

# MINI-PROJECT REPORT ON

---

---

## *OCR*

### *Optical Character Recognition to Text And Audio*

---

---

Submitted in partial completion of the requirement for the award of the degree of

**BACHELOR OF TECHNOLOGY IN  
COMPUTER SCIENCE & ENGINEERING**

**Submitted by:**

**Student Name: Kajal Rawat**

**University Roll No.:2118668**

*Under the Mentorship of*  
**Ms. Preeti Chaudhary**  
Assistant Professor



**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**GRAPHIC ERA HILL UNIVERSITY, DEHRADUN**



## **CANDIDATE'S DECLARATION**

I hereby certify that the work which is being presented in the project report entitled “**OCR- Optical Character Recognition to Text And Audio**” in partial fulfillment of the requirements for the award of the Degree of Bachelor of Technology in Computer Science and Engineering of the Graphic Era Hill University, Dehradun shall be carried out by myself under the mentorship of **Ms. Preeti Chaudhary, Assistant Professor**, Department of Computer Science and Engineering, Graphic Era Hill University, Dehradun Campus.

**Name: Kajal Rawat**

**University Roll no.: 2118668**

# Table of Contents

---

| Chapter No. | Description       | Page No. |
|-------------|-------------------|----------|
| Chapter 1   | Introduction      | 1        |
| Chapter 2   | Literature Review | 2        |
| Chapter 3   | Methodology       | 3-5      |
| Chapter 4   | Results           | 6-9      |
| Chapter 5   | Sample Code       | 10-14    |
| Chapter 6   | Future Work       | 15       |
| Chapter 7   | Conclusion        | 16       |
| Chapter 8   | References        | 17       |

## Chapter 1

# **INTRODUCTION**

"Embracing AI and ML for Enhanced Information Accessibility" is about using smart computer technology to make information easier to get. It's like having clever tools that can read and understand things like text in pictures or documents. This technology helps take messy information and turns it into something computers can understand. By using these smart tools, it becomes easier for everyone to understand and use information, even if they speak different languages or have different ways of getting information.

Our project is dedicated to leveraging AI and ML technologies to significantly enhance information accessibility and usability across diverse user scenarios. Through 'AI in Voice Technology,' we aim to enable seamless voice-based interactions, breaking language barriers, and fostering improved information access through spoken language. This facet of our project focuses on:

- Facilitating voice-based communication and interactions
- Overcoming language barriers
- Enhancing information accessibility via spoken language

Concurrently, 'Machine Learning for Human Assistance' plays a pivotal role in our project. By harnessing technologies like Optical Character Recognition (OCR) and other ML tools, it aims to:

- Assist in efficiently processing and extracting information
- Improve data accessibility and understandability

These integrated applications align with our overarching goal of rendering information more user-friendly and accessible, benefiting a diverse range of users in their interaction with various forms of data. By leveraging AI and ML advancements, our project strives to make information universally accessible and comprehensible."

## Chapter 2

### LITERATURE REVIEW

Our project implementation involves a fusion of AI and ML technologies, notably employing Optical Character Recognition (OCR), Natural Language Processing (NLP), and Voice Technology (gTTS) to enhance information accessibility. This implementation aligns with existing research, where similar AI-driven approaches have been utilized:

- **Efficient Data Extraction:** Through OCR and image processing, our approach echoes established practices, enabling precise text extraction from diverse sources like images and documents.
- **Multilingual Accessibility:** Leveraging translation tools and gTTS aligns with research emphasizing the importance of multilingual access, catering to users who prefer content in their native language.
- **Personalized Interaction:** Implementing ML-driven personalized content mirrors findings emphasizing the significance of tailored recommendations, enhancing user engagement and satisfaction.
- **Enhanced User Experience:** Our project adheres to accessible UI design principles, catering to varying user needs and conforming to established standards.
- **Improved Efficiency:** Utilizing NLP techniques and automated summarization aligns with research advocating for quicker comprehension and navigation through large volumes of data, thereby enhancing overall efficiency in information processing.
- **Inclusivity and Fairness:** The project's emphasis on ethical AI practices, avoiding biases, and ensuring fairness in data processing resonates with established literature stressing the importance of unbiased algorithms for inclusive accessibility.

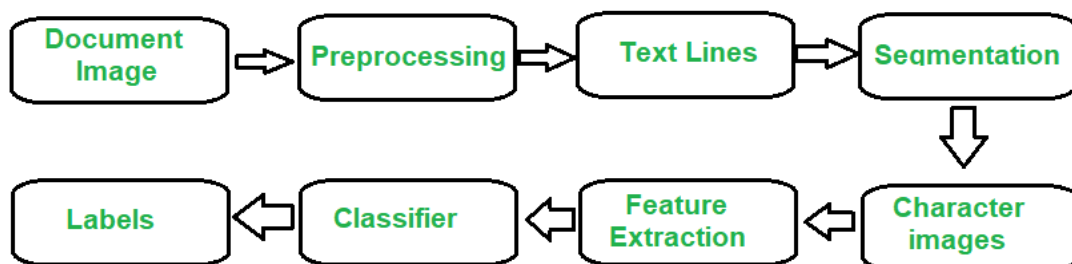
By integrating these findings into our project implementation, we aim to achieve tangible benefits such as enhanced data accessibility, improved user experiences, and efficient information processing, ultimately aligning with the broader goal of making information universally accessible and comprehensible.

## Chapter 3

### METHODOLOGY

The OCR system includes the following steps:

- **User Interaction:**
  - Using HTML, CSS, Flask.
  - Users upload PDF or image files through an intuitive interface (as shown in fig.3).
  - File types are validated for PDFs or images.
- **File Processing and Text Extraction:**
  - Receive uploaded files and save them to a designated location.
  - Check the file type (PDF or image) and proceed accordingly.
  - For PDFs:
    - Utilize PyMuPDF (fitz) to load the PDF and extract text from text-based elements on each page.
    - Extract text from images within the PDF using Pytesseract OCR.
  - For images:
    - Use OpenCV (cv2) to preprocess the image (grayscale, thresholding, noise removal).
    - Apply Pytesseract OCR to extract text from the preprocessed image



*Fig.1: Flow Diagram of OCR*

- **Translation and Audio Generation:**
  - Enable user input for text translation.
  - Utilize gTTS for language translation and convert translated text into MP4 audio.
- **Output Management:**
  - Store translated text and audio in TXT and MP4 formats, respectively.
  - Manage file storage using os module and provide download links (as shown in Fig.8).
- **Error Handling and Feedback:**
  - Implement error handling for invalid file formats and processing errors.
  - Provide user feedback on successful translations and error notifications.
- **Optimization and Testing:**
  - Continuously refine image processing and text extraction techniques.
  - Ensure accuracy and efficiency in translation and text-to-speech conversion.
- **Deployment:**
  - Test functionalities across browsers and devices.
  - Deploy the functional web application for public access.

### **PesudoCodes:**

- Below is a basic pseudocode representation of image preprocessing:

```
function preprocess(image):
    # Convert the image to grayscale
    grayscale_image = convert_to_grayscale(image)

    # Apply thresholding to the grayscale image
    thresholded_image = apply_thresholding(grayscale_image)

    # Additional preprocessing steps (if needed)
    processed_image = additional_processing(thresholded_image)

    return processed_image

function convert_to_grayscale(image):
    # Convert the input image to grayscale
    grayscale_image = convert_to_grayscale_algorithm(image)

    return grayscale_image

function apply_thresholding(image):
    # Apply thresholding technique to the image
    thresholded_image = apply_thresholding_algorithm(image)

    return thresholded_image
```

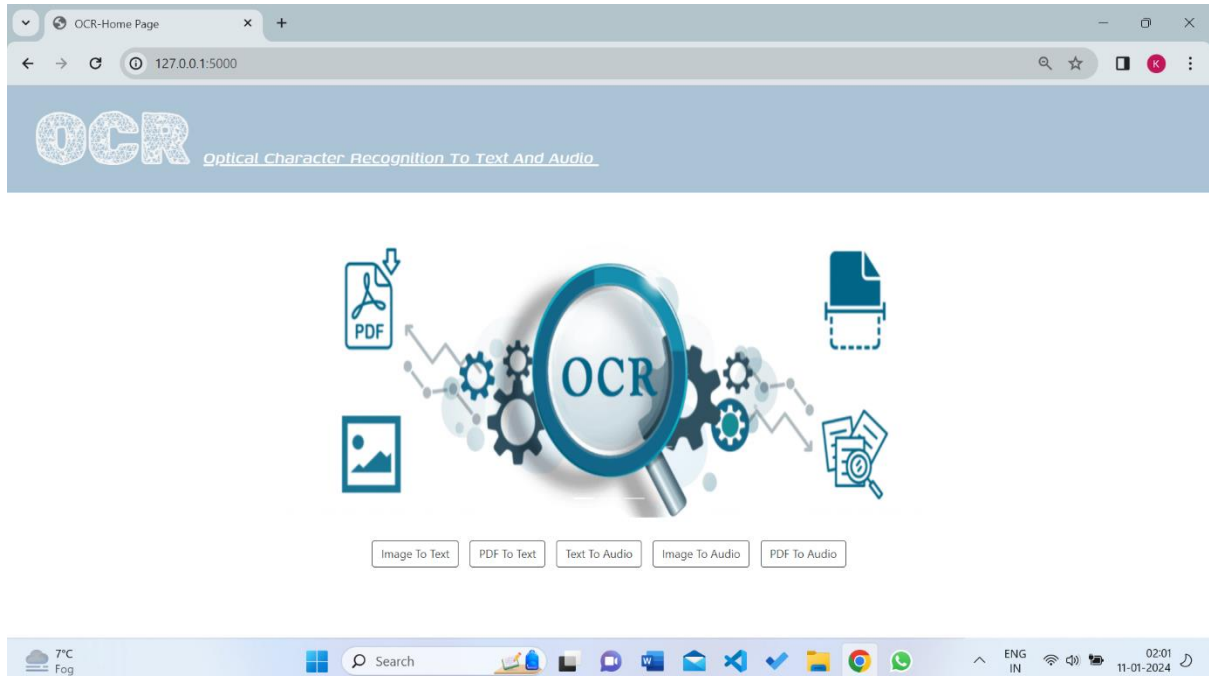
- Below is a basic pseudocode representation of an OCR algorithm:

```
function OCR_algorithm(preprocessed_image):  
    # Initialize variables to store extracted text  
    extracted_text = ""  
  
    # Define OCR processing steps  
    for each text_region in preprocessed_image:  
        # Apply OCR on each text region within the preprocessed image  
        text = apply_ocr_on_region(text_region)  
  
        # Append extracted text from the region to the overall text  
        extracted_text += text + " "  
  
    return extracted_text  
  
function apply_ocr_on_region(text_region):  
    # Call OCR function to extract text from the given text region  
    extracted_text = OCR_function(text_region)  
  
    return extracted_text
```

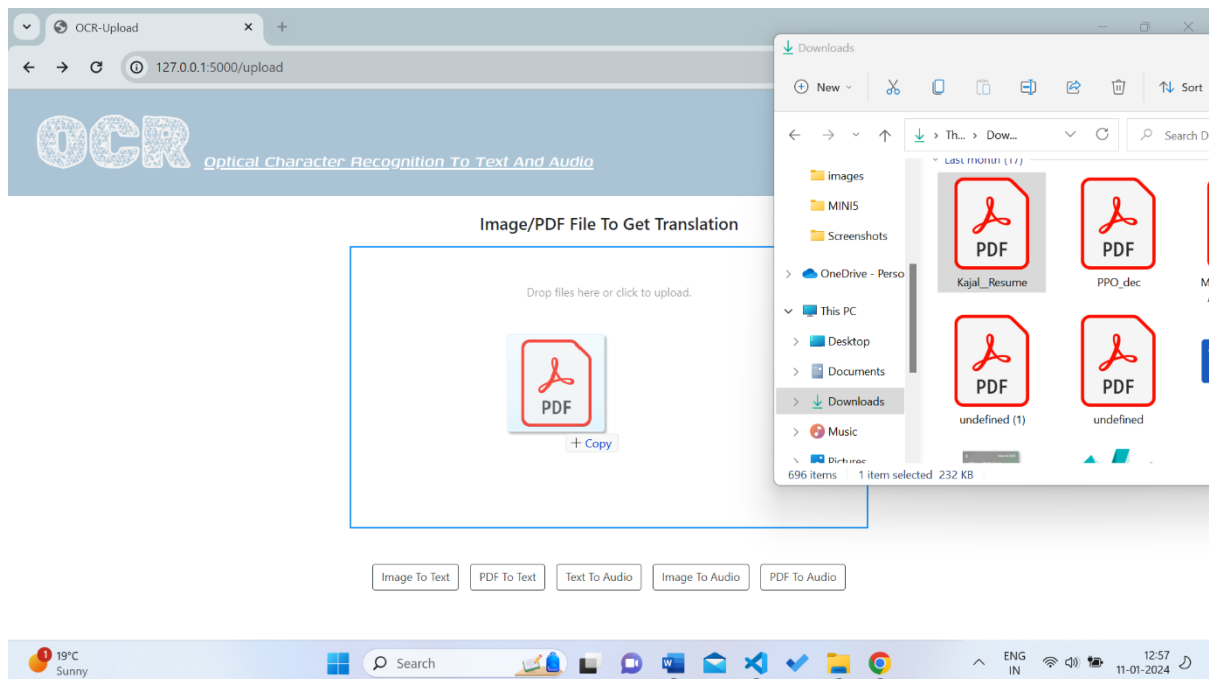


## Chapter 4

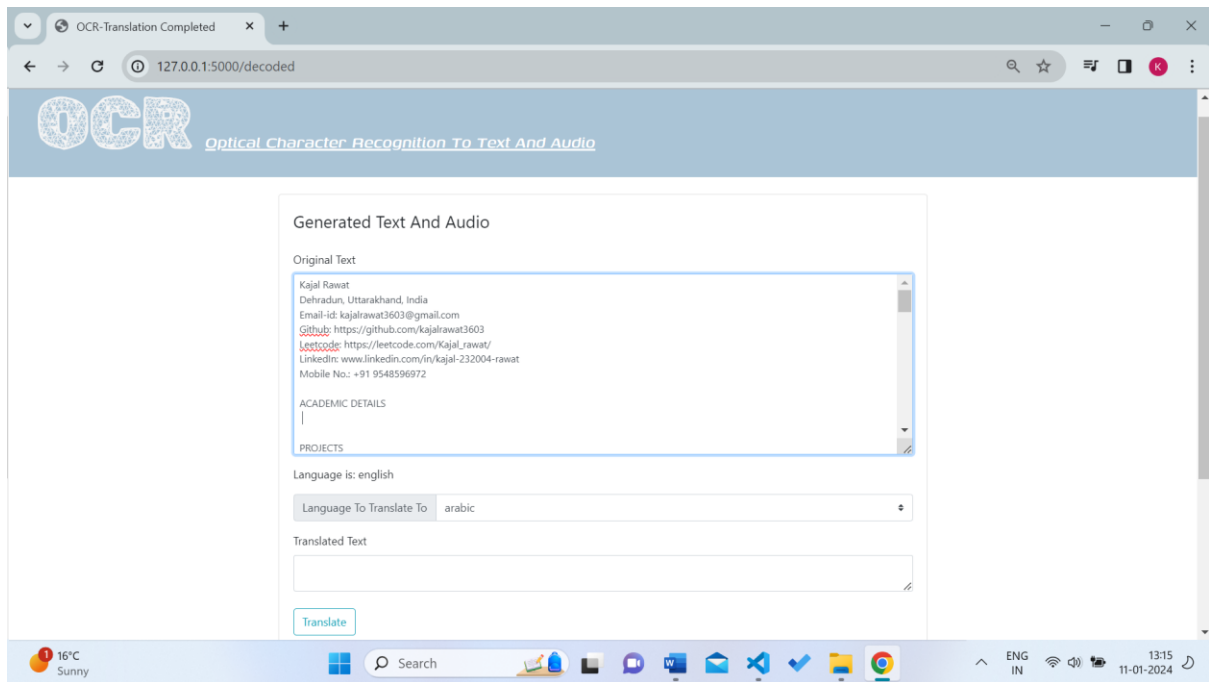
### RESULT



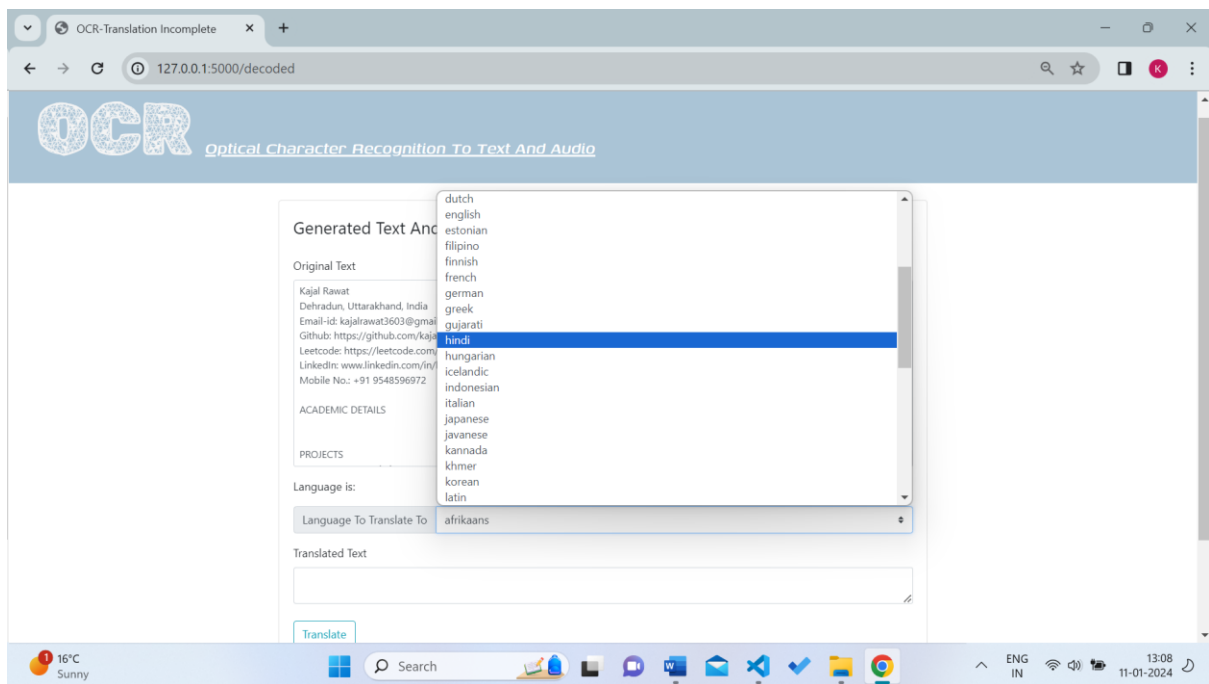
*Fig.2: User Interface*



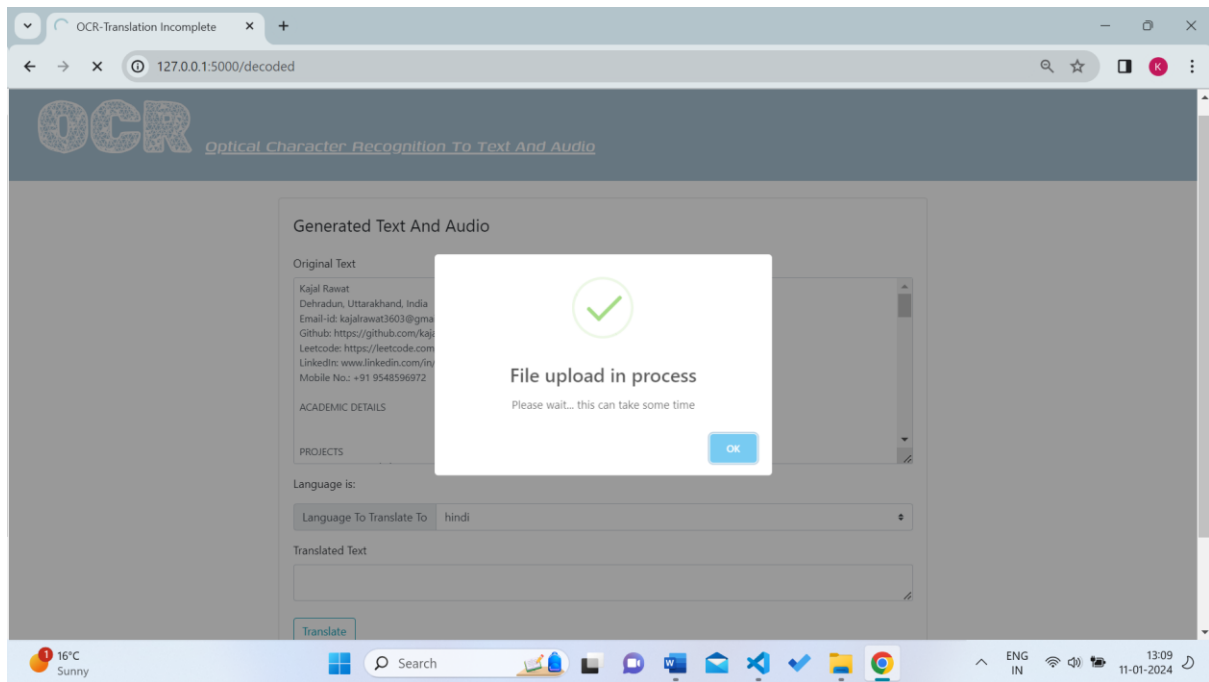
*Fig.3: Files Uploading*



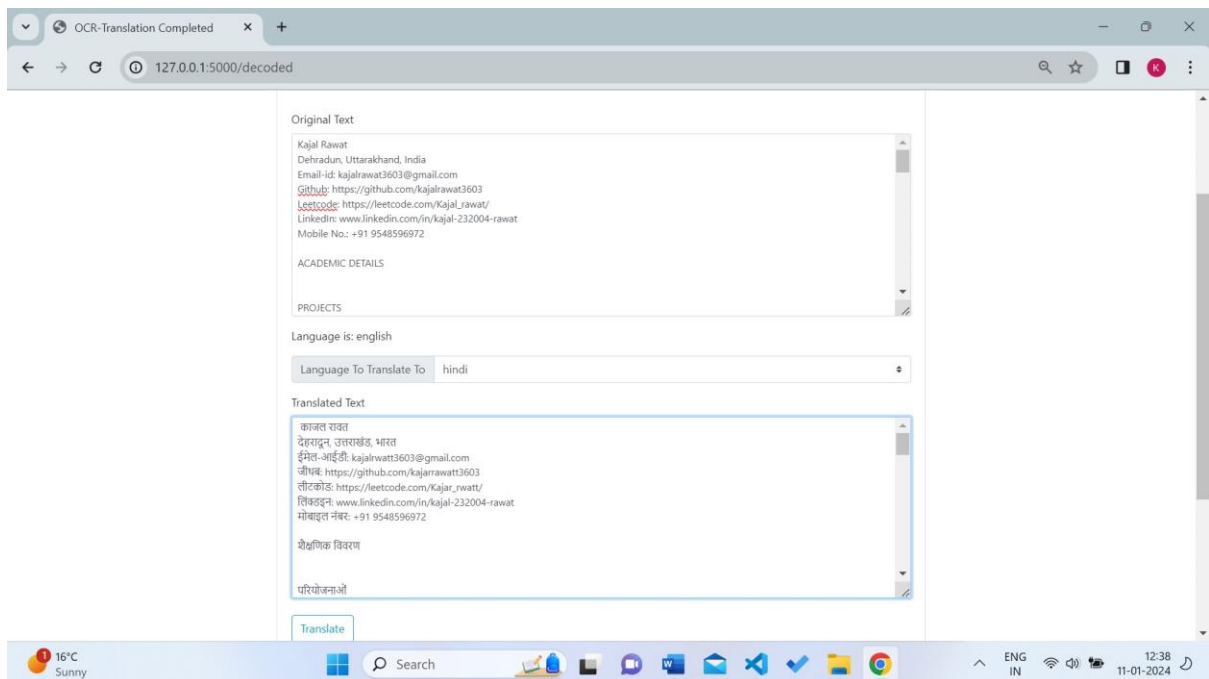
*Fig.4: Extracted Text*



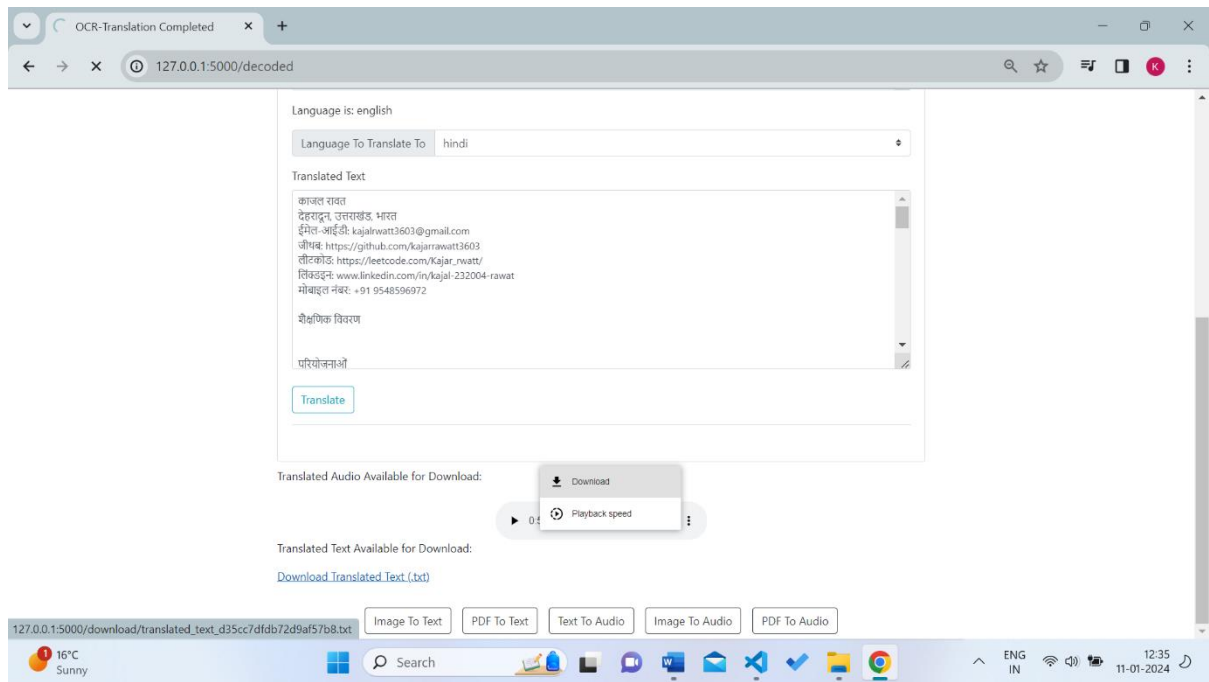
*Fig.5: language Selection*



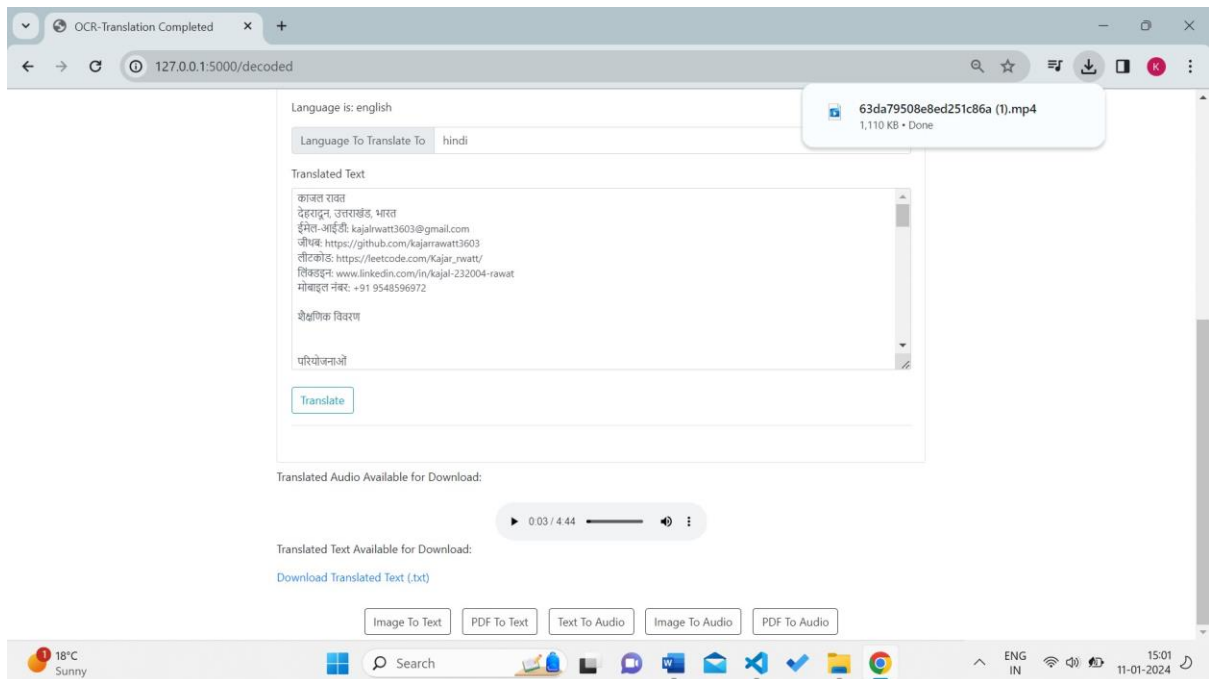
*Fig.6: Translation Phase*



*Fig.7: Translation Completed*



*Fig.8: Audio and Text files generated for download.*



*Fig.9: Downloaded Text/Audio files*

## Chapter 5

### SAMPLE CODE

Environment Settings/Requirments:

- IDE == VSCode
- Python Version== 3.11.0

| Package        | Version   | Description   |
|----------------|-----------|---|
| certifi        | 2021.10.8 | Certificate authority bundle for secure connections.        |
| chardet        | 3.0.4     | Detects the character encoding of text for proper decoding. |
| click          | 8.1.7     | Creates command-line interfaces for user interaction.       |
| Flask          | 3.0.0     | A web framework for building the OCR application.           |
| Flask-Dropzone | 1.6.0     | Integrates DropzoneJS for handling file uploads in Flask.   |
| gTTS           | 2.2.3     | Google Text-to-Speech API wrapper for generating speech.    |
| h11            | 0.9.0     | Supports HTTP/1.1 communication in Python.                  |
| h2             | 3.2.0     | Implements the efficient HTTP/2 protocol in Python.         |
| httpcore       | 0.9.1     | Powers the HTTP client for making requests.                 |
| httpx          | 0.13.3    | A feature-rich HTTP client for Python 3.7 and above.        |
| hyperframe     | 5.2.0     | Handles the framing layer for HTTP/2 in Python.             |
| opencv-python  | 4.5.5.62  | Computer vision library for image processing in Python.     |
| Pillow         | 9.3.0     | Image processing library essential for OCR tasks.           |
| playsound      | 1.3.0     | Plays sound, useful for notifications in the application.   |
| pytesseract    | 0.3.8     | OCR engine for recognizing text in images.                  |
| requests       | 2.27.1    | Simplifies making HTTP requests for data retrieval.         |
| urllib3        | 1.26.8    | Manages HTTP connections and requests in Python.            |

*Table.1: Required Pakages and Versions*

## Code SnapShots:

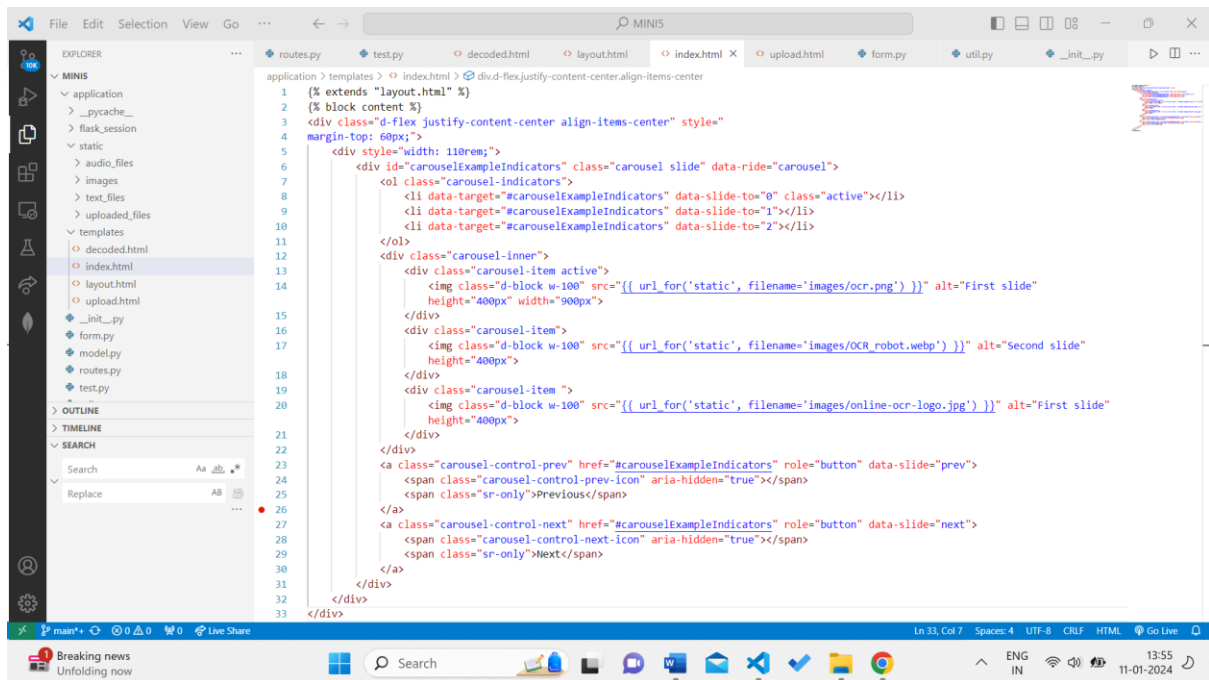


Fig.10: Sample HTML page (index.html)

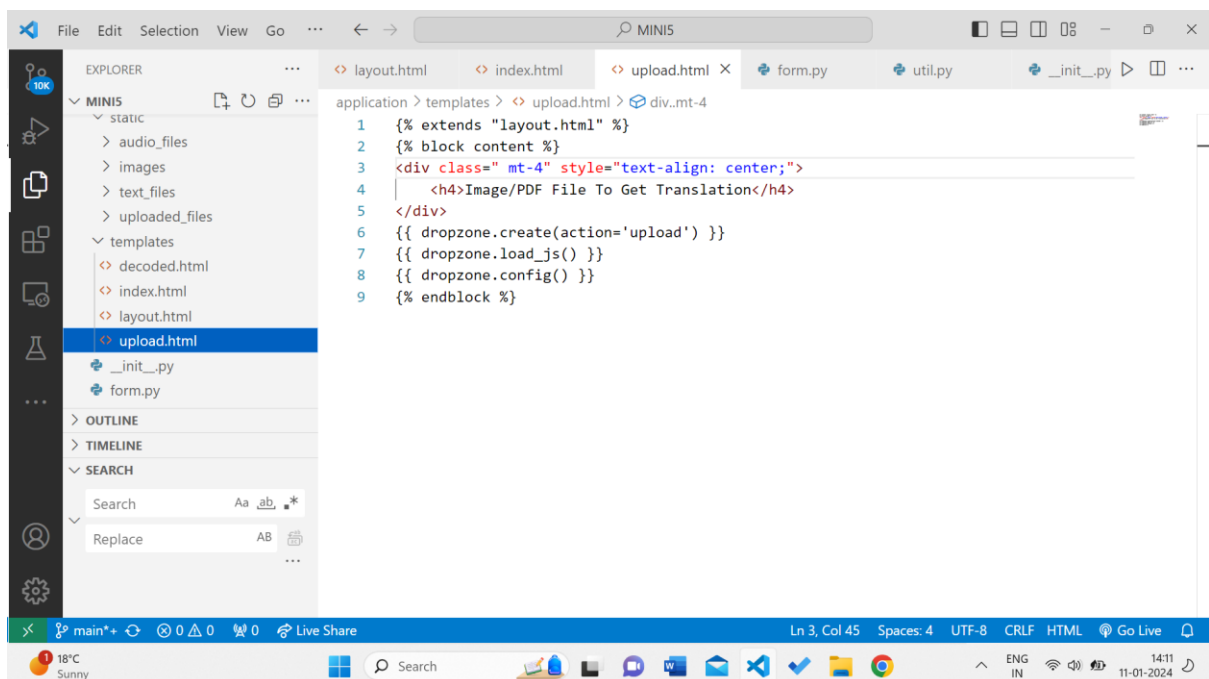


Fig.11: For uploading files using Drag and Drop (upload.html)

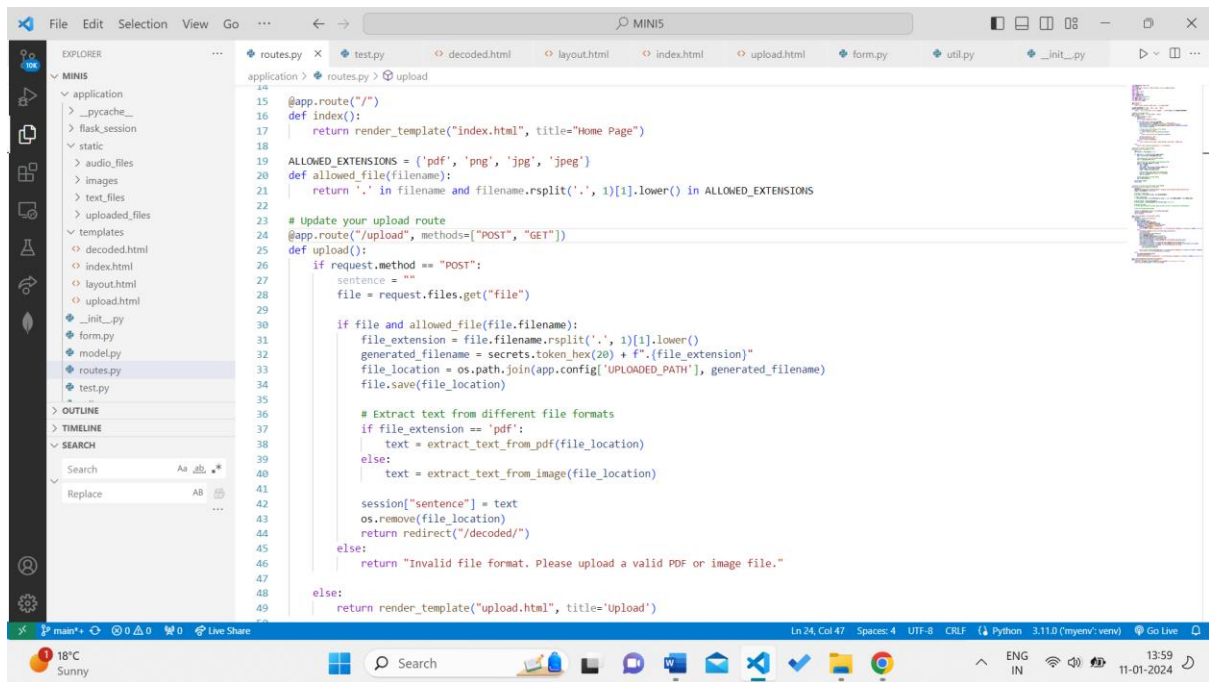


Fig.12: Request for upload Document (routes.py)

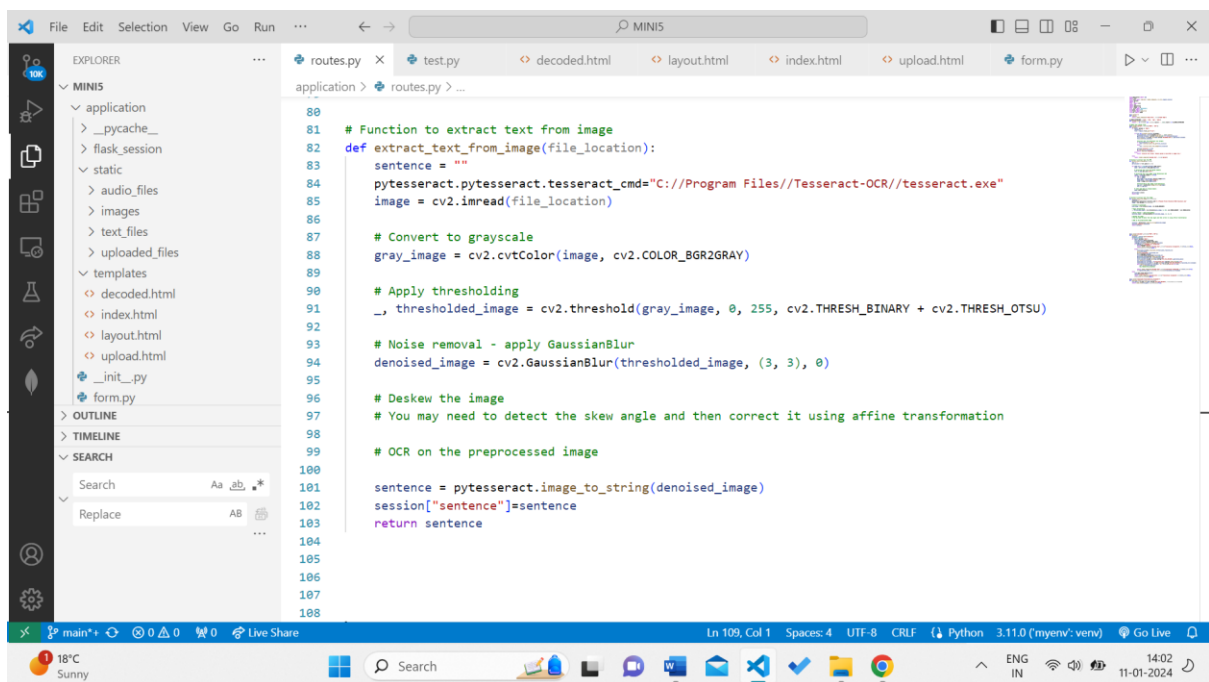


Fig.13: Text Extraction (routes.py)

```

108
109
110 @app.route("/decoded", methods=["POST", "GET"])
111 def decoded():
112     sentence = session.get("sentence")
113     form = QRCodeData()
114     if request.method == "POST":
115         text_data = form.data_field.data
116         translate_to = form.language_field.data
117         language = util.detect_language(text_data)
118         if language == "Language not recognized":
119             form.language = "Unknown"
120             return render_template("decoded.html", title="Translation Incomplete", form=form, audio=None, translated_txt_filename=None)
121         else:
122             translated_txt = util.translate_txt(text_data, translate_to)
123             print(text_data)
124             print("translated", translated_txt)
125             form.language = language
126             form.translated_field.data = translated_txt
127             generated_audio = secrets.token_hex(10) + ".mp4"
128             tts = gTTS(translated_txt, lang=translate_to)
129             file_location = os.path.join(app.config["TXT_FILE_UPLOAD"], generated_audio)
130             tts.save(file_location)
131             translated_txt_filename = f"translated_text_{secrets.token_hex(10)}.txt"
132             translated_txt_file_path = os.path.join(app.config["TXT_FILE_UPLOAD"], translated_txt_filename)
133             with open(translated_txt_file_path, "w", encoding="utf-8") as file:
134                 file.write(translated_txt)
135             # os.remove(file_location)
136
137             return render_template("decoded.html", title="Translation Completed", form=form, audio=True, file=generated_audio,
138                                   translated_txt_filename=translated_txt_filename)
139     else:
140         form.data_field.data = sentence
141         session["sentence"] = ""
142         return render_template("decoded.html", title="Translation Incomplete", form=form, audio=None, translated_txt_filename=None)

```

Fig.14: Post request to send translated files (routes.py)

```

1 from wtforms import TextAreaField, SubmitField, SelectField
2 from flask_wtf import FlaskForm
3 from wtforms.validators import DataRequired, Length
4
5 # import utils
6 from . import util
7
8
9 languages_choice = []
10 for key, value in util.languages.items():
11     languages_choice.append((key, value))
12
13
14 class QRCodeData(FlaskForm):
15     data_field = TextAreaField('Original Text', validators=[DataRequired(), Length(min=1, max=250)])
16     language = ""
17     language_field = SelectField("Language to translate to", choices=languages_choice)
18     translated_field = TextAreaField('Translated Text')
19     submit = SubmitField('Translate')

```

Fig.15: Form to get the translated files (form.py).



```

1 from flask import Flask
2 from flask_dropzone import Dropzone
3 import os
4 from flask_session import Session
5
6 app=Flask(__name__)
7
8 app.config['SECRET_KEY']='5f395692341b60f5f364e40cf4358c3c916c5a27285b4a958f85e6c29119'
9
10 SESSION_TYPE="filesystem"
11 app.config.from_object(__name__)
12 Session(app)
13 dir_path=os.path.dirname(os.path.realpath(__file__))
14
15
16 app.config.update(
17     UPLOADED_PATH=os.path.join(dir_path,"static/uploaded_files"),
18     DROPEZONE_ALLOWED_FILE_TYPE='image',
19     DROPZONE_MAX_FILE_SIZE=3,
20     DROPZONE_MAX_FILES=1,
21     AUDIO_FILE_UPLOAD=os.path.join(dir_path,"static/audio_files"),
22     TXT_FILE_UPLOAD=os.path.join(dir_path,"static/text_files")
23 )
24
25 app.config['DROPZONE_REDIRECT_VIEW']='decoded'
26 dropzone = Dropzone(app)
27 from application import routes
28
29

```

Fig.16: (\_\_init\_\_.py)

```

1 # pip install googletrans==3.1.0a0
2 from cv2 import triangulatePoints
3 from googletrans import Translator
4
5 translator = Translator()
6
7 def detect_language(text):
8     # get language used
9     detected_lang_data = translator.detect(text)
10    # print(lang)
11    lang = languages.get(detected_lang_data.lang,"Language not recognized")
12    conf = detected_lang_data.confidence
13
14    return lang, conf
15
16 def translate_txt(text, dest):
17     translated_text = translator.translate(text, dest=dest)
18     return translated_text.text
19
20
21 languages = {
22     'af': 'afrikaans',
23     'sq': 'albanian',
24     'en': 'english',
25     'es': 'spanish',
26     'fr': 'french',
27     'de': 'german',
28     'it': 'italian',
29     'ja': 'japanese',
30     'ko': 'korean',
31     'pt': 'portuguese',
32     'ru': 'russian',
33     'zh': 'chinese'
34 }

```

Fig.17: Translation of text

## Chapter 6

### **FUTURE SCOPE**

- **Improved Language Support:** Expanding language options for translation and enhancing accuracy, especially for less common languages.
- **Advanced Image Processing:** Implementing sophisticated algorithms and deep learning models for better text extraction from complex or degraded images.
- **AI-based Document Understanding:** Developing systems to analyze document structures, extract meaningful content, and automate summarization.
- **Enhanced Accessibility:** Improving the user interface for accessibility and incorporating features for users with disabilities.
- **Real-time Translation and Interaction:** Enabling live translation and voice-based interaction for instant communication.
- **Cloud Integration and Scalability:** Moving towards cloud-based solutions for scalability and better performance.

## Chapter 7

### CONCLUSION

Through OCR and translation technology, this project takes a big step in making information easy to get. The web app we built can read and translate text from PDFs and images smoothly. We used tools like Pytesseract, and gTTS to make it work. This helps people understand and use information, no matter the language.

During the project, we added key features like handling different file types, pulling out text, translating, and creating audio. I also made sure the app can handle errors well and get feedback from users, making it user-friendly even when things get tough.

While we achieved a lot, there's more enhancement that we can do like to add more languages, improve how we process images, and use AI to understand documents better. Making the app work on mobiles, doing real-time translations, and making it even friendlier with a better look are on our to-do list.

In conclusion, this project isn't just about solving today's problems; it's about starting something new. We're paving the way for more ideas in AI, making information accessible to everyone. This could change how people connect with and get information, making our digital world more inclusive.

## Chapter 8

### **REFERENCES**

- [ibm.com/blog/optical-character-recognition/](https://ibm.com/blog/optical-character-recognition/)
- <https://gtts.readthedocs.io/en/latest/>
- <https://flask-dropzone.readthedocs.io/en/latest/>
- <https://flask.palletsprojects.com/en/3.0.x/>
- <https://pillow.readthedocs.io/en/latest/installation.html>
- [https://docs.opencv.org/4.x/d7/dbd/group\\_\\_imgproc.html](https://docs.opencv.org/4.x/d7/dbd/group__imgproc.html)
- <https://py-googletrans.readthedocs.io/en/latest/>
- <https://ieeexplore.ieee.org/document/9935961>
- <https://ieeexplore.ieee.org/document/6993174>
- [youtube](#)