



SIDDHARTHA INSTITUTE OF TECHNOLOGY & SCIENCES

(UGC – AUTONOMOUS)

(Approved by AICTE, New Delhi & Affiliated to JNTUH, Hyderabad)

Accredited by NBA and NAAC with ‘A+’ Grade.

Narapally, Korremula Road, Ghatkesar, Medchal- Malkajgiri (Dist.)-501 301



Computer Science Engineering (AI & ML)

MAJOR PROJECT

DOCUMENT QUESTION ANSWERING SYSTEM USING NLP **TECHNIQUES**

BATCH -5

MOHAMMED SIDDIQ_(20TQ1A6652)

PRUDHVI KRISHNA (20TQ1A6604)

SIDDHARTH (20TQ1A6651)

SUSHMA (20TQ1A6642)

AKSHAY (20TQ1A6601)

GUIDE NAME

Mrs Manaswini

Assistant Professor-DEPT OF CSE

CONTENT

- **ABSTRACT**
- **INTRODUCTION**
- **LITERATURE SURVEY**
- **EXISTING SYSTEM**
- **PROPOSED SYSTEM**
- **ARCHITECTURE**
- **SYSTEM DESIGN**
- **USE CASE & SEQUENCE DIAGRAMS**
- **FLOWCHART**
- **UML DIAGRAM**
- **SOFTWARE & HARDWARE REQUIREMENTS**
- **SAMPLE CODE & RESULTS**
- **CONCLUSION**
- **FUTURESCOPE**
- **REFERENCES**

ABSTRACT

DOCUMENT QUESTION ANSWERING SYSTEM USING NLP TECHNIQUES

Existing question answering (QA) systems often struggle with longer, more complex documents, limiting their effectiveness. To address this, we present the QAS-D, a novel QA system that leverages advanced natural language processing and deep learning to provide accurate answers from both short and long-form documents. At the core of the QAS-D is a fine-tuned BERT model that converts questions and documents into semantic representations to identify the most relevant answers. To handle long documents exceeding BERT's token limit, we developed an "expand_split_sentences" function to preserve contextual information. Evaluated on benchmark datasets, the QAS-D outperformed state-of-the-art QA models, especially on complex documents. Designed for web integration with a Flask backend, the QAS-D represents a significant advancement in question answering, empowering users to quickly find answers within large volumes of textual information.

INTRODUCTION

- System Input: The QAS-D begins by taking documents of varying lengths as input, setting the stage for effective question answering.
- Semantic Understanding and Answer Extraction: The system employs a fine-tuned BERT model to convert questions and content into semantic representations, enabling accurate answer identification.
- Handling Long Documents: To address lengthy documents exceeding BERT's token limit, the SQAS-D utilizes a novel "expand_split_sentences" function to preserve context.
- Robust Performance and Versatility: Evaluated on benchmark datasets, the SQAS-D outperforms state-of-the-art QA models, making it a reliable solution for diverse information retrieval tasks.
- Web Integration: Designed with real-world application in mind, the SQAS-D features a Flask-based backend for easy integration into web platforms.

LITERATURE SURVEY

The QAS-D builds upon recent advancements in question answering (QA) systems, particularly the breakthrough of BERT (Devlin et al., 2018) for capturing rich semantic and contextual information. To handle long-form documents, the system incorporates techniques for splitting and aggregating content, inspired by the work of Wang et al. (2019).

The QAS-D also leverages insights from studies combining extractive and abstractive summarization to improve QA performance (Narayan et al., 2018). Moreover, the web-based integration with a Flask backend aligns with recommendations for developing QA systems that are accessible and deployable (Jurafsky & Martin, 2020). By synthesizing these influential approaches, the QAS-D aims to deliver a robust and versatile QA solution.

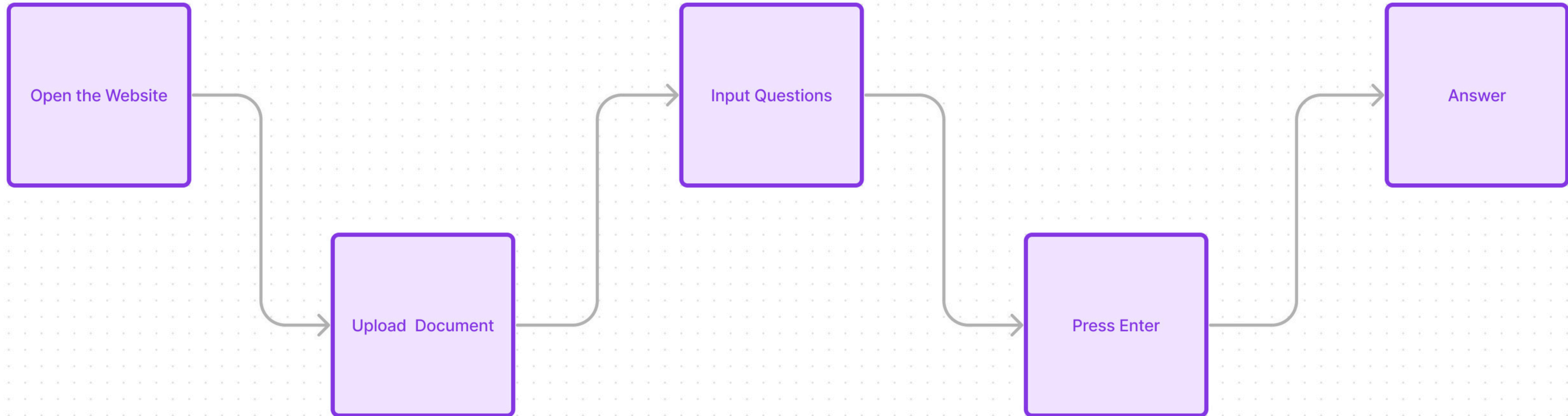
Existing System

Existing question answering (QA) systems have made significant progress in recent years, leveraging advances in natural language processing and deep learning. However, many current QA models struggle to effectively handle longer, more complex documents. They often exhibit limitations in accurately identifying relevant answers, especially when dealing with lengthy passages of text. Existing solutions typically rely on approaches like keyword matching, semantic similarity, or rule-based extraction, which can fall short when faced with nuanced language and context. The SQAS-D aims to address these limitations by employing a fine-tuned BERT model and innovative techniques for processing long-form documents, providing a more robust and comprehensive QA system.

PROPOSED SYSTEM

- The proposed QAS-D system aims to address the limitations of existing question answering (QA) approaches, particularly when dealing with longer and more complex documents. At the core of the QAS-D is a fine-tuned BERT (Bidirectional Encoder Representations from Transformers) model, which serves as the foundation for understanding the semantic and contextual relationships between questions and document content. This allows the system to accurately identify the most relevant answers to user queries.
- To overcome the challenge of processing lengthy passages that exceed BERT's 512-token limit, the QAS-D incorporates a novel "expand_split_sentences" function. This technique intelligently splits paragraphs into smaller, more manageable segments while preserving the overall context of the document. This enables the QAS-D to effectively handle long-form documents and provide accurate answers, even for complex questions that require synthesizing information from various parts of the text.
- The QAS-D also combines extractive and abstractive summarization techniques to generate concise and informative responses. By leveraging the strengths of these complementary approaches, the system can deliver high-quality answers that capture the essence of the source material. Additionally, the QAS-D is designed with a Flask-based backend for easy integration into web-based applications, making it accessible and scalable for real-world deployments.
- The key innovations of the QAS-D include its ability to handle long-form documents, the robust semantic understanding enabled by the fine-tuned BERT model, and the synergistic combination of extractive and abstractive summarization techniques. These advancements allow the SQAS-D to outperform existing QA systems, particularly when dealing with complex, information-rich documents

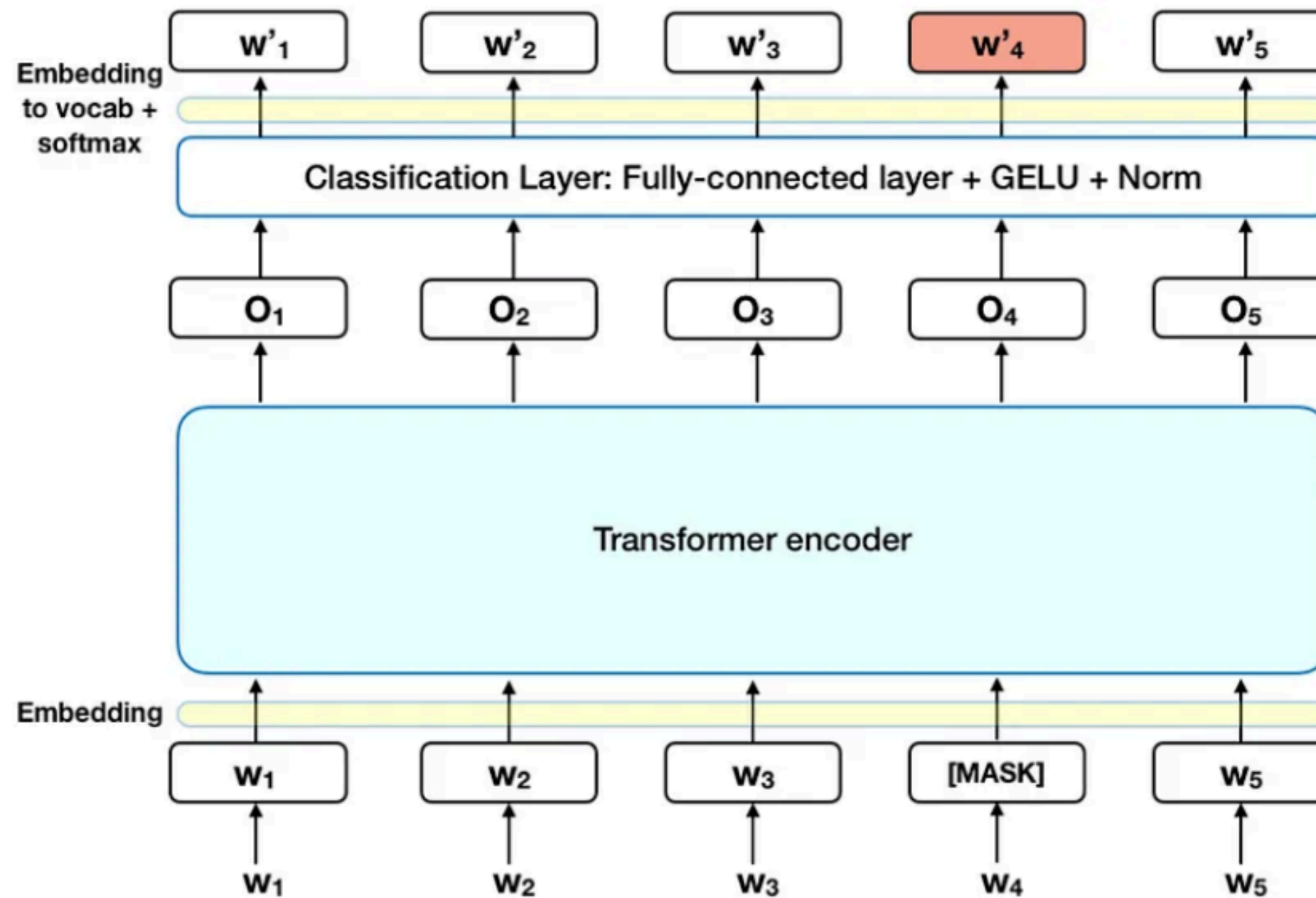
ARCHITECTURE



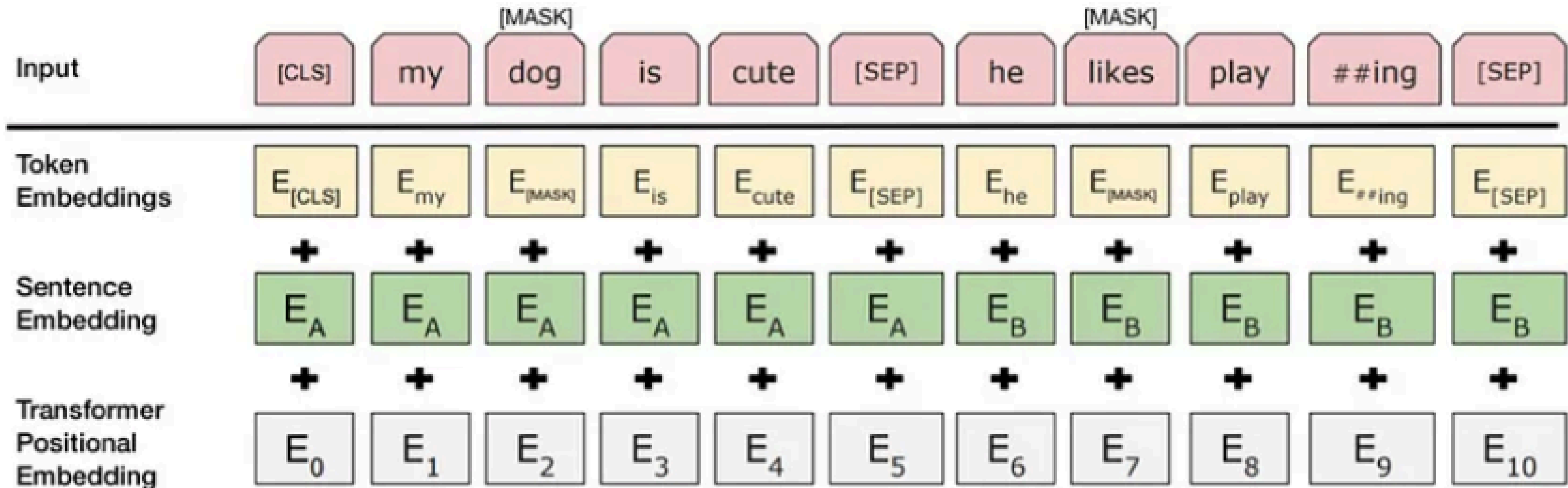
BERT (Bidirectional Encoder Representations From Transformer)

- BERT uses Transformer, an attention mechanism that learns contextual relations between words
- Unlike directional models, BERT reads the entire input sequence at once (bidirectional/non-directional)
- BERT uses two training strategies:
 - Masked Language Modeling (MLM): 15% of words are masked, and the model predicts the original values
 - Next Sentence Prediction (NSP): the model predicts if the second sentence is connected to the first
- The [CLS] token output is used for the NSP task
- The combined loss function of MLM and NSP is optimized during training
- This non-directional approach allows BERT to develop a deeper understanding of language

BERT Architecture (Masked Language Modelling)

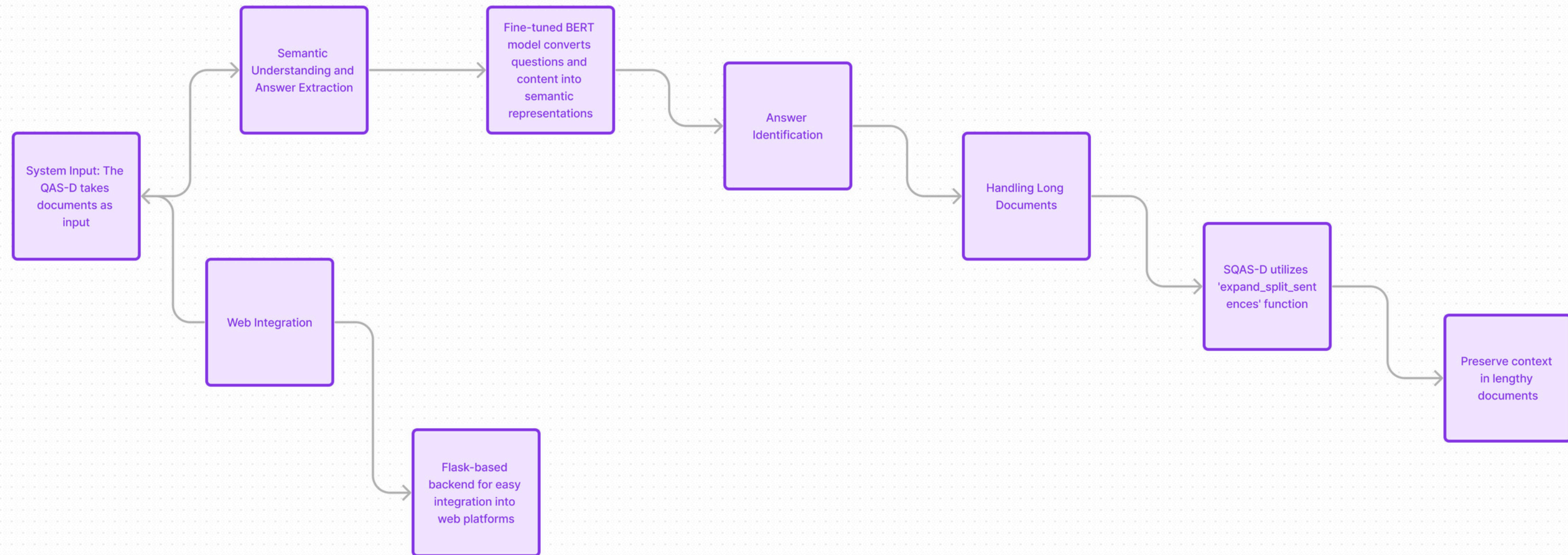


BERT Architecture (Next Sentence Prediction)

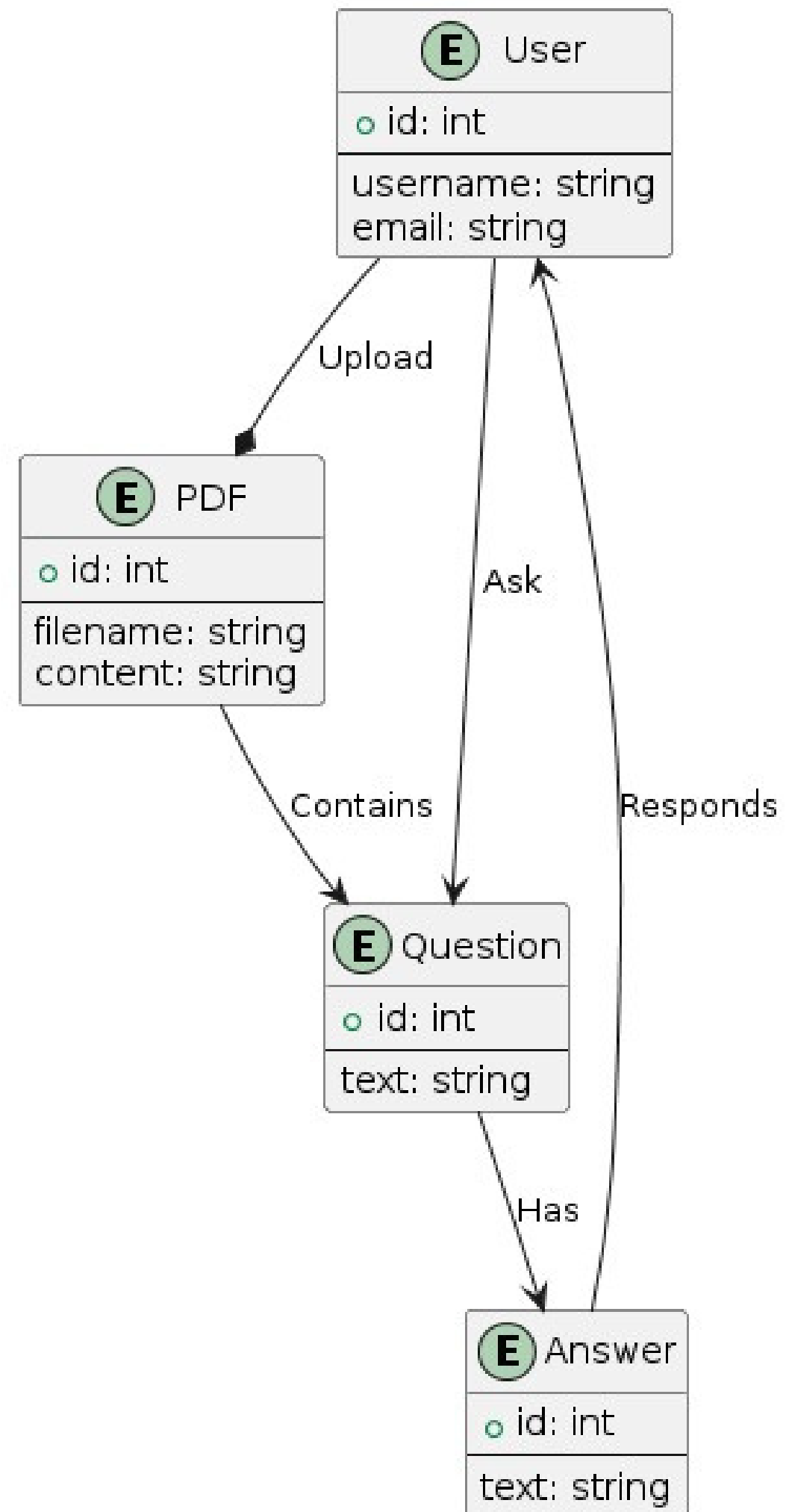


Source: [BERT](#) [Devlin et al., 2018], with modifications

SYSTEM DESIGN



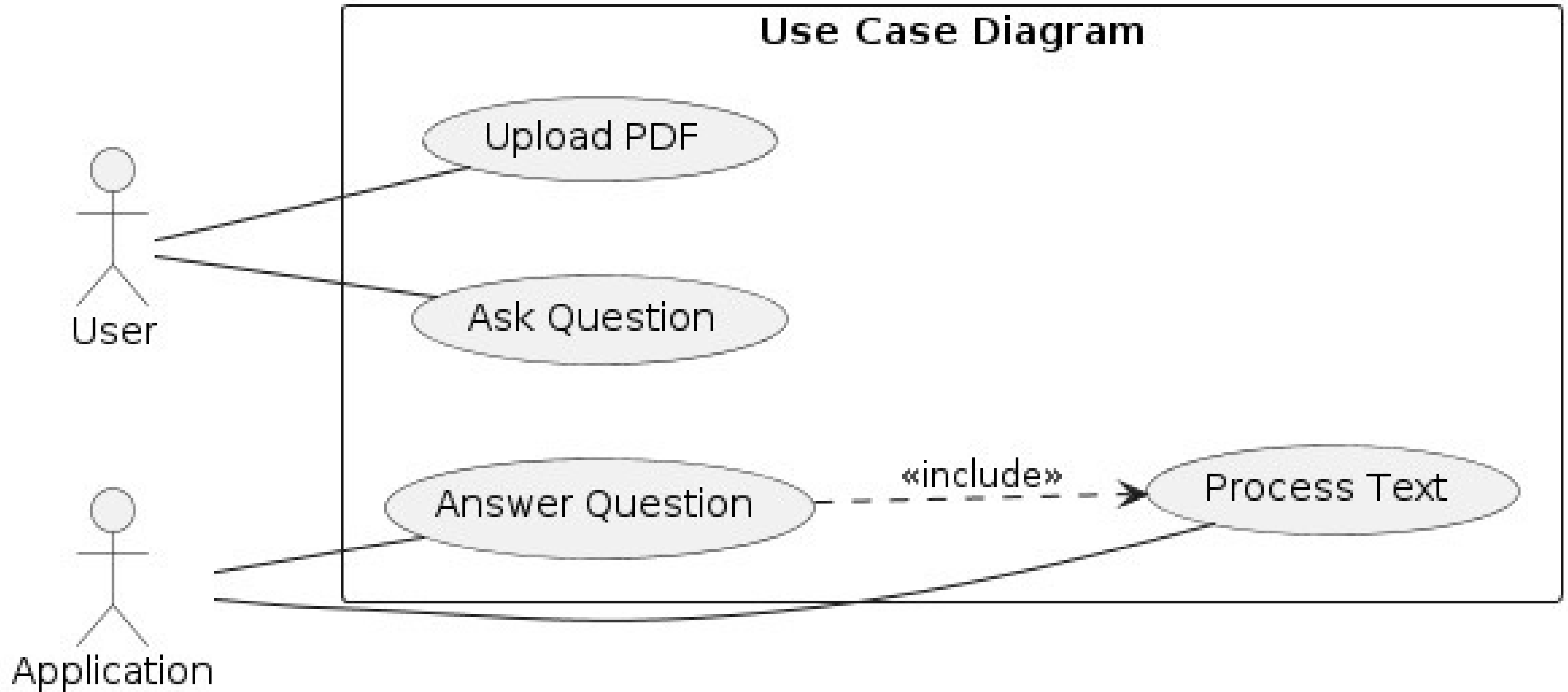
ER DIAGRAM



USE CASE DIAGRAM

- A Use Case diagram is a graphical depiction of a user's possible interactions with a system. A use case diagram shows various use cases and different types of users the system has and will often be accompanied by other types of diagrams as well.
A use case diagram doesn't describe the order in which the use cases are carried out but the behavior and the interactions between the user and the system.

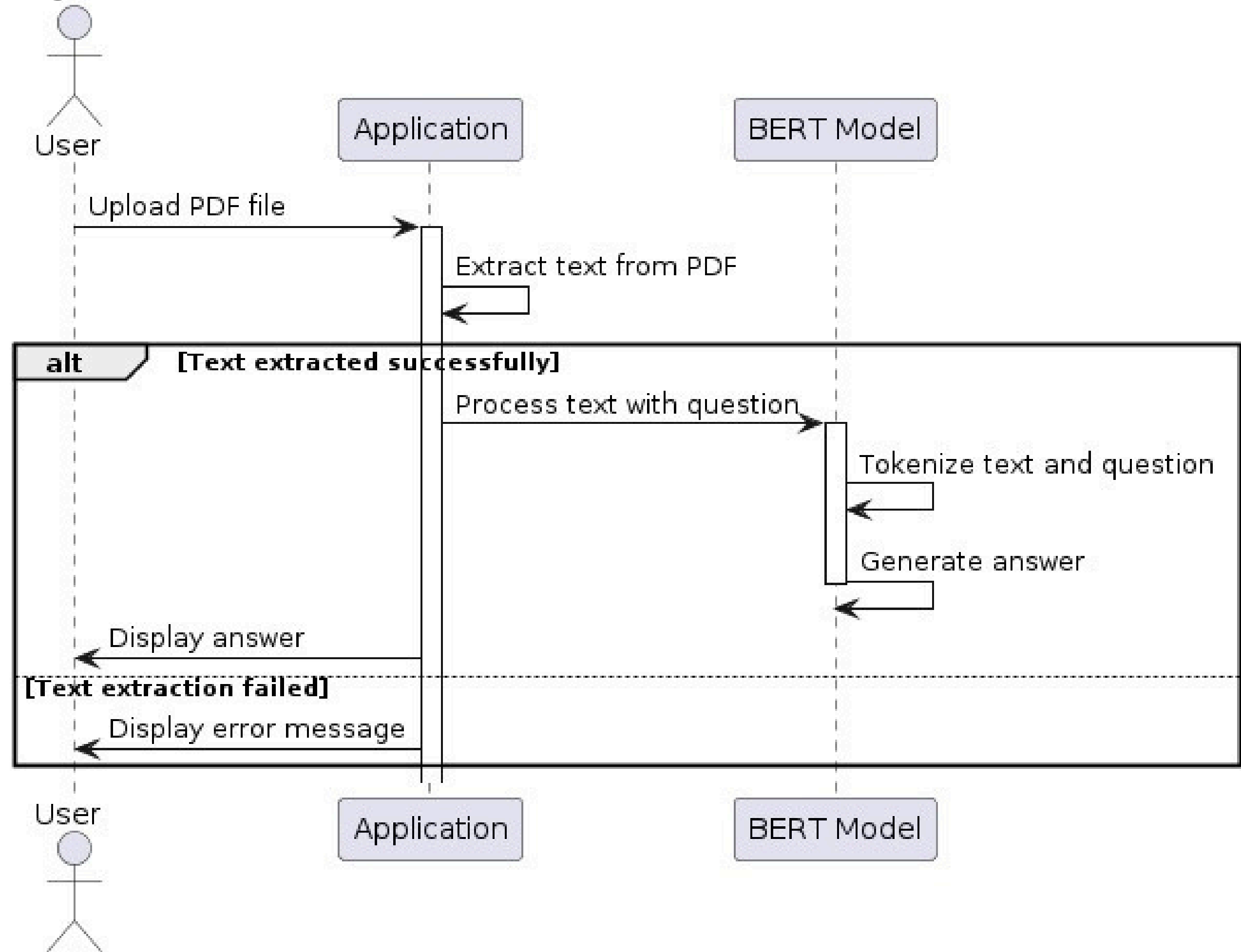
USE CASE DIAGRAM



SEQUENCE DIAGRAM

A sequence diagram is a visual representation that illustrates the interactions and order of messages exchanged among objects or components in a system or process. It is commonly used in system design to depict the flow of actions or events over time.

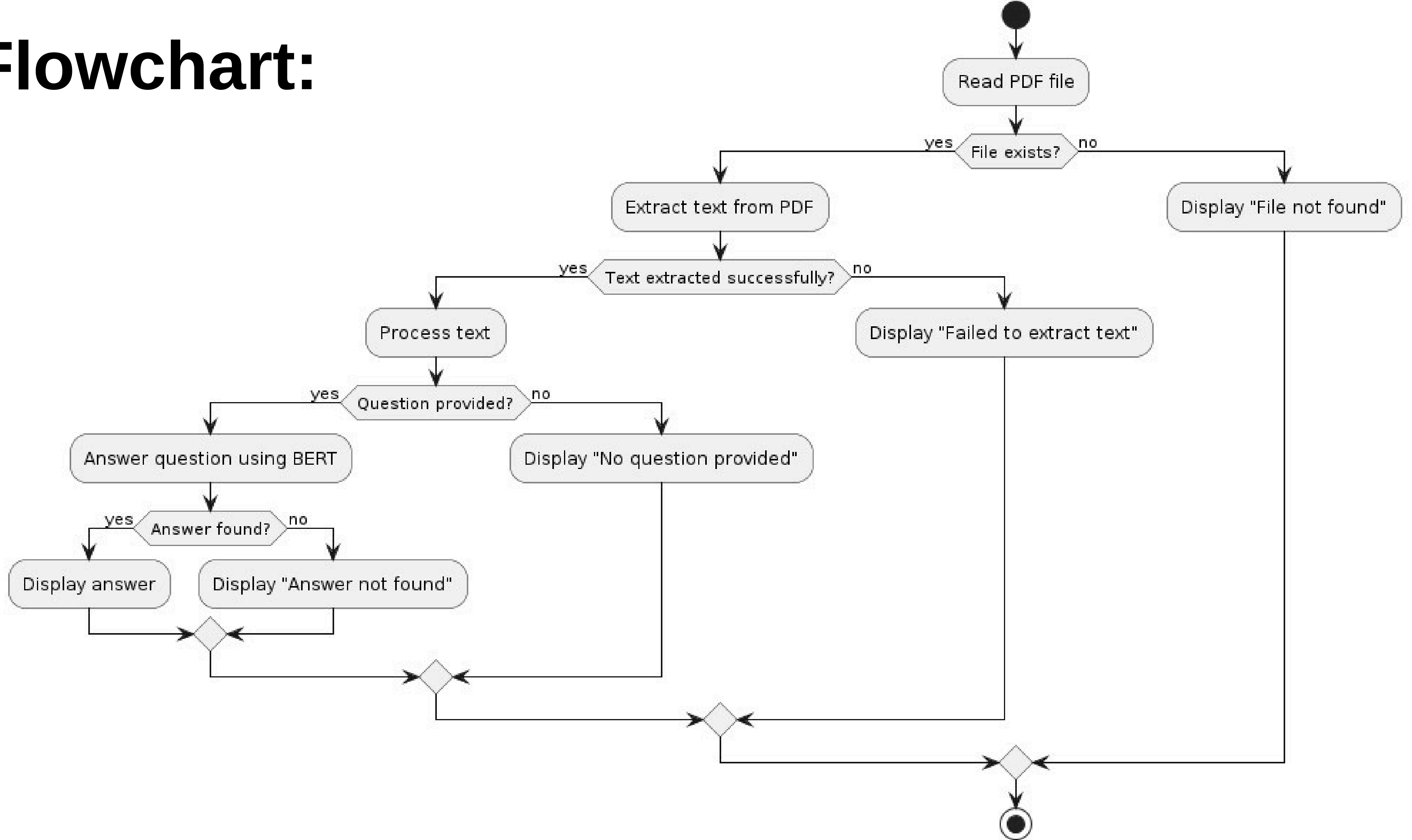
SEQUENCE DIAGRAM



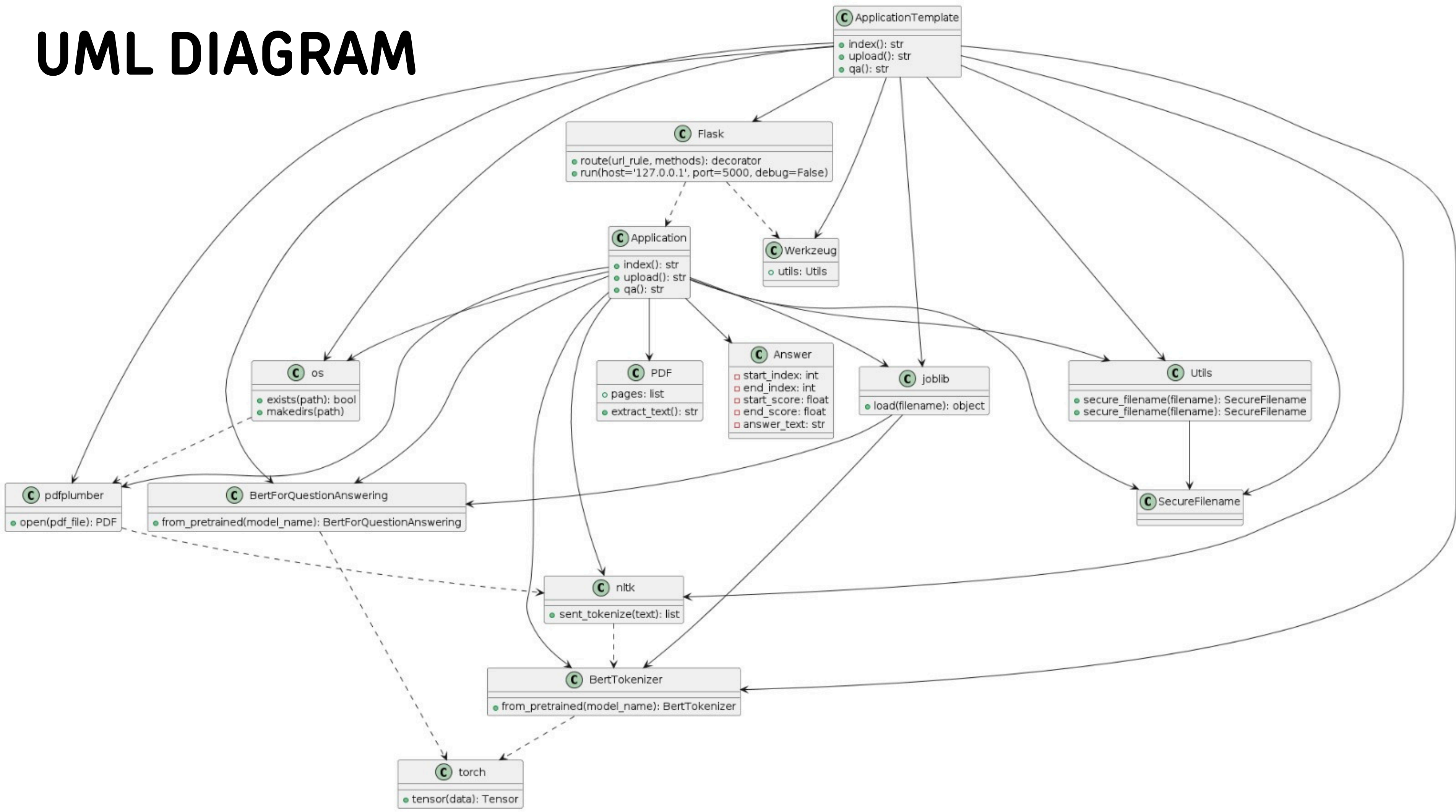
FLOWCHART

- It is a type of diagram that represents the flow of process. It can be a pictorial or diagrammatical representation of an algorithm or a system or progression through a procedure. It is represented using multiple conventional symbols and lines that connect the symbols.

Flowchart:



UML DIAGRAM



Software Requirements :

1. Operating System: Compatible with Windows, macOS, and Linux.
2. Python: Requires Python 3.7 or later.
3. Libraries: Utilizes Flask, pandas, transformers, numpy, nltk, and other necessary libraries.
4. Flask: Framework for web integration.
5. Internet Connectivity: Needed for downloading additional resources.
6. Development Environment: Supports various IDEs like Visual Studio Code.
7. Benchmark Datasets: Required for performance evaluation.
8. Dependencies: Install any additional dependencies as needed.

Hardware Requirements :

- Processor (CPU): A multi-core processor with decent processing power is recommended for handling text processing tasks efficiently.
- Memory (RAM): At least 4GB of RAM is recommended for smooth execution, especially when working with large datasets or running complex summarization models.
- Storage: Sufficient disk space to store the application code, libraries, and any generated data. This requirement can vary depending on the size of the dataset and the models used.

Sample Code :

The image shows a VS Code editor interface with a dark theme. The Explorer sidebar on the left displays the project structure for 'MAJOR-PROJECT', including folders like '.vercel', 'docs', 'Flask', and 'Flask-QA-System'. The 'Flask-QA-System' folder is expanded, showing files like 'Abstract.pdf', 'Amazon-Q4-2019-...', 'AMZN-Q1-2020-E...', 'AMZN-Q3-2020-E...', 'fdi.pdf', 'Q2-2020-Amazon-...', 'Q3-2019-Amazon-...', 'Resume.pdf', 'templates', 'bert_model.pkl', 'bert_tokenizer...', 'Flask_Deploy...', 'IMG-20240514...', 'node_modules', and '.gitignore'. The main editor area shows the 'app.py' file, which is a Flask application. The code includes imports for 'os', 'pdfplumber', 'numpy', 'torch', 'transformers', 'flask', 'werkzeug', and 'joblib'. It defines a 'pdf_extract' function that reads a PDF file and extracts its text, and a 'bert_drive' function that uses a BERT model to answer a question based on the extracted text. The application is configured to run on 'http://127.0.0.1:5000/'.

```
1  # -*- coding: utf-8 -*-
2
3  import os
4
5  if not os.path.exists('docs'):
6      os.makedirs('docs')
7  import pdfplumber
8  import numpy as np
9  import torch
10 from transformers import BertForQuestionAnswering, BertTokenizer
11 from flask import Flask, request, render_template, redirect, url_for
12 from werkzeug.utils import secure_filename
13 import joblib
14
15 app = Flask(__name__)
16
17 def pdf_extract(file_name):
18     pdf_txt = ""
19     try:
20         with pdfplumber.open(os.path.join("docs", file_name)) as pdf:
21             for pdf_page in pdf.pages:
22                 single_page_text = pdf_page.extract_text()
23                 pdf_txt += single_page_text
24     except FileNotFoundError:
25         # Handle case where the file is not found
26         return "File not found: {}".format(file_name)
27     return pdf_txt
28
29 def bert_drive(file_name, question):
30     text = pdf_extract(file_name)
31     if text:
32         max_score = 0
33         final_answer = ""
34         new_df = expand_split_sentences(text)
35
36         for new_context in new_df:
37             answer = get_answer(question, new_context)
38             score = answer['score']
39             if score > max_score:
40                 max_score = score
41                 final_answer = answer['answer']
42
43     return final_answer
44
45 else:
46     return "Sorry, couldn't retrieve the PDF text."
```

Sample Code :

```
File Edit Selection View Go Run Terminal Help Major-Project

EXPLORER
MAJOR-PROJECT
> .vercel
> docs
> Flask
Flask-QA-System
  docs
    Abstract.pdf
    Amazon-Q4-2019-...
    AMZN-Q1-2020-E...
    AMZN-Q3-2020-E...
    fdi.pdf
    Q2-2020-Amazon-...
    Q3-2019-Amazon-...
    Resume.pdf
  templates
    bert_model.pkl
    bert_tokenizer....
    Flask_Deploy...
    IMG-20240514...
    IMG-20240514...
    IMG-20240514...
    IMG-20240514...
    IMG-20240514...
    IMG-20240514...
    node_modules
    .gitignore
    app.py
    Base-Paper.pdf
    Basic-Architecture.jpg
    credentials.json
    LICENSE
    MLP-BERT.png
    Next-Sentence-Predic...
    package-lock.json
    package.json
    Pipfile
    README.md
  OUTLINE
  OPEN EDITORS
  TIMELINE

app.py
46
47 def expand_split_sentences(pdf_txt):
48     import nltk
49     nltk.download('punkt', quiet=True)
50     new_chunks = nltk.sent_tokenize(pdf_txt)
51     new_df = []
52     for i in range(len(new_chunks)):
53         paragraph = ""
54         for j in range(i, len(new_chunks)):
55             tmp_token = tokenizer.encode(paragraph + new_chunks[j])
56             if len(tmp_token) < 510:
57                 paragraph += new_chunks[j]
58             else:
59                 break
60         new_df.append(paragraph)
61     return new_df
62
63 # BERT-related setup
64 # model = BertForQuestionAnswering.from_pretrained('bert-large-uncased-whole-word-masking-finetuned-squad')
65 # tokenizer = BertTokenizer.from_pretrained('bert-large-uncased-whole-word-masking-finetuned-squad')
66
67 model = joblib.load('bert_model.pkl')
68 tokenizer = joblib.load('bert_tokenizer.pkl')
69
70
71 def get_answer(question, context):
72     return model(question=question, context=context)
73
74 #Tokenize input question and passage
75 #Add special tokens - [CLS] and [SEP]
76 input_ids = tokenizer.encode (question, context, max_length= max_len, truncation=True)
77
78
79 #Getting number of tokens in question and context passage that contains the answer
80 sep_index = input_ids.index(102)
81 len_question = sep_index + 1
82 len_context = len(input_ids)- len_question
83
84
85 #Separate question and context
86 #Segment ids will be 0 for question and 1 for context
87 segment_ids = [0]*len_question + [1]*(len_context)
88
89 #Converting token ids to tokens
90 tokens = tokenizer.convert_ids_to_tokens(input_ids)
```

Sample Code :

The image shows a Visual Studio Code editor interface with a dark theme. The Explorer sidebar on the left displays the project structure for 'MAJOR-PROJECT', including files like 'app.py', 'Flask_Deployment_Question_Answering_System.ipynb', and various PDFs and images. The main editor area shows the code for 'app.py'.

```
127
130 @app.route('/', methods=['GET', 'POST'])
131 def index():
132     if request.method == 'POST':
133         if request.form.get('btn') == 'index':
134             if 'upload' not in request.files:
135                 return "No file selected!"
136             upload = request.files['upload']
137             if upload.filename == '':
138                 return "No file selected!"
139             file_name = secure_filename(upload.filename)
140             upload.save(os.path.join("docs", file_name))
141             return redirect(url_for('qa', file_name=file_name))
142         elif request.form.get('btn') == 'qa':
143             question = request.form.get('question')
144             file_name = request.form.get('file_name')
145             if file_name is None:
146                 return "No file selected!"
147             answer, s_scores, e_scores, tokens = bert_drive(file_name, question)
148             return render_template('qa.html', answer=answer, question=question, file_name=file_name)
149         return render_template('index.html')
150
151
152 @app.route('/upload/', methods=['GET', 'POST'])
153 def upload():
154     return render_template('upload.html')
155
156 @app.route('/qa/', methods=['GET', 'POST'])
157 def qa():
158     file_name = request.args.get('file_name')
159     if file_name is None:
160         return "No file selected"
161     file_names = [f for f in os.listdir("docs")]
162     return render_template('qa.html', file_names=file_names, file_name=file_name)
163
164 if __name__ == '__main__':
165     app.run(debug=False)
```

CONCLUSION

- QAS-D presents a groundbreaking solution, adept at handling complex documents with finesse. Leveraging sophisticated semantic understanding and seamless web integration, it ensures efficient information retrieval. Users can swiftly navigate extensive texts, marking a substantial leap in question answering systems. QAS-D stands as a testament to advancements in natural language processing, empowering users with accurate, rapid access to knowledge within vast textual resources.

FUTURE SCOPE

- The QAS-D project holds immense potential for future development and expansion:
- Domain-Specific Adaptation: Customizing models for specific domains like finance or healthcare would enhance relevance and accuracy. Fine-tuning on domain-specific data would deepen understanding.
- Real-Time Capabilities: Implementing real-time processing for up-to-date answers. Integration with news APIs or RSS feeds could provide timely responses.
- User Feedback Integration: Incorporating user feedback mechanisms for continuous model improvement. Adjusting based on user assessments ensures relevance and performance enhancement over time.

REFERENCES

1. Original Research Papers:

- Luhn, H. P. (1958). "The Automatic Creation of Literature Abstracts."
- Edmundson, H. P. (1969). "New Methods in Automatic Extracting."

2. Open-Source Libraries and Frameworks:

- Streamlit Documentation (2023)
- NLTK Documentation (2023)
- Gensim Documentation (2023)
- Hugging Face Transformers Documentation (2023)

3. Official Documentation:

- Streamlit Official Documentation (2023)
- NLTK Official Documentation (2023)
- Gensim Official Documentation (2023)
- Hugging Face Transformers Official Documentation (2023)

4. Research Papers on Evaluation Metrics:

- Lin, C. (2004). "ROUGE: A Package for Automatic Evaluation of Summaries."

THANK YOU..