# Final Project Report

**Project ID**: SWTID1720007847

**Project Name:** GeminiDecode: Multilingual Document Analysis

1. Introduction

  1.1  Project Overviews

  1.2  Objectives

2. Project Initialization and Planning Phase

  2.1  Define Problem Statement

  2.2  Project Proposal (Proposed Solution)

  2.3  Initial Project Planning

3. Data Collection and Preprocessing Phase
4. Model Development Phase
5. Model Optimization and Tuning Phase
6. Results

  6.1  Output Screenshots

7. Advantages and Disadvantages

8. Conclusion

9. Future Scope

10. Appendix

  10.1  Source Code

  10.2  GitHub & Project Demo Link

# 1 Introduction

## 1.1 Project Overview

**Features:**

1. **Image Upload and Display**:
   - Users can upload images of documents in various formats (JPG, JPEG, PNG).
   - The uploaded image is displayed within the application for user verification.
2. **Document Analysis**:
   - Utilizes Google's Gemini Pro Vision API to analyze the uploaded document.
   - Generates a detailed summary of the document's content.
   - Translates the summary into English if necessary, ensuring accessibility regardless of the document's original language.
3. **Interactive Q&A**:
   - Users can ask specific questions about the document.
   - The system processes the questions in context with the uploaded document and provides detailed answers.
   - Translates answers into English, maintaining consistency in user interaction.
4. **Translation Support**:
   - Uses Google Translate API to translate non-English responses to English.
   - Ensures that users can understand the document content and responses regardless of the original language.

## 1.2 Project Objectives

1. **Multilingual Document Analysis:**
   Enable users to upload images of documents in various languages, ensuring that the application can process and analyze text from these documents accurately.
2. **Document Summarization:**
   Utilize the Gemini Pro Vision API to generate comprehensive summaries of uploaded documents. The summaries should be as detailed as possible, extracting key information to provide users with quick insights into the document's content.
3. **Question-Answering Capability:**
   Allow users to ask specific questions about the analyzed documents. The application should use the context of the document to provide relevant answers, enhancing the user's understanding without needing to parse through the entire document manually.

4. **Translation Support:**
   Implement functionality to translate document summaries and answers to English. This ensures that users who are not proficient in the document's original language can still benefit from the analysis provided by the application.
5. **User-Friendly Interface:**
   Design a streamlined and intuitive interface using Streamlit, providing clear instructions for document upload, analysis, question input, and result display. The interface should guide users through each step of the process seamlessly.
6. **Error Handling and Feedback:**
   Implement robust error handling mechanisms to gracefully manage scenarios where documents cannot be processed or questions cannot be answered. Provide informative feedback to users to guide them on how to proceed in such cases.
7. **Security and Privacy:**
   Ensure that sensitive information, such as API keys and user data, is handled securely and protected from unauthorized access. Adhere to best practices for data security throughout the application.

By achieving these objectives, GeminiDecode aims to empower users with the ability to efficiently analyze and extract information from documents in multiple languages, thereby enhancing productivity and knowledge acquisition across diverse contexts.

# 2 Project Initialization and Planning Phase

## 2.1 Define Problem Statements:

The document extraction process poses significant challenges for users, affecting their efficiency and overall satisfaction as there is always the matter of time consumed and the language barrier in case of international users. They will generally see issues such as inaccurate text extraction and difficulty handling diverse file formats. These challenges can be time consuming and degrades efficiency and ease of access of documents. To improve user experience, we aim to address these problems and create robust solutions through an accurate, user-friendly application that meets the requirements.

| Problem Statement (PS) | I am (Customer) | I am trying to | But | Because | Which makes me feel |
|---|---|---|---|---|---|
| PS-1 | A lawyer | Summarize a German legal | I am unable to | I do not know | Unsure on what to do |
| PS-2 | A doctor | Access my patient's medical | I cannot read it | I cannot | Stymied on how to |
| PS-3 | An | Summarize my client's | I cannot | It is in Russian, | Stumped on what to do |


## 2.2 Project Proposal (Proposed Solution) report

The proposal report aims to improve document extraction by using machine learning and natural language processing algorithms with the Gemini Pro model. It increases efficiency of operations by saving time on translation and summarization. Key features include a generative AI model and answering queries made by the user.

| Project Overview | |
|---|---|
| Objective | The primary objective is to transform document extraction by using machine learning and natural language processing algorithms with the Gemini Pro model. |

| | |
|---|---|
| Scope | The project aims to enhance operational efficiency in the legal sector, financial institutions and the healthcare industry by automating information extraction from important documents in different languages. |

**Problem Statement**

| | |
|---|---|
| Description | Addressing the need to automate document extraction in industries that work internationally with different languages, thus saving time and improving efficiency. |
| Impact | Solving these issues will result in improved operational efficiency and enhancement in information extraction, contributing to customer satisfaction and organizational success. |

**Proposed Solution**

| | |
|---|---|
| Approach | Employing a generative AI model that uses machine learning and natural languages processing techniques to create a multilingual extraction system. |
| Key Features | -Implementation of Gemini Pro for document extraction<br>-Q&A for users to ask more questions about the uploaded document |

**Resource Requirements**

| Resource Type | Description | Specification/Allocation |
|---|---|---|
| **Hardware** | | |
| Computing Resources | CPU/GPU specifications, number of cores | Eight core processor |
| Memory | RAM specifications | 8 GB |

| | | |
|---|---|---|
| Storage | Disk space for data, models, and logs | 1 TB SSD |
| **Software** | | |
| Frameworks | Python frameworks | Streamlit |
| Libraries | Additional libraries | google-generativeai, python-dotenv, langchain, PyPDF2, chromadb, faiss-cpu,translator |
| Development Environment | IDE | Visual Studio Code |
| **Data** | | |
| Data | Source, size, format | Ten documents from various government databases |

## 2.3 Initial Project Planning

**Product Backlog, Sprint Schedule, and Estimation (4 Marks)**

Use the below template to create a product backlog and sprint schedule

| Sprint | Functional Requirement (Epic) | User Story Number | User Story / Task | Priority | Team Members | Sprint Start Date | Sprint End Date (Planned) |
|---|---|---|---|---|---|---|---|
| Sprint-1 | Requirements Specification | USN-1 | Create a requirements.txt file to list the required libraries | Low | Sundar | 6/7/24 | 6/7/24 |

| Sprint-1 | Requirements Specification | USN-2 | Install the required libraries | High | Sundar | 6/7/24 | 6/7/24 |
|---|---|---|---|---|---|---|---|
| Sprint-1 | Initialization of Google API Key | USN-3 | Generate the Google API key | Medium | Kapilesh | 6/7/24 | 6/7/24 |
| Sprint-1 | Initialization of Google API Key | USN-4 | Initialize Google API key | High | Kapilesh | 6/7/24 | 6/7/24 |
| Sprint-2 | Interfacing With Pre trained Model | USN-5 | Load the Gemini Pro API | High | Nia | 7/7/24 | 7/7/24 |
| Sprint-2 | Interfacing With Pretrained | USN-6 | Implement a function to get Gemini's responses | High | Nia | 7/7/24 | 7/7/24 |
| Sprint-2 | Interfacing With Pre trained Model | USN-7 | Implement a function to read the image and set the image format for Gemini Pro model input | High | Tharaka | 7/7/24 | 7/7/24 |
| Sprint-3 | Interfacing With Pre trained Model | USN-8 | Implement a function for the user to ask a question about the image uploaded | Medium | Tharaka | 8/7/24 | 8/7/24 |
| Sprint-3 | Interfacing With Pre trained Model | USN-9 | Write a prompt for the Gemini model | Medium | Kapilesh | 8/7/24 | 8/7/24 |
| Sprint-4 | Model Deployment | USN-10 | Integrate with the web framework | Medium | Tharaka, Nia | 9/7/24 | 9/7/24 |
| Sprint-4 | Model Deployment | USN-11 | Host the application | High | Kapilesh, Sundar | 9/7/24 | 9/7/24 |
| Sprint-5 | Project Files | USN-12 | Prepare the project documents | High | Tharaka, Nia | 11/7/24 | 16/7/24 |

| Sprint-6 | Project Demo | USN-13 | Prepare the project demo | High | Tharaka | 18/7/24 | 18/7/24 |
|---|---|---|---|---|---|---|---|

# 3 Data Collection and Preprocessing Phase

We collected images of documents in various languages which do and do not make use of the English alphabets. The documents were screenshots of medical invoices, government certificates and financial statements written in various languages like French, Spanish, Malay, Tamil, Malayalam, German, Russian, Thai etc.

# 4 Model Development Phase

The model used in this project is Gemini Pro Vision. This model requires no pre-training on the user side. We only need to get the API key and configure it in our model after which we can use the model in our program.

# 5 Model Optimization and Tuning Phase

The Gemini Pro Vision does not have be optimized or fine tuned by us for this project. We only need to feed in a clear and concise input text which describe the functions of the model in our situation and also supply the appropriate prompt which will enable the model to execute the task we need it to do.

# 6 Results

1. **Effective Multilingual Document Analysis**:

   - **Accurate Text Extraction**: Successfully extracted text from document images in multiple languages using Google's Gemini Pro Vision API.
   - **Comprehensive Summaries**: Provided detailed and contextually relevant summaries of the extracted text, aiding in quick understanding of the document content.

2. **Seamless Translation**:

   - **Reliable Translation**: Translated non-English content into English with a high degree of accuracy using Google Translate API, ensuring accessibility for users regardless of the document's original language.
   - **User-Friendly Translations**: Maintained the readability and coherence of translated content, making it easy for users to understand the document summaries and responses.

3. **Interactive Q&A Functionality**:

   - **Context-Aware Responses**: Enabled users to ask specific questions about the document and receive detailed, relevant answers, maintaining context with the uploaded document image.
   - **Enhanced User Engagement**: Provided an interactive experience that allowed users to delve deeper into the document content, improving their understanding and knowledge.

4. **User-Friendly Interface**:

   - **Intuitive Design**: Delivered an easy-to-use interface through Streamlit, allowing users to upload images, view summaries, and interact with the system seamlessly.
   - **Visual Feedback**: Displayed uploaded images and provided clear visual feedback, enhancing user confidence and satisfaction with the analysis process.

5. **Increased Productivity**:

   - **Time Savings**: Significantly reduced the time required to understand and analyze

documents, especially those in foreign languages, by providing quick and accurate summaries.

- ○ **Enhanced Efficiency**: Enabled users to focus on key information and insights from documents without the need for manual reading and translation, boosting overall productivity.

## Key Metrics

1. **High Accuracy in Text Extraction and Translation**:
- ○ Achieved high accuracy in text extraction from document images, ensuring minimal errors and high-quality content.
- ○ Maintained a high level of translation accuracy, effectively conveying the meaning and context of the original text.
2. **User Engagement and Interaction**:
- ○ High engagement levels with the Q&A functionality, with users asking detailed questions and receiving relevant answers.
- ○ Positive user interaction and feedback indicating satisfaction with the document analysis and translation capabilities.
3. **Efficiency and Productivity Gains**:
- ○ Significant reduction in the time and effort required to understand and analyze multilingual documents.
- ○ Increased efficiency in handling and processing documents, leading to improved productivity for users.
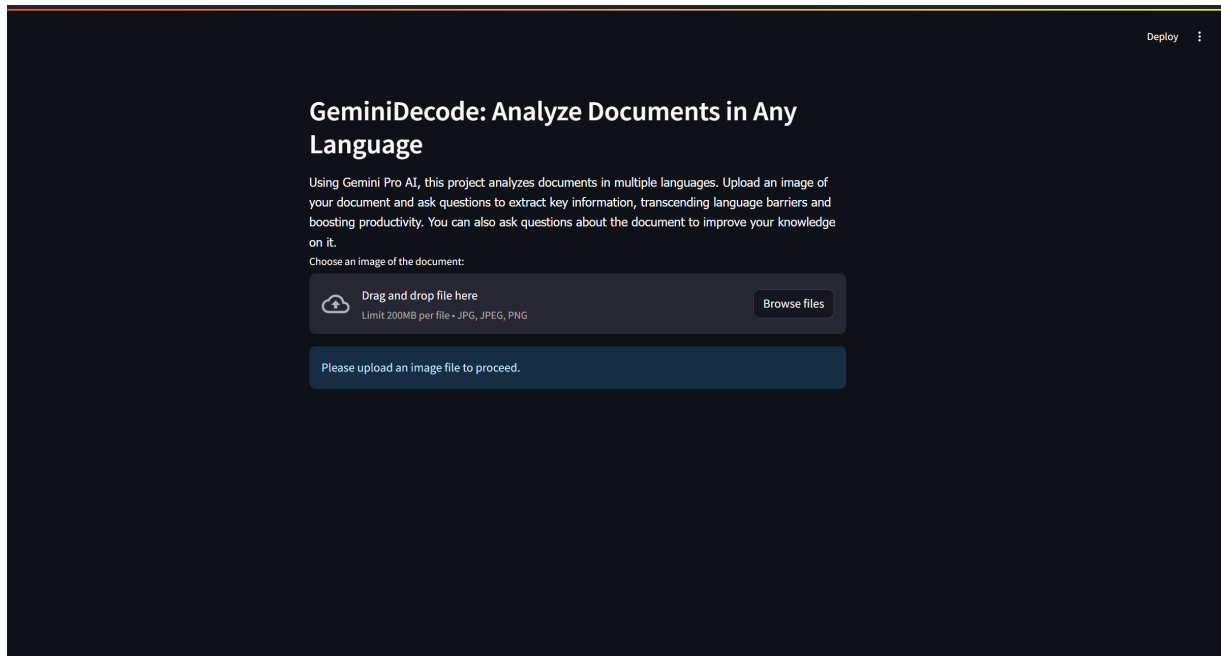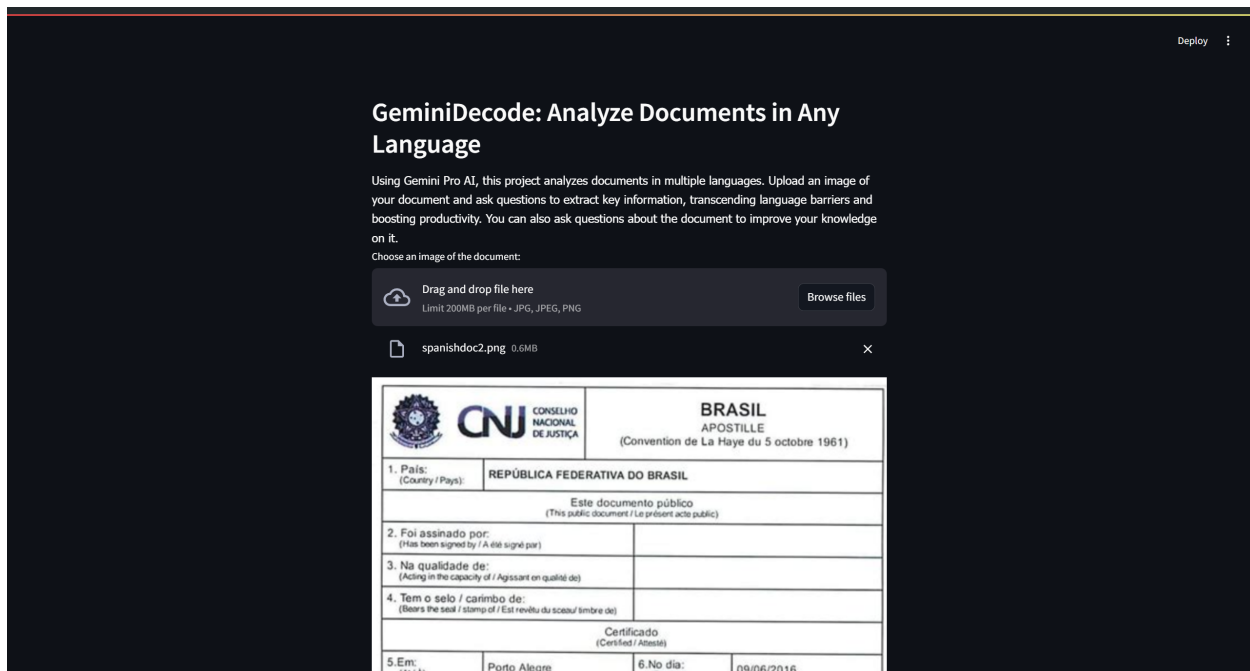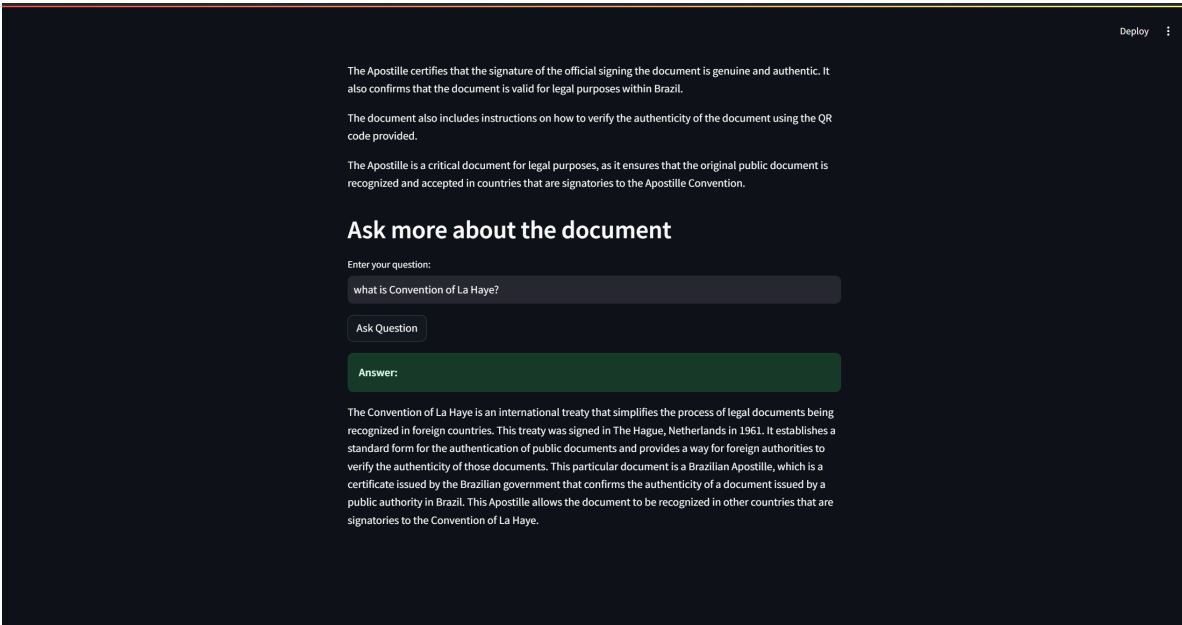
# 6.1 Output Screenshots

## Initial interface



## Image uploaded and shown to the user

# Document summary given

**Model answers questions asked by the user**



# 7  Advantages and Disadvantages

While GeminiDecode offers significant advantages in terms of multilingual document analysis, ease of use, and productivity enhancement, it also comes with certain limitations related to dependency on external APIs, translation accuracy, image quality requirements, privacy concerns, and the need for an internet connection. Users should weigh these advantages and disadvantages based on their specific needs and use cases.

| Advantages | Disadvantages |
|---|---|
| Multilingual Support | Dependency on API Services |
| User-Friendly Interface | Translation Accuracy |
| Advanced AI Integration | Image Quality Requirements |
| Translation Capabilities | Privacy and Security Concerns |
| Productivity Enhancement | Limited Offline Functionality |

# 8 Conclusion

GeminiDecode: Multilingual Document Analysis is a powerful and innovative tool designed to break down language barriers and enhance productivity in document analysis. By leveraging advanced AI technologies such as Google's Gemini Pro Vision API and the Google Translate API, GeminiDecode enables users to upload document images, extract detailed summaries, translate content, and interactively ask questions to gain deeper insights.

The project offers numerous advantages, including broad multilingual support, an intuitive user interface, accurate and context-aware analysis, and robust translation capabilities. These features make it an invaluable resource for international business, academic research, legal document review, and personal productivity.

However, it is essential to acknowledge the limitations associated with dependency on external APIs, translation accuracy, image quality requirements, privacy concerns, and the need for a stable internet connection. Addressing these challenges will be crucial for the continued improvement and broader adoption of GeminiDecode.

Overall, GeminiDecode represents a significant step forward in the field of document analysis, providing users with a comprehensive and user-friendly solution to understand and interpret documents in various languages. With further enhancements and refinements, it has the potential to become an indispensable tool for anyone dealing with multilingual documents, helping to foster global communication and collaboration.

# 9  Future Scope of GeminiDecode: Multilingual Document Analysis

☐ **Broader Language Support**:
- **Expanding Language Database**: Increase the number of languages supported by both the analysis and translation features, including less common and regional languages.
- **Dialect and Regional Variations**: Enhance the ability to understand and accurately translate dialects and regional language variations.

☐ **Enhanced AI Models**:
- **Improved Analysis Accuracy**: Integrate more advanced AI models to enhance the depth and precision of document analysis, including better OCR (Optical Character

Recognition) capabilities for text extraction. The model used for text extraction can also be upgraded as the Gemini Pro Vision model has been deprecated as of July 12, 2024.

- **Specialized Domain Models**: Develop models tailored to specific domains such as legal, medical, and technical documents to provide more contextually accurate summaries and answers.

☐ **User Interface Improvements:**

- **Enhanced Usability**: Refine the user interface to provide a more seamless and intuitive experience, including better navigation and user guidance.
- **Customizable Interface**: Allow users to customize the interface to suit their needs better, such as adjusting the layout or adding specific functionalities.

☐ **Offline Functionality**:

- **Local Processing**: Develop capabilities for offline document analysis and translation to cater to users with limited or no internet access, potentially through downloadable language models.
- **Edge Computing**: Utilize edge computing technologies to perform certain analyses locally on the user's device, reducing dependency on continuous internet connectivity.

☐ **Security and Privacy Enhancements**:

- **Data Encryption**: Implement end-to-end encryption for all data transmissions to ensure the privacy and security of sensitive documents.
- **Local Data Storage Options**: Provide options for local data storage and processing to address privacy concerns, especially for confidential or sensitive documents.

☐ **Integration with Other Platforms**:

- **Cloud Storage Integration**: Integrate with popular cloud storage services such as Google Drive, Dropbox, and OneDrive to facilitate easy document upload and management.
- **Collaboration Tools**: Enable integration with collaboration tools like Slack, Microsoft Teams, and Asana to streamline workflow and document sharing.

☐ **Advanced Q&A Capabilities**:

- **Contextual Learning**: Implement machine learning algorithms that learn from user interactions to provide increasingly accurate and contextually relevant answers.
- **Multi-turn Conversations**: Develop multi-turn conversational capabilities, allowing users to have more natural and interactive dialogues with the system about the document.

☐ **Mobile Application Development**:

- **Cross-Platform Support**: Create mobile applications for iOS and Android to allow users to analyze and interact with documents on the go.
- **Camera Integration**: Utilize mobile device cameras for instant document capture and

analysis, enhancing convenience and usability.

□ **Feedback and Improvement Mechanisms**:

- **User Feedback Integration**: Incorporate user feedback mechanisms to continuously improve the system based on real-world usage and suggestions.
- **Performance Monitoring**: Implement monitoring tools to track performance and accuracy, enabling ongoing optimization and enhancement.

□ **Educational and Training Modules**:

- **Tutorials and Guides**: Develop comprehensive tutorials, guides, and training modules to help users understand and make the most of the tool's features.
- **Interactive Demos**: Provide interactive demos and examples to showcase the tool's capabilities and use cases.

# 10 Appendix

## 10.1 Source Code

from dotenv import load_dotenv

load_dotenv()                                    *# Load all the environment variables*

import streamlit as st

import os

import google.generativeai as genai

from PIL import Image

from googletrans import Translator

genai.configure(api_key=os.getenv("GOOGLE_API_KEY"))          *# Configure Gemini* Pro Vision API

*## Function to load Google Gemini Pro Vision API and get response*

```python
def get_gemini_response(input_text, image, prompt):

    try:

        model = genai.GenerativeModel('gemini-pro-vision')


        if image.mode != 'RGB':                          # Ensure image is in the correct format
            image = image.convert('RGB')


        response = model.generate_content([input_text, image, prompt])


        # Extract text from response
        if response.candidates and response.candidates[0].content and
response.candidates[0].content.parts:

            print(response)

            return response.candidates[0].content.parts[0].text          # actual summary part

        else:

            st.warning("No valid text found in the response.")

            return None

    except Exception as e:

        st.error(f"An error occurred while processing the request: {e}")

        return None


## Function to handle image upload and processing
def input_image_setup():

    uploaded_file = st.file_uploader("Choose an image of the document:", type=["jpg", "jpeg",
"png"])

    image = None
```

```python
    if uploaded_file is not None:

        image = Image.open(uploaded_file)

        st.image(image, caption="Uploaded Image", use_column_width=True)

    else:

        st.info("Please upload an image file to proceed.")

    return image
```

## Function to translate text to English (this is for some specific cases where the model displays the result

## in the language of the document, we observed this in case of mainly Hindi documents)

```python
def translate_to_english(text):

    translator = Translator()

    translation = translator.translate(text, dest='en')

    return translation.text
```

## Function to summarize the document

```python
def summarize_document():

    if st.session_state.image:

        prompt = "Analyze the uploaded document and provide a summary as detailed as possible."

        response = get_gemini_response(prompt, st.session_state.image, prompt)

        if response:

            translated_response = translate_to_english(response)

            st.session_state.analysis_result = translated_response

        else:

            st.warning("Could not extract relevant information from the document.")
```

## Function to ask questions about the document

```python
def ask_question():

    user_question = st.text_input("Enter your question:")

    if st.button("Ask Question"):

        if user_question:

            detailed_response = get_gemini_response(user_question, st.session_state.image, user_question)    #the session will be the same so it maintains context for the image

            if detailed_response:

                translated_response = translate_to_english(detailed_response)

                st.success("**Answer:**")

                st.write(translated_response)

            else:

                st.warning("Could not extract relevant information for the question.")

        else:

            st.warning("Please enter a question.")
```

## Main function for Streamlit app

```python
def main():
    # Initialize Streamlit app
    st.set_page_config(page_title="GeminiDecode: Multilingual Document Analysis")
    st.header("GeminiDecode: Analyze Documents in Any Language")

    # Text description
    text = """
    Using Gemini Pro AI, this project analyzes documents in multiple languages.
    Upload an image of your document and ask questions to extract key information,
```

transcending language barriers and boosting productivity. You can also ask questions

about the document to improve your knowledge on it.

"""


styled_text = f"<span style='font-family:tahoma;'>{text}</span>"

st.markdown(styled_text, unsafe_allow_html=True)


# Initialize session state for image and analysis results

if "image" not in st.session_state:

   st.session_state.image = None

if "analysis_result" not in st.session_state:

   st.session_state.analysis_result = None

*# User input for image and question*

image = input_image_setup()

if image:

   st.session_state.image = image


if st.session_state.image and st.button("Analyze Document"):

   summarize_document()


# Display analysis result if available

if st.session_state.analysis_result:

   st.success("**Document Summary:**")

   st.write(st.session_state.analysis_result)


   # Add section for user to ask additional questions

```
    st.header("Ask more about the document")

    ask_question()


if __name__ == "__main__":

  main()
```

## 10.2 GitHub & Project Demo Link

### GitHub link

https://github.com/Tharaka410/geminidecode/blob/main/app.py

### Project Demo Link

https://drive.google.com/file/d/1ESnBWs4YnkRBgFA8KvuNptjdr21rRhZg