

A PROJECT REPORT

On

Retrieval-Augmented Generation (RAG) Based Question Answering System

For the parliament fulfilling for the award of the degree of

BACHELOR OF TECHNOLOGY

In

Computer Science and Engineering

Submitted by

Priyanshi Sharma (202349624)

Vishal Dubey (2023000836)

Abhishek Ku(2024464096)

Under the supervision of

Mr. Ayush Singh



DEPARTMENT OF COMPUTER SCIENCE ENGINEERING

SHARDA UNIVERSITY

GREATER NOIDA

Declaration

We, the undersigned the work embodied this project work hereby, "**“Retrieval-Augmented Generation (RAG) Based Question Answering System”**" forms my own contribution to the research work carried out under the guidance of "**“Mr. Ayush Singh”**" is a result of my research work and his experience previously submitted in this project. We are Bachelor of Technology in Computer Science & Engineering Submitted to Sharda University, Greater Noida, is based on my own Work carried Out at the Department of Computer Science & and Engineering Sharda University, Greater Noida. The Work Contained in the report is original and the project work reported in this report has not been submitted by us for award of any other degree or diploma.

Signature:

Name: Priyanshi Sharma

System ID: 2023492624

Signature:

Name: Vishal Dubey

System ID: 2023343293

Signature:

Name: Abhishek Kumar Tiwari

System ID: 2024464096

Acknowledgement

The merciful guidance bestowed to us by the almighty made us stick to this project to a successful end. We humbly pray with sincere heart for his guidance to continue forever.

We pay thanks to our project guide **Mr. Ayush Singh** who has given guidance and light to us during this project. His versatile knowledge has helped us in the critical times during the span of this project.

We pay special thanks to our Head of Department **Dr. Sudeep Varshney** who has been always present as a support and help us in all possible ways during this project.

We also take this opportunity to express our gratitude to all those who have been directly and indirectly with us during the completion of the project.

We want to thank our friends who have always encouraged us during this project.

At last, thanks to all the faculty of CSE department who provided valuable suggestions during the period of project.

Abstract

Large Language Models often generate inaccurate or generic responses when answering domain-specific queries. Retrieval-Augmented Generation (RAG) improves response accuracy by retrieving relevant information from external documents before generating answers.

This project implements a RAG-based Question Answering system using a PDF document as the knowledge source. The system performs text extraction, chunking, embedding generation, similarity search using FAISS, and answer generation using a lightweight transformer model. A Streamlit interface is developed to allow users to interact with the system.

The proposed system improves reliability, reduces hallucinations, and enables efficient document-based question answering.

Keywords: *RAG, FAISS, Embeddings, NLP, Streamlit, Transformer*

Table of Content

Declaration	(ii)
Acknowledgement.....	(iii)
Abstract	(iv)
Chapter1. Introduction.....	6
Chapter2. Literature review.....	9
Chapter3. System design	14
Chapter4. Hardware and software requirements	19
Chapter5. Research methodology/Algorithm used.....	23
Chapter6. Implementation-coding	28
Chapter7. Results.....	34
Chapter8. Future enhancements.....	38
Chapter 9. Conclusion.....	41

Chapter-1

Introduction

Artificial Intelligence and Natural Language Processing (NLP) have experienced rapid growth in recent years, transforming the way humans interact with machines and access information. Large Language Models (LLMs) such as GPT, BERT, and T5 have demonstrated remarkable capabilities in text generation, summarization, translation, and question answering. These models are trained on massive datasets and are capable of understanding context, generating coherent responses, and performing a wide range of language-related tasks. As a result, AI-powered conversational systems, virtual assistants, and automated knowledge systems are becoming increasingly common in education, business, healthcare, and research.

Despite their impressive capabilities, traditional Large Language Models suffer from a fundamental limitation. These models rely primarily on the knowledge they learned during training and do not have direct access to external or updated information at runtime. When users ask domain-specific or factual questions outside the training distribution, the model may generate incorrect, outdated, or fabricated responses. This phenomenon, commonly referred to as *hallucination*, reduces the reliability and trustworthiness of AI systems, especially in applications where accuracy is critical.

Another important challenge arises when organizations need AI systems that work with their own private documents, research papers, manuals, or institutional data. Retraining large language models on custom datasets is computationally expensive, time-consuming, and often impractical. Furthermore, static training does not allow the system to dynamically update its knowledge when new information becomes available. Therefore, there is a growing need for AI systems that can combine the language understanding capabilities of LLMs with real-time access to external knowledge sources.

Retrieval-Augmented Generation (RAG) has emerged as an effective solution to these challenges. RAG is a hybrid architecture that integrates information retrieval techniques with generative language models. Instead of generating responses solely based on pre-trained knowledge, a RAG system first retrieves relevant information from an external knowledge base or document collection. The retrieved content is then provided as context to the language model, which generates a response grounded in the retrieved information. This approach significantly improves factual accuracy, reduces hallucinations, and enables domain-specific question answering without the need for extensive retraining.

The core idea behind RAG is simple yet powerful. When a user submits a query, the system converts the query into a semantic vector representation using an embedding model. This vector is compared with pre-computed vectors of document chunks stored in a vector database. The most relevant chunks are retrieved based on similarity search and passed to a language model along with the original query. The language model then generates a response that is informed by the retrieved

context. This ensures that the output is based on actual document content rather than generic or imagined knowledge.

With the advancement of embedding models such as Sentence Transformers and efficient similarity search libraries like FAISS (Facebook AI Similarity Search), it is now possible to build lightweight and efficient RAG systems that run locally without requiring large-scale cloud infrastructure. These systems enable fast semantic search over large document collections and support real-time question answering applications.

The importance of RAG systems is growing across multiple domains. In education, they can serve as intelligent study assistants capable of answering questions from textbooks or lecture notes. In enterprises, they can function as knowledge management tools that allow employees to query internal documentation. In research environments, they can help users quickly extract relevant information from large volumes of technical papers. Similarly, customer support systems can use RAG to provide accurate responses based on product manuals and FAQs.

This project focuses on the design and implementation of a lightweight Retrieval-Augmented Generation system for document-based question answering. The system uses a Machine Learning educational PDF as the knowledge source. The document is processed through text extraction and divided into smaller overlapping chunks to preserve contextual continuity. Each chunk is converted into vector embeddings using a pre-trained Sentence Transformer model. These embeddings are stored in a FAISS vector database to enable efficient similarity search.

When a user enters a query through the interface, the system generates an embedding for the query and retrieves the most relevant document chunks using FAISS. The retrieved context is then combined with the user query and passed to a lightweight transformer-based language model (FLAN-T5 Small) to generate a context-aware answer. A Streamlit-based web interface is developed to provide an interactive and user-friendly environment for testing the system.

The proposed system is designed with the following objectives:

- Improve the accuracy of AI-generated responses by grounding them in document content
- Reduce hallucinations and irrelevant outputs
- Enable domain-specific question answering without retraining large models
- Provide a lightweight solution that runs locally without external APIs
- Develop an interactive interface for real-time user interaction

By integrating document retrieval with generative AI, the project demonstrates how modern NLP techniques can be used to build reliable and practical knowledge-based systems. The implementation highlights the effectiveness of semantic search, vector databases, and transformer models working together in a unified pipeline.

As organizations increasingly rely on data-driven decision-making and intelligent information systems, Retrieval-Augmented Generation represents a significant step toward more trustworthy and adaptable AI solutions. This project contributes to the growing field of applied AI by providing a practical implementation of RAG architecture that can be extended to multiple real-world applications such as educational assistants, enterprise search systems, and domain-specific conversational agents.

Chapter-2

Literature review

The field of question answering and information retrieval has evolved significantly over the past few decades. With the rapid growth of digital information, the need for intelligent systems capable of understanding natural language queries and retrieving accurate information has become increasingly important. Researchers have explored a wide range of approaches, beginning with traditional information retrieval methods and progressing toward modern deep learning and Retrieval-Augmented Generation (RAG) architectures.

2.1 Traditional Information Retrieval Systems

Early question answering systems were primarily based on keyword matching and rule-based retrieval techniques. These systems relied on classical Information Retrieval (IR) models such as:

- Boolean Retrieval Model
- Vector Space Model
- TF-IDF (Term Frequency–Inverse Document Frequency)
- BM25 ranking algorithm

In these approaches, documents were indexed based on word frequency, and user queries were matched using lexical similarity. Although these methods were efficient and easy to implement, they suffered from several limitations:

- Inability to understand semantic meaning
- Poor performance when synonyms or paraphrases were used
- Limited contextual understanding
- Dependence on exact keyword matches

As a result, traditional IR systems often failed to retrieve relevant information when the wording of the query differed from the document text.

2.2 Neural and Embedding-Based Retrieval

With the advancement of deep learning, researchers introduced neural embedding techniques to capture semantic relationships between words and sentences. Models such as Word2Vec, GloVe, and FastText enabled the representation of words in continuous vector spaces, allowing systems to measure semantic similarity rather than relying solely on exact matches.

Later, transformer-based models such as BERT (Bidirectional Encoder Representations from Transformers) significantly improved semantic understanding. Sentence-level embedding models, including Sentence-BERT (SBERT), allowed entire sentences and paragraphs to be converted into dense vector representations.

Embedding-based retrieval systems offer several advantages:

- Semantic search capability
- Improved handling of synonyms and paraphrasing
- Better contextual understanding
- Higher retrieval accuracy compared to keyword-based methods

To support large-scale similarity search, efficient vector indexing libraries such as FAISS (Facebook AI Similarity Search) were developed. FAISS enables fast nearest-neighbor search over millions of high-dimensional vectors and has become a standard component in modern semantic search systems.

2.3 Transformer-Based Language Models for Question Answering

The introduction of transformer architectures revolutionized Natural Language Processing. Models such as GPT, BERT, T5, and FLAN-T5 demonstrated strong performance in language understanding and generation tasks.

Pre-trained Large Language Models (LLMs) can:

- Generate human-like responses
- Understand context and intent
- Perform zero-shot and few-shot learning
- Answer questions based on learned knowledge

However, these models have several limitations:

- Knowledge is fixed at training time
- Cannot access external or updated information
- May generate incorrect or fabricated answers (hallucination)
- High computational cost for retraining

These challenges highlighted the need for systems that combine generative models with external knowledge sources.

2.4 Retrieval-Augmented Generation (RAG)

Retrieval-Augmented Generation was introduced as a hybrid approach that integrates information retrieval with text generation. Instead of relying solely on pre-trained knowledge, RAG systems retrieve relevant documents and use them as context for response generation.

The general architecture of RAG includes:

1. Document preprocessing and chunking
2. Embedding generation
3. Vector database storage
4. Query embedding
5. Similarity search
6. Context-based response generation

Research studies have shown that RAG significantly improves:

- Factual accuracy
- Domain adaptation
- Response reliability
- Reduction in hallucinations

RAG has been successfully applied in:

- Open-domain question answering
- Enterprise knowledge systems
- Educational assistants
- Customer support automation

2.5 Vector Databases and Similarity Search

Efficient retrieval is a critical component of RAG systems. Several vector databases have been developed for large-scale semantic search, including:

- FAISS
 - Pinecone
-

- Weaviate
- Milvus

Among these, FAISS is widely used for local implementations due to:

- High-speed similarity search
- Low memory overhead
- CPU-based operation support
- Easy integration with Python

Studies indicate that approximate nearest-neighbor search methods significantly reduce retrieval time while maintaining high accuracy.

2.6 Lightweight Language Models for Local Deployment

While large models such as GPT-4 provide high-quality responses, they require significant computational resources and cloud-based access. For local or resource-constrained environments, lightweight transformer models such as FLAN-T5 Small and DistilBERT are preferred.

These models offer:

- Faster inference
- Lower memory requirements
- Offline execution capability
- Suitable performance for domain-specific tasks

Recent research emphasizes the importance of lightweight AI systems for educational and small-scale applications where cloud dependency is not desirable.

2.7 Research Gap

Although significant progress has been made in retrieval and generation techniques, several challenges remain:

- Many systems rely on cloud-based APIs
- Large models require high computational resources
- Limited availability of lightweight, locally deployable RAG systems
- Need for simple interfaces for non-technical users

This project addresses these gaps by developing a lightweight RAG-based question answering system that runs locally using FAISS, Sentence Transformers, and a small transformer model integrated with a Streamlit interface.

Chapter-3

System design/architecture

System Design and Architecture

The development of a Retrieval-Augmented Generation (RAG) based question answering system requires a structured architecture capable of handling document processing, semantic search, and context-based response generation. The proposed system is designed as a modular pipeline where each component performs a specific task in transforming raw documents into an intelligent interactive knowledge system.

The architecture ensures efficient document indexing, fast similarity search, accurate context retrieval, and reliable response generation while maintaining low computational requirements for local execution.

3.1 Overall System Architecture

The system follows a sequential pipeline consisting of the following major components:

- Document Input
- Text Extraction
- Text Chunking
- Embedding Generation
- Vector Storage (FAISS)
- User Query Processing
- Similarity Search
- Context Retrieval
- Response Generation
- Streamlit Interface Output

This layered architecture ensures scalability, modularity, and efficient processing

3.2 Document Processing Layer

3.2.1 Document Input

The knowledge source for the system is a PDF document containing machine learning educational content. The document is stored locally and used as the primary data source for retrieval.

3.2.2 Text Extraction

The PDF is processed using a document parsing library (PyPDF). The extracted text is cleaned to remove unnecessary formatting characters, extra spaces, and special symbols.

This step converts the unstructured document into machine-readable text format.

3.3 Text Chunking Layer

Large documents cannot be directly processed efficiently for semantic search. Therefore, the extracted text is divided into smaller overlapping segments called **chunks**.

Configuration:

- Chunk size: 500 characters
- Overlap: 50 characters

Purpose of Chunking

- Maintains contextual continuity between segments
- Prevents loss of information at boundaries
- Improves retrieval precision
- Optimizes memory usage

Each chunk represents a meaningful piece of information that can be independently retrieved during query processing.

3.4 Embedding Generation Layer

Each text chunk is converted into a numerical vector representation using a pre-trained sentence embedding model.

Model used:

sentence-transformers/all-MiniLM-L6-v2

Advantages

- Lightweight and fast
- Good semantic understanding
- Works efficiently on CPU
- Suitable for local deployment

These embeddings capture semantic meaning rather than just keyword frequency, enabling accurate similarity matching.

3.5 Vector Database Layer

The generated embeddings are stored in a vector database using **FAISS (Facebook AI Similarity Search)**.

Role of FAISS

- Stores high-dimensional vectors
- Performs fast nearest-neighbor search
- Enables real-time semantic retrieval
- Works locally without cloud dependency

The index structure allows efficient comparison between query embeddings and stored document embeddings.

3.6 Query Processing Layer

When a user enters a question through the interface, the system performs the following steps:

1. Accept user input
2. Convert the query into an embedding using the same MiniLM model
3. Search the FAISS index for the most similar document chunks
4. Retrieve top-k relevant chunks (context)

This ensures that only relevant information is passed to the language model.

3.7 Context-Based Response Generation

The retrieved chunks are combined with the user query and provided as input to a lightweight transformer-based language model.

Model used:

google/flan-t5-small

Function of the Language Model

- Understand the query
- Interpret retrieved context
- Generate a coherent and accurate answer
- Avoid hallucination by grounding response in document content

The use of a lightweight model ensures faster inference and local execution.

3.8 User Interface Layer

The system includes a web-based interface built using **Streamlit**.

Features

- Simple text input box for user queries
- Display of generated answers
- Real-time interaction
- Local browser access

Run command:

`streamlit run app.py`

The interface allows non-technical users to interact easily with the system.

3.9 System Workflow Summary

The complete workflow of the system can be summarized as:

1. Load PDF document
2. Extract and clean text

3. Divide text into chunks
4. Generate embeddings
5. Store embeddings in FAISS
6. Accept user query
7. Convert query into embedding
8. Retrieve similar chunks
9. Generate answer using FLAN-T5
10. Display result via Streamlit

3.10 Advantages of the Proposed Architecture

- Modular and scalable design
- Lightweight and efficient
- Works without external APIs
- Reduces hallucinations
- Enables domain-specific question answering
- Fast response time
- Easy to extend for multiple documents

3.11 Scope for Extension

The architecture can be extended to support:

- Multiple document sources
- Persistent vector storage
- Chat-based conversational memory
- Cloud deployment
- Integration with enterprise knowledge systems

Chapter-4

Hardware and software requirements

The successful implementation of a Retrieval-Augmented Generation (RAG) system requires an appropriate combination of hardware resources and software tools. Since the proposed system is designed as a lightweight, locally executable solution, it does not require high-end computational infrastructure. However, sufficient processing capability and memory are necessary to handle document processing, embedding generation, vector indexing, and model inference efficiently.

This chapter describes the minimum and recommended hardware and software requirements for developing and running the system.

4.1 Hardware Requirements

The hardware requirements depend on the size of the document collection and the complexity of the models used. Since this project uses lightweight embedding and generation models, it can run efficiently on standard personal computers.

4.1.1 Minimum Hardware Requirements

These specifications are sufficient for running the system with a single or small number of documents:

- Processor: Intel Core i3 (7th generation or above) / AMD equivalent / Apple M-series (base model)
- RAM: 4 GB
- Storage: 2–5 GB free disk space
- Graphics: Integrated GPU (not mandatory)
- System Type: 64-bit operating system

Suitable for:

- PDF processing
- Embedding generation (CPU-based)
- FAISS indexing
- Running Streamlit interface

4.1.2 Recommended Hardware Requirements

For smoother performance and faster response times, the following configuration is recommended:

- Processor: Intel Core i5/i7 / AMD Ryzen 5/7 / Apple M-series
- RAM: 8–16 GB
- Storage: SSD with at least 10 GB free space
- GPU: Optional (not required for this project)
- System Type: 64-bit operating system

Suitable for:

- Faster embedding generation
- Handling larger documents or multiple files
- Improved Streamlit performance
- Reduced query response time

4.1.3 High-Performance Setup (Optional)

If the system is extended to support:

- Multiple large documents
- Large language models
- Real-time multi-user deployment

Then the following may be required:

- Processor: Multi-core CPU
- RAM: 16–32 GB
- GPU: CUDA-enabled GPU (NVIDIA)
- Cloud platforms: AWS, Google Cloud, or Azure

4.2 Software Requirements

4.2.1 Operating System

The system can be implemented on any modern operating system:

- Windows 10/11

- Linux (Ubuntu recommended for advanced deployments)
- macOS (Intel and Apple Silicon supported)

4.2.2 Programming Language

- Python 3.8 or above

Python is used for:

- Document processing
- Embedding generation
- Vector indexing
- Model integration
- Web interface development

4.2.3 Python Libraries Used

Data Processing

- NumPy
- Pandas

Document Handling

- PyPDF

Embedding Model

- sentence-transformers

Vector Database

- faiss-cpu

Language Model

- transformers
- torch

Web Interface

- streamlit

4.2.4 Installation of Dependencies

Required libraries can be installed using:

```
pip install streamlit sentence-transformers faiss-cpu transformers pypdf torch
```

4.2.5 Development Tools

- VS Code / PyCharm – Code development
- Jupyter Notebook (optional) – Testing and experimentation
- Git (optional) – Version control

4.3 System Environment

Execution Steps:

1. Place the PDF file in the project folder
2. Install required dependencies
3. Run the application using:

```
streamlit run app.py
```

4. Open browser at:

<http://localhost:8501>

4.4 Advantages of the Setup

- Works entirely offline
- No external API or cloud dependency
- Low computational cost
- Easy to install and run
- Suitable for educational and small-scale applications

Chapter-5

Research methodology/Algorithm used

The proposed Retrieval-Augmented Generation (RAG) system follows a structured methodology to enable accurate and context-aware question answering from a document. The methodology integrates document processing, semantic embedding, vector-based retrieval, and transformer-based text generation into a unified pipeline.

The overall workflow consists of the following stages:

Document Collection
→ Text Extraction
→ Text Chunking
→ Embedding Generation
→ Vector Storage
→ Query Processing
→ Similarity Search
→ Context Retrieval
→ Answer Generation

Each stage is described in detail below.

5.1 Data Source

The knowledge base for the system is a Machine Learning educational PDF document. This document serves as the domain-specific information source from which answers are generated.

Characteristics:

- Unstructured text format
- Multiple topics and sections
- Large size, requiring segmentation for efficient processing

5.2 Document Preprocessing

5.2.1 Text Extraction

The PDF document is processed using the PyPDF library to extract textual content.

Steps:

- Load PDF file
-

- Read text page by page
- Combine extracted text into a single document
- Remove unnecessary line breaks and special characters

This step converts the document into machine-readable format.

5.3 Text Chunking

Large documents cannot be directly used for semantic search due to memory and performance limitations. Therefore, the extracted text is divided into smaller overlapping segments called **chunks**.

Configuration:

- Chunk size: 500 characters
- Overlap: 50 characters

Purpose of Chunking

- Maintains contextual continuity
- Prevents loss of important information at boundaries
- Improves retrieval accuracy
- Reduces computational load

Each chunk acts as an independent retrieval unit.

5.4 Embedding Generation

To enable semantic search, each text chunk is converted into a numerical vector representation.

Model Used

sentence-transformers/all-MiniLM-L6-v2

Working Principle

- Converts text into dense vector embeddings
- Captures semantic meaning rather than keyword frequency
- Similar texts produce similar vectors

Advantages

- Lightweight and fast
- High semantic accuracy
- Efficient CPU performance
- Suitable for local deployment

5.5 Vector Database Creation

The generated embeddings are stored using **FAISS (Facebook AI Similarity Search)**.

Role of FAISS

- Stores high-dimensional vectors
- Enables fast nearest-neighbor search
- Supports large-scale similarity comparison
- Works locally without internet access

Indexing Process

1. Convert all chunk embeddings into a matrix
2. Create FAISS index
3. Add embeddings to index
4. Store mapping between vectors and original text

5.6 Query Processing

When a user enters a question, the system performs the following steps:

1. Accept query input
2. Convert query into embedding using the same MiniLM model
3. Compare query vector with stored vectors in FAISS
4. Retrieve top-k most similar chunks

This ensures that only relevant document content is selected.

5.7 Similarity Search Algorithm

FAISS performs similarity search using distance metrics such as:

- Cosine similarity
- Euclidean distance

The algorithm identifies the nearest vectors to the query embedding and returns the corresponding text chunks.

Configuration:

- Top-k retrieval (typically 3–5 chunks)

5.8 Context Construction

The retrieved chunks are combined to form the **context**.

Context = Top retrieved chunks + User query

This context is passed to the language model to generate a grounded response.

Purpose:

- Improves factual accuracy
- Reduces hallucination
- Ensures answer relevance

5.9 Response Generation

Model Used

[google/flan-t5-small](#)

Working

1. Input: Query + Retrieved Context
2. Model interprets context
3. Generates natural language response

Advantages

- Lightweight transformer model
- Instruction-tuned for question answering

- Fast inference
- Works locally without GPU

5.10 User Interaction

The system uses a **Streamlit interface** for real-time interaction.

Features:

- Text input for queries
- Instant response generation
- Simple web interface
- Local browser execution

5.11 Algorithm Workflow Summary

Step 1: Load PDF

Step 2: Extract text

Step 3: Split into chunks

Step 4: Generate embeddings

Step 5: Store embeddings in FAISS

Step 6: Accept user query

Step 7: Convert query to embedding

Step 8: Retrieve similar chunks

Step 9: Generate answer using FLAN-T5

Step 10: Display result

5.12 Advantages of the Methodology

- Combines retrieval and generation
- Reduces incorrect or fabricated answers
- Works offline without external APIs
- Lightweight and efficient
- Easily extendable to multiple documents

Chapter-6

Implementation

This chapter describes the practical implementation of the Retrieval-Augmented Generation (RAG) based question answering system. The system is developed using Python and various Natural Language Processing and machine learning libraries. The implementation follows a modular approach consisting of document processing, embedding generation, vector indexing, query handling, and response generation.

6.1 Development Environment

Programming Language: Python 3.8+

Development Tools:

- Visual Studio Code
- Jupyter Notebook (for testing)

Libraries Used:

- streamlit
- sentence-transformers
- faiss-cpu
- transformers
- torch
- pypdf
- numpy

6.2 Document Loading and Text Extraction

The first step is to load the PDF document and extract text content.

Steps:

1. Import PyPDF library
 2. Open the PDF file
 3. Read text page by page
-

4. Store extracted text in a single variable

Example:

```
from pypdf import PdfReader
```

```
reader = PdfReader("ml_intro.pdf")
```

```
text = ""
```

```
for page in reader.pages:
```

```
    text += page.extract_text()
```

This converts the document into machine-readable format for further processing.

6.3 Text Chunking

The extracted text is divided into smaller chunks to improve retrieval efficiency.

Configuration:

- Chunk size: 500 characters
- Overlap: 50 characters

Example:

```
def split_text(text, chunk_size=500, overlap=50):
```

```
    chunks = []
```

```
    start = 0
```

```
    while start < len(text):
```

```
        end = start + chunk_size
```

```
        chunks.append(text[start:end])
```

```
        start += chunk_size - overlap
```

```
    return chunks
```

```
chunks = split_text(text)
```

Chunking ensures better semantic matching and context preservation.

6.4 Embedding Generation

Each text chunk is converted into a vector using Sentence Transformers.

Model used:

all-MiniLM-L6-v2

Example:

```
from sentence_transformers import SentenceTransformer
```

```
model = SentenceTransformer('all-MiniLM-L6-v2')
```

```
embeddings = model.encode(chunks)
```

These embeddings capture semantic meaning of the text.

6.5 FAISS Vector Index Creation

The generated embeddings are stored in a FAISS index for fast similarity search.

Example:

```
import faiss
```

```
import numpy as np
```

```
dimension = embeddings.shape[1]
```

```
index = faiss.IndexFlatL2(dimension)
```

```
index.add(np.array(embeddings))
```

FAISS allows efficient nearest-neighbor search during query processing.

6.6 Query Processing and Retrieval

When a user enters a query:

Steps:

1. Convert query into embedding

2. Search FAISS index
3. Retrieve top relevant chunks

Example:

```
def retrieve(query, k=3):  
    query_embedding = model.encode([query])  
    distances, indices = index.search(query_embedding, k)  
    return [chunks[i] for i in indices[0]]
```

This retrieves the most relevant document sections.

6.7 Response Generation

The retrieved chunks are combined with the query and passed to a language model.

Model used:

google/flan-t5-small

Example:

```
from transformers import pipeline
```

```
generator = pipeline("text2text-generation", model="google/flan-t5-small")
```

```
def generate_answer(query, context):  
    input_text = f"Context: {context} \n Question: {query}"  
    result = generator(input_text, max_length=200)  
    return result[0]['generated_text']
```

This generates a context-based answer.

6.8 Streamlit User Interface

A simple web interface is created using Streamlit.

Features:

- Text input box
- Answer display
- Real-time interaction

Example:

```
import streamlit as st
```

```
st.title("RAG-based QA System")
```

```
query = st.text_input("Enter your question")
```

```
if query:
```

```
    retrieved_chunks = retrieve(query)
    context = " ".join(retrieved_chunks)
    answer = generate_answer(query, context)
    st.write(answer)
```

Run the application:

```
streamlit run app.py
```

Open in browser:

```
http://localhost:8501
```

6.9 Implementation Workflow Summary

1. Load PDF document
2. Extract text
3. Split into chunks
4. Generate embeddings
5. Store embeddings in FAISS
6. Accept user query

7. Retrieve relevant chunks
8. Generate answer using FLAN-T5
9. Display result via Streamlit

6.10 Advantages of Implementation

- Lightweight and efficient
- Runs locally without external APIs
- Fast response time
- Easy to deploy and extend
- User-friendly interface

Results and Discussion

This chapter presents the results obtained from the implementation of the Retrieval-Augmented Generation (RAG) based Question Answering System. It evaluates the system's performance in terms of retrieval accuracy, response quality, execution efficiency, and overall usability.

The system was tested using a Machine Learning PDF document as the knowledge source and multiple user queries related to the document content.

7.1 Experimental Setup

The system was executed in a local environment with the following configuration:

- Platform: Local machine
- Embedding Model: all-MiniLM-L6-v2
- Vector Database: FAISS
- Language Model: FLAN-T5 Small
- Interface: Streamlit

The document was processed into text chunks and indexed before testing user queries.

7.2 Document Processing Results

- Total document pages processed: Multiple pages
- Total text chunks generated: Approximately 300–600 (depending on document size)
- Chunk size: 500 characters
- Overlap: 50 characters

The chunking process ensured that contextual continuity was maintained and that relevant information could be retrieved efficiently.

7.3 Retrieval Performance

The FAISS vector database was tested for similarity search performance.

Observations:

- Query embedding generation time: < 1 second
- Retrieval time for top-3 chunks: < 0.1 seconds
- Relevant content was retrieved accurately for most domain-specific queries

The semantic embedding approach successfully retrieved context even when the query wording differed from the document text.

Example:

Query: “*What is supervised learning?*”

Retrieved chunks contained the definition and explanation from the document.

7.4 Response Generation Results

The FLAN-T5 Small model generated responses using the retrieved context.

Observations:

- Response generation time: 1–3 seconds
- Answers were context-based and relevant
- Minimal hallucination observed
- Responses were concise and understandable

Example Interaction:

User Query:

What is machine learning?

System Output:

Machine learning is a subset of artificial intelligence that enables systems to learn patterns from data and make predictions or decisions without being explicitly programmed.

7.5 Accuracy Evaluation (Qualitative)

Since the system is document-based, performance was evaluated qualitatively based on:

Criteria	Observation
-----------------	--------------------

Context relevance	High
-------------------	------

Criteria	Observation
----------	-------------

Answer correctness	High
--------------------	------

Hallucination	Low
---------------	-----

Response clarity	Good
------------------	------

Retrieval precision	High
---------------------	------

The system consistently generated answers grounded in the document content.

7.6 System Performance

Operation	Time Taken
-----------	------------

Document indexing (one-time)	10–30 seconds
------------------------------	---------------

Query processing	< 1 second
------------------	------------

Retrieval	< 0.1 seconds
-----------	---------------

Answer generation	1–3 seconds
-------------------	-------------

The system performs efficiently on CPU without requiring GPU acceleration.

7.7 User Interface Evaluation

The Streamlit interface provided:

- Simple query input
- Instant response display
- Local browser access
- Easy usability

The interface was tested with multiple queries and functioned smoothly without errors.

7.8 Advantages Observed

- Accurate document-based answers
 - Reduced hallucination compared to standalone language models
-

- Fast response time
- Fully offline execution
- Lightweight and easy to deploy

7.9 Limitations Observed

- Answers limited to document content only
- Performance depends on chunk quality
- Limited response length due to small model
- No conversational memory

7.10 Discussion

The results demonstrate that integrating semantic retrieval with a lightweight language model significantly improves answer reliability. The embedding-based retrieval ensures that relevant context is provided to the model, which reduces incorrect or fabricated responses.

The system performs well for domain-specific question answering and can be extended for multiple documents, enterprise knowledge bases, or educational assistants.

Chapter-8

Future Enhancements

Although the developed Retrieval-Augmented Generation (RAG) system performs effectively for document-based question answering, several improvements can be made to enhance its scalability, performance, and real-world usability. The following enhancements can be considered for future work.

8.1 Multi-Document Support

The current system processes a single PDF document as the knowledge source. Future versions can be extended to support:

- Multiple PDF files
- Word documents and text files
- Entire document repositories

This would allow the system to function as a complete knowledge management solution.

8.2 Persistent Vector Storage

Currently, embeddings are generated each time the application is restarted. Future implementation can include:

- Saving FAISS index to disk
- Loading pre-built index at startup

This will reduce initialization time and improve system efficiency.

8.3 Chat-Based Interaction

The current system follows a single question–answer format. It can be enhanced by:

- Adding conversational memory
- Maintaining previous context
- Supporting multi-turn dialogue

This will improve user experience and make the system behave like a conversational assistant.

8.4 Hybrid Retrieval Techniques

Future improvements may include:

- Combining keyword search with semantic search
- Implementing hybrid ranking methods
- Using reranking models for improved retrieval accuracy

This can further enhance the relevance of retrieved content.

8.5 Use of Advanced Language Models

The current system uses a lightweight model for local execution. Future work may include:

- Integration with larger transformer models
- Use of cloud-based APIs for higher-quality responses
- Fine-tuning models for domain-specific tasks

This would improve response quality and generation capability.

8.6 Web and Cloud Deployment

The system can be deployed as:

- A cloud-based web application
- An enterprise knowledge assistant
- A public question-answering platform

Cloud deployment would enable multi-user access and better scalability.

8.7 Support for Additional Data Sources

Future versions may integrate:

- Databases
- Web content
- Real-time APIs

This would allow the system to provide updated and dynamic information.

8.8 Security and Access Control

For enterprise environments, the system can be enhanced with:

- User authentication
- Role-based document access
- Data privacy and encryption

This will make the system suitable for organizational use.

8.9 Performance Optimization

Further optimizations may include:

- GPU acceleration
- Advanced FAISS indexing methods
- Distributed vector storage

This will improve performance for large-scale deployments.

Chapter-9

Conclusion

The rapid growth of Artificial Intelligence and Natural Language Processing has created a strong demand for intelligent systems capable of understanding user queries and providing accurate information. However, traditional Large Language Models often generate responses based only on pre-trained knowledge, which may result in outdated or incorrect answers when dealing with domain-specific queries.

This project presented the design and implementation of a Retrieval-Augmented Generation (RAG) based Question Answering System that integrates semantic search with a transformer-based language model. The system processes a PDF document, converts the text into semantic embeddings, and stores them in a FAISS vector database. When a user submits a query, the system retrieves the most relevant document sections and generates a response grounded in the retrieved context.

The implementation demonstrates that combining document retrieval with text generation significantly improves response reliability and reduces hallucination. The use of lightweight models and local processing enables the system to operate efficiently without requiring high computational resources or external APIs. The Streamlit interface provides a simple and user-friendly platform for real-time interaction.

The system was tested using multiple domain-specific queries, and the results showed that it successfully retrieved relevant content and generated accurate, context-aware responses within a short processing time. The overall performance indicates that the proposed approach is suitable for educational assistants, document search systems, and small-scale knowledge management applications.

Although the current implementation is limited to a single document and lightweight models, the architecture is modular and scalable. With further enhancements such as multi-document support, conversational capabilities, and cloud deployment, the system can be extended for real-world enterprise and large-scale applications.

In conclusion, the project successfully demonstrates the practical implementation of Retrieval-Augmented Generation and highlights its importance in building reliable, domain-specific AI systems. The integration of semantic retrieval and generative models represents a significant step toward more trustworthy and efficient intelligent information systems.
