio (56%), iTunes (53%) and CDs (50%). The popularity of YouTube is dramatically higher among teens than all adults. Only 7 percent of all respondents said they discover music most frequently through the site. Strangely, the report doesn't address teens' listening habits on music-streaming sites such as Pandora, Spotify or Grooveshark. "The accessibility of music has seen tremendous expansion and diversification," said David Bakula, a Nielsen senior vice president. "While younger listeners opt for technologically advanced methods, traditional methods of discovery like radio and word-of-mouth continue to be strong drivers." The study also shows the changing landscape of music sales. Just over one in three teens (36 percent) said they'd bought a CD in the past year, while more than half (51 percent) said they'd paid to download at least one song. The report came from an online survey of 3,000 people in the United States. Among all respondents, finding new music still relies upon a mix of new media and more traditional means. In addition to the 48 percent who favor radio, 10 percent said they rely on friends or relatives for suggestions. But more than half (54 percent) said those personal recommendations make them more likely to make a purchase. YouTube, which is owned by Google, claims more than 800 million unique visitors every month. More than 3 billion hours of videos are viewed each month and 72 hours worth of video are uploaded every minute, according to the site. Streaming debate strikes chord with music fans. Anontune: the new social music platform from 'Anonymous'

Summary:
YouTube is now most popular way teens listen to music .
64 percent of teens say they listen to music through the site .
```

Figure 3.2: Filtered dataset inputs from CNN articles

3.1.2 T5 Model Summarization and Accuracy Evaluation

For summarization, we employed the T5 model, which performs similarly to a pipeline for text summarization but involves manual tokenization. The process was implemented as follows

```
def summarize_with_t5(model_checkpoint: str, articles: list, batch_size: int=8) -> list:
    """
        Compute summaries using a T5 model.
        This is similar to a 'pipeline' for a T5 model but does tokenization manually.

        :param model_checkpoint: Name for a model checkpoint in Hugging Face, such as "t5-small" or "t5-base"
        :param articles: List of strings, where each string represents one article.
        :return: List of strings, where each string represents one article's generated summary
    """
    if torch.cuda.is_available():
        device = "cuda:0"
    else:
        device = "cpu"

    model = T5ForConditionalGeneration.from_pretrained(
        model_checkpoint,
        cache_dir=DA.paths.datasets
    ).to(device)
    tokenizer = AutoTokenizer.from_pretrained(
        model_checkpoint,
        model_max_length=1024,
        cache_dir=DA.paths.datasets)

    def perform_inference(batch: list) -> list:
        inputs = tokenizer(
            batch,
            max_length=1024,
            return_tensors="pt",
```

Figure 3.3: 5 Model Summarization

We initially evaluated the T5 model using a basic 0/1 accuracy metric, but the result was zero. This highlights that accuracy, as a generic metric, is not suitable for summarization tasks. Small variations in wording may not significantly impact the quality of a summary, and multiple valid summaries may exist for the same content. Therefore, more nuanced evaluation methods like ROUGE metrics are necessary for assessing the performance of summarization models.

You may see some warning messages in the output above. While pipelines are handy, they provide less control over the tokenizer and model; we will dive deeper later. But first, let's see how our summarization pipeline does! We'll compute 0/1 accuracy, a classic ML evaluation metric.

```
Cmd 18 Python ▶ ▶ ▶ - x

1 accuracy = 0.0
2 for i in range(len(reference_summaries)):
3     generated_summary = t5_small_summaries[i]
4     if generated_summary == reference_summaries[i]:
5         accuracy += 1.0
6 accuracy = accuracy / len(reference_summaries)
7
8 print(f"Achieved accuracy {accuracy}!")

Achieved accuracy 0.0!
Command took 0.10 seconds -- by sam.raymond@databricks.com at 5/25/2023, 3:01:16 PM on LLM Course
```

Figure 3.4: Accuracy evaluation shows limitations for summarization tasks

3.2 Model Performance

Rouge is primarily used for comparing a computer-generated summary of text with a reference(i.e. human-generated) summary. Rouge is mostly used for evaluating text-summarization tasks. The value of Rouges metrics ranges from 0 to 1. 1 is a higher score, indicating that computer-generated summary and reference summary have a high similarity.

We evaluated three models on their summarization capabilities using ROUGE metrics: T5-small, FLAN-T5 Base, and GPT-3. The table below summarizes their performance.

Table 3.1: *Model Performance Summary using ROUGE Metrics*

Model	ROUGE-1	ROUGE-2	ROUGE-L
T5-small	0.45	0.32	0.40
FLAN-T5	0.47	0.34	0.42
GPT-3	0.50	0.28	0.35

3.2.1 Summary of the ROUGE Metrics

- **ROUGE-1:** Measures the overlap of single words between the generated summary and the reference summary.
- **ROUGE-2:** Evaluates the overlap of consecutive word pairs (bigrams).
- **ROUGE-3:** Assesses the overlap of three consecutive words (trigrams).
- **ROUGE-L:** Measures the longest common subsequence between the generated text and the reference, assessing the model's ability to preserve the overall structure and coherence of the content.

3.3 Model Analysis

- **T5-small:** This model, part of the T5 family, is a text-to-text transformer with 60 million parameters. It has been trained on a diverse mixture of unsupervised and supervised tasks, making it well-suited for summarization. The results show that T5-small performs well across all ROUGE metrics, indicating its effectiveness in generating high-quality summaries.

t5_small
a magnitude 6.7 earthquake rattles Papua new Guinea early Friday afternoon. no tsunami warning was issued, according to the Tsunami Warning Center. papu
the two-Test cricket series is being played in England due to security issues. the series is being played in England due to security issues in Pakistan. the series is being played in
federal prosecutors want jared Lee Loughner to submit a handwriting sample. prosecutors want the sample to compare with handwritten notes found in his residence. the federal grand
new: "he tried to kill people," a 17-year-old student says. new: "you could feel in the room all the anger that everyone had for him," she
double-amputee sprinter Oscar Pistorius will compete at the 2012 able bodied Olympics. the four-time paralympic gold medalist won a silver
new: a grand jury indicted the governor on charges of coercion of a public servant and abuse of his capacity. new: a grand jury indicted the

Figure 3.5: T5-small: A 60M parameter transformer excelling in ROUGE metrics for summarization

- **FLAN-T5 Base:** The FLAN-T5 Base model, an improved version of T5, showed slightly better performance in ROUGE-1 and ROUGE-L compared to T5-small.

This improvement suggests that FLAN-T5 Base may have better fine-tuning or additional training that enhances its summarization capabilities.

t5_base

the quake was centered about 200 miles north-northeast of Port Moresby. no tsunami warning was issued. Papua New Guinea is on the

Pakistan reach 148-3 on opening day of two-Test series against australia. teenagers mohammad aamer and mohammad asif take three wickets each

federal prosecutors want a handwriting sample from a man accused of killing six people. they want it to compare with notes that include references to guns and bullets. the government

"he's just this freak lookin' dude with some orange hair," teen says. "i think he needs to be killed," she says. she says she

double-amputee sprinter Oscar Pistorius will compete at the 2012 olympics. he was named in the individual 400m and 4x400m

new: governor's attorney calls indictment a "political abuse of the court system" new: he can continue to serve as governor while under indictment.

Figure 3.6: FLAN-T5 Base: A robust transformer with strong ROUGE performance, tailored for high-quality summarization

- **GPT-3:** GPT-3's higher ROUGE-1 score indicates its strength in generating content that aligns well with individual word choices, which can be beneficial for text completion tasks where predicting the next word is crucial. However, its lower ROUGE-2 and ROUGE-L scores compared to FLAN-T5 suggest that while GPT-3 excels in word-level accuracy, it may not preserve contextual relationships and summary coherence as effectively as FLAN-T5. This highlights GPT-3's suitability for tasks focused on word prediction and general text generation, rather than detailed summarization." GPT-3 employs an autoregressive transformer architecture, which focuses on the next-token prediction task. During training, it learns to predict the probability distribution of the next word (or token) in a sequence, given the previous words. This process is also known as causal or masked language modeling

GPT-3	
	"I'm not a racist, I'm a good player." The Czech Republic has also been fined \$1,000 for racist chantir against Italy on Saturday. The Czech Republic has also
	The U.S. has been working with Mexico on drug trafficking since the 1980s, and the U.S. has been w on drug trafficking since the 1990s. The new rules
	The escape was not discovered until a resident and Covington police reported seeing what appeared uniforms walking down a street. The escape was not discovered until a resident and Covington
	The storm is expected to bring heavy snow and heavy rain to the area. The National Weather Service storm is expected to bring heavy snow and heavy rain to the area. The National Weather
	The game will be played at the University of Texas at Austin. -- Louisiana Gov. Bobby Jindal said he w saddened" by the flooding in Louisiana. -- Louisiana Gov. Bobby Jindal said
	The film was nominated for best supporting actor, best supporting actress, best supporting actress, t actress, best supporting actress, best supporting actress, best supporting actress, best supporting a

Figure 3.7: GPT-3: A versatile model with high ROUGE-1 scores, excelling in content generation but less focused on summarization coherence

Chapter 4

CONCLUSION

The **GenQuery** exemplifies a successful application of advanced technologies to streamline academic research processes. By integrating LangChain’s Unstructured URL Loader, OpenAI’s embeddings, and FAISS for similarity search, the tool offers a robust solution for managing and retrieving extensive textual data. The architecture leverages transformer models, particularly for generating high-quality embeddings that capture the semantic nuances of the content, thereby facilitating accurate and relevant information retrieval.

The successful implementation of this project highlights several key technological and architectural achievements:

- **Efficient Data Retrieval:** The tool automates the extraction of content from URLs and text files, ensuring that users can quickly access large volumes of information without manual intervention. LangChain’s UnstructuredURL Loader plays a crucial role in processing raw web content into a structured format suitable for further analysis.
- **Accurate Embedding Generation:** Using OpenAI’s embeddings, the tool creates precise vector representations of the content. These embeddings capture the semantic essence of the text, which is crucial for effective similarity searches and query responses.
- **Swift Information Retrieval:** FAISS, a powerful similarity search library, is employed to enable fast and efficient retrieval of relevant information from the embedding vectors. This ensures that users receive prompt and accurate answers to their queries, enhancing the overall efficiency of the research process.

- **Seamless User Interaction:** The tool is designed to be user-friendly, providing an intuitive interface for querying and interacting with the embedded data. The integration of transformer models facilitates contextually relevant responses, making it easier for users to obtain the information they need without extensive manual effort.

In addition to these technological advancements, our evaluation of various models highlights the specific strengths and weaknesses relevant to different tasks. For example, GPT-3’s higher ROUGE-1 score indicates its strength in generating content that aligns well with individual word choices, which can be beneficial for text completion tasks where predicting the next word is crucial. However, its lower ROUGE-2 and ROUGE-L scores compared to FLAN-T5 suggest that while GPT-3 excels in word-level accuracy, it may not preserve contextual relationships and summary coherence as effectively as FLAN-T5. This highlights GPT-3’s suitability for tasks focused on word prediction and general text generation, rather than detailed summarization.

This nuanced analysis helps to align the strengths of GPT-3 with its application in text completion, emphasizing its suitability for tasks requiring precise word generation.

Bibliography

- [1] <https://platform.openai.com/docs/models>
- [2] *K. Pandya and M. Holia*, ‘Automating Customer Service using LangChain: Building custom open-source GPT Chatbot for organizations’, arXiv [cs.CL]. 2023
- [3] *M. Douze et al.*, ‘The Faiss library’, arXiv [cs.LG]. 2024.
- [4] <https://js.langchain.com/v0.2/docs/introduction/>
- [5] *A. Vaswani et al.*, ‘Attention is All you Need’, in Neural Information Processing Systems, 2017.
- [6] <https://docs.streamlit.io/>