**Department of Artificial Intelligence and Data Sciences**

# AI-Driven Document Intelligence: A Comprehensive Approach to Classification, Q&A Handling, and Fraud Detection

**Submitted By:  Group 1**

Shambhavi Rai      00901172020

Sonanshi Goel      01601172020

Princy Singhal      05101172020

Mehak Aggarwal     03401172020

Kanika Kanojia      06301172020

**Supervisor :**

Ms. Ritika Kumari

Assistant professor

AI & DS Dept.

# INDEX

- Introduction
- Literature Review
- Research Gaps
- Research Objectives
- Proposed Methodology
- Experimental Results
- Conclusion and Future Scope
- References

# INTRODUCTION

With the growth in dependency on electronic verification systems used by both government and private organizations, there has been an increase in the usage of digitized manuscripts. This demands user-friendly and efficient document interaction tools.

Optimization of the document management process has become AI-driven with features like identity document classification, authenticity checking algorithms, and interactive question answering.

If we effectively identify the document type, we can streamline workflows and enable secure verification processes. The project employs Deep Learning techniques to classify documents like PAN cards, Aadhar cards, etc.

Additionally, we employ a Similar Document Template Matching Algorithm which can seamlessly extract features, match the template with a real document, and thereby detect fraud.

Furthermore, Large Language Models (LLMs) have been leveraged to enable interactive question answering over documents. Alongside these primary features, we have integrated additional functionalities such as paraphrasing, grammar checking, read aloud, and summarization. These enhancements aim to improve user interaction, making the system more versatile and accessible.

Experimental findings in our study showed the effectiveness of our work in improving efficiency, accuracy, and usability in document management.

# LITERATURE REVIEW

| REF NO | METHODOLOGY ADOPTED | RESEARCH GAPS | YEAR |
|--------|---------------------|---------------|------|
| 01 | CNN, VGG-16 and YOLO | Long training time and image-containing documents require feature extraction. Deep learning uses the image directly for classification but is time-consuming and requires hyperparameter tuning. | 2021 |
| 02 | OCR and Similarity Score | Only looks at only one saved sample to predict the outputs, can be made to predict based on all saved samples of the specific template to generalize better and improve overall accuracy | 2019 |
| 03 | OCR and SSIM | Concerns persist due to variations in document layouts, necessitating comprehensive labeled datasets, scalability for large data, and adaptability to real-world scenarios. | 2023 |

| REF NO | METHODOLOGY ADOPTED | RESEARCH GAP | YEAR |
|---|---|---|---|
| 04 | Closed-book Generation (T5) and (BART)<br><br>Retrieval-augmented Generation(RAG)<br><br>LLM-based Generation models (gpt-3.5-turbo-0613) and (LLaMA2-13B-Chat) | The study lacks a comparative analysis of the integrated UniGen framework against separate models for Generative Document Retrieval (GDR) and Grounded Answer Generation (GAR), hindering understanding of its comparative advantages. | 2024 |
| 05 | Automated question generation with human-guided templates. | ToolQA heavily relies on external tools for question answering, which could introduce biases or limitations based on the effectiveness and coverage of these tools. | 2024 |

# RESEARCH GAPS

1. **Integrated Technology Approach:**
   a. **Existing Gap**: Prior research often isolates document classification,OCR etc. Hence missing the synergy of combined technologies.
   b. **Our Contribution**: Seamlessly integrates classification, NLP, and speech synthesis etc for a unified solution.
2. **Robust Performance in Varied Conditions:**
   a. **Existing Gap**: Challenges in achieving consistent accuracy across different imaging conditions.
   b. **Our Contribution**: Employs advanced data augmentation and fine-tuning to ensure reliable classification and fraud detection in diverse environments and document types.
3. **Efficiency & Accuracy:**
   a. **Existing Gap**:  Traditional CNN architectures, while powerful, are hampered by lengthy training times and intensive hyperparameter tuning when processing image-heavy documents.
   b. **Our Contribution**: Achieves an optimal balance between high accuracy and computational efficiency, leveraging pre-trained models.
4. **Integration Challenges in AI-Driven Q&A Systems:**
   a. **Existing Gap**: Recent tools are either too complex and resource-heavy or fail to connect information smoothly, making them inefficient for fast and large-scale applications.
   b. **Our Contribution**: Implements a streamlined, efficient Q&A system that balances resource usage with effective information retrieval, ensuring scalability and speed.

# RESEARCH OBJECTIVES

1. **Optimizing Architectures:**
   - Improve performance by modifying custom CNN architecture.
   - Evaluate the impact of pretrained models (e.g., VGG16, VGG19).

2. **Dataset Formulation:**
   - Utilize advanced data augmentation for robust classification model training.
   - Create test cases to detect subtle document fraud variations.

3. **Real-Time Document Processing:**
   - Develop efficient models for real-time classification, interactive query answering, and fraud detection.

4. **Efficient Information Retrieval Systems:**
   - Implement sophisticated question answering systems using NLP and semantic search for real-time information extraction from unstructured PDFs.
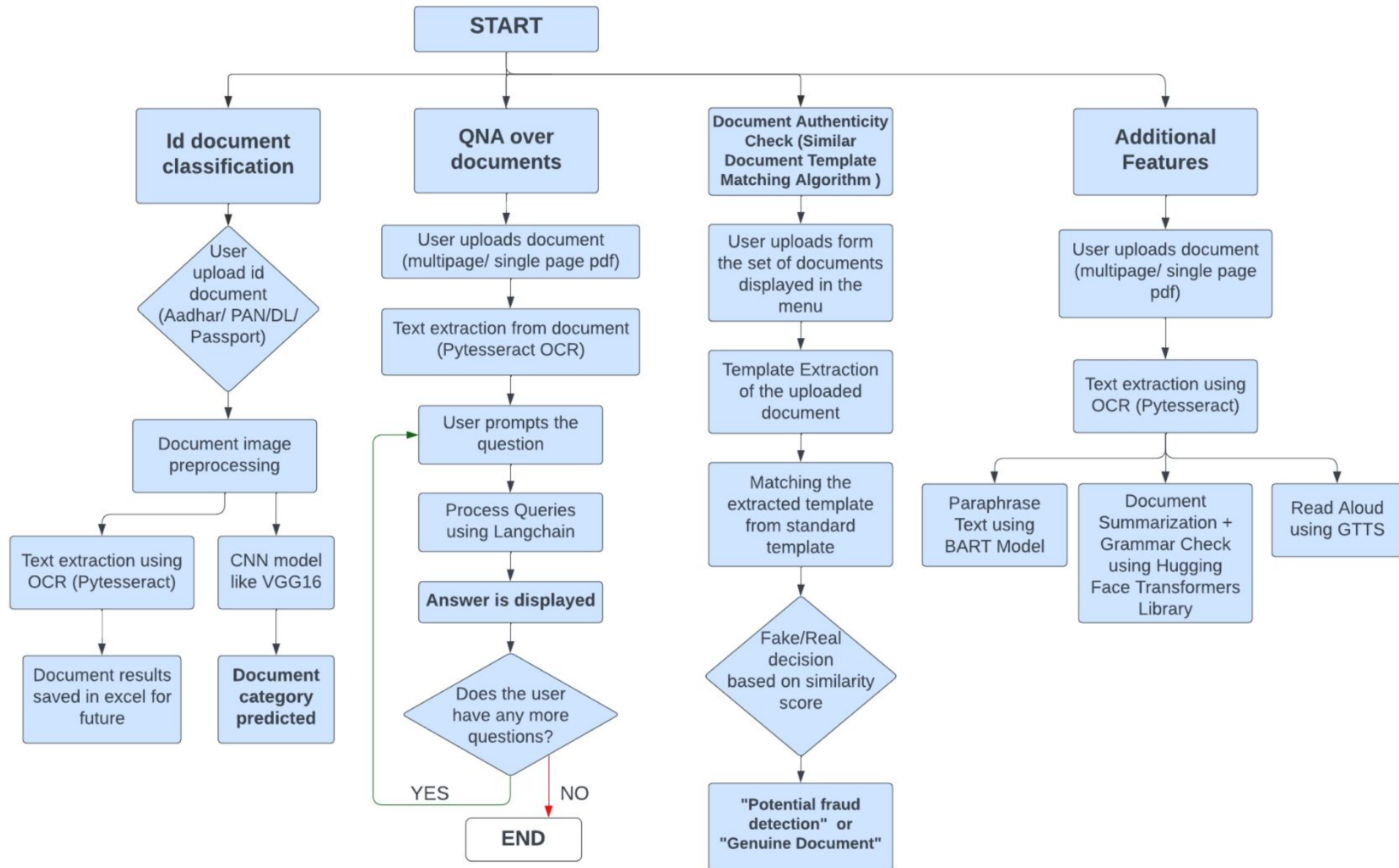
5. **Quality Fraud Detection Mechanisms:**
   - Create robust template matching algorithms with ORB feature descriptors to detect subtle identity document forgeries.

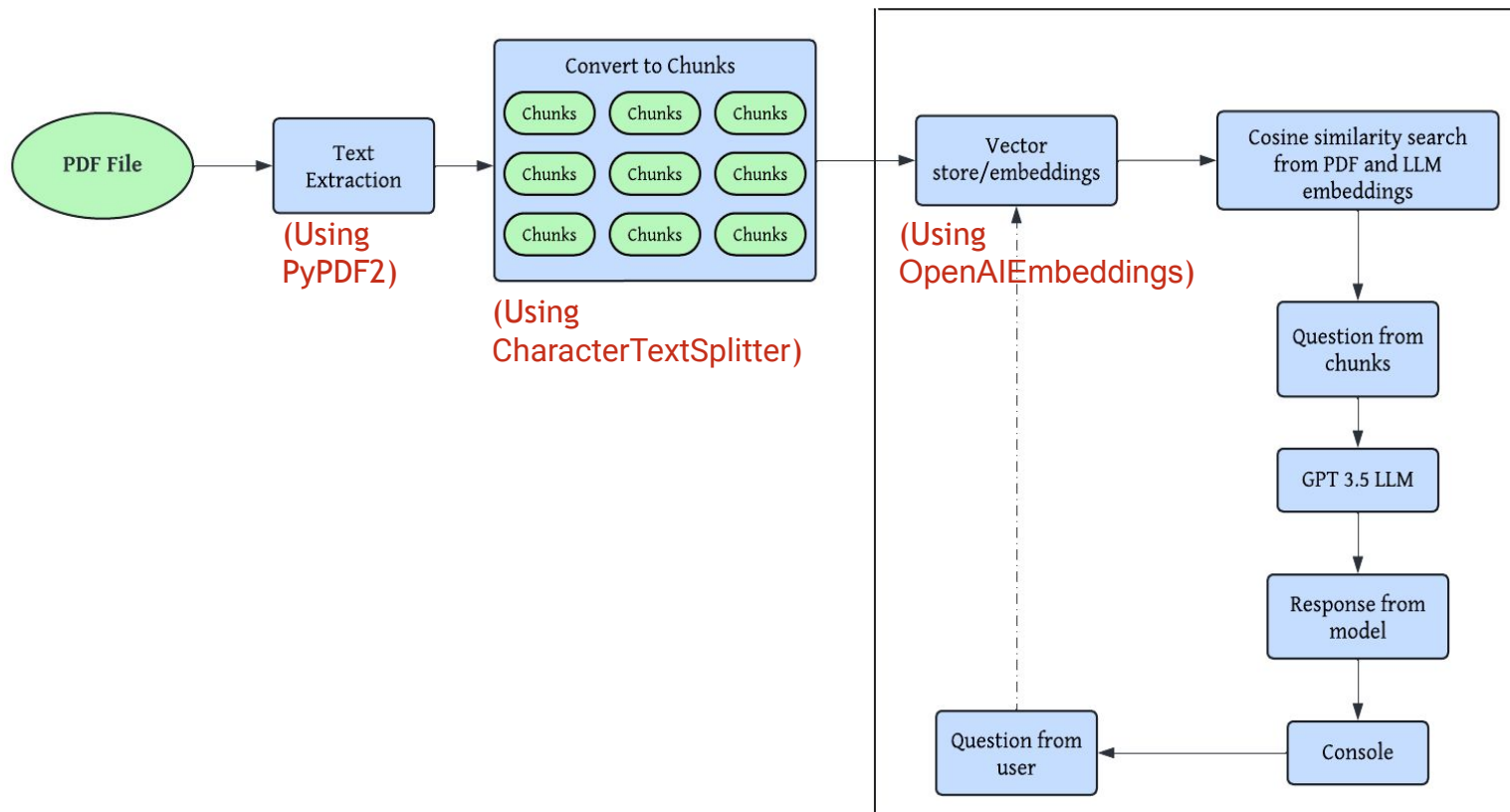6. **Enhancement of Document Accessibility and Quality:**
   - Incorporate auxiliary features like text-to-speech, summarization, grammar checking, and paraphrasing to improve accessibility, comprehension, and quality.

# PROPOSED METHODOLOGY

# I) INTERACTIVE QUERY ANSWERING

Leverages advanced Language Learning Models (LLMs) combined with OpenAI embeddings to iteratively extract and provide accurate answers from documents. This approach ensures enhanced context understanding and precise information retrieval, improving the overall efficiency and reliability of the QnA process
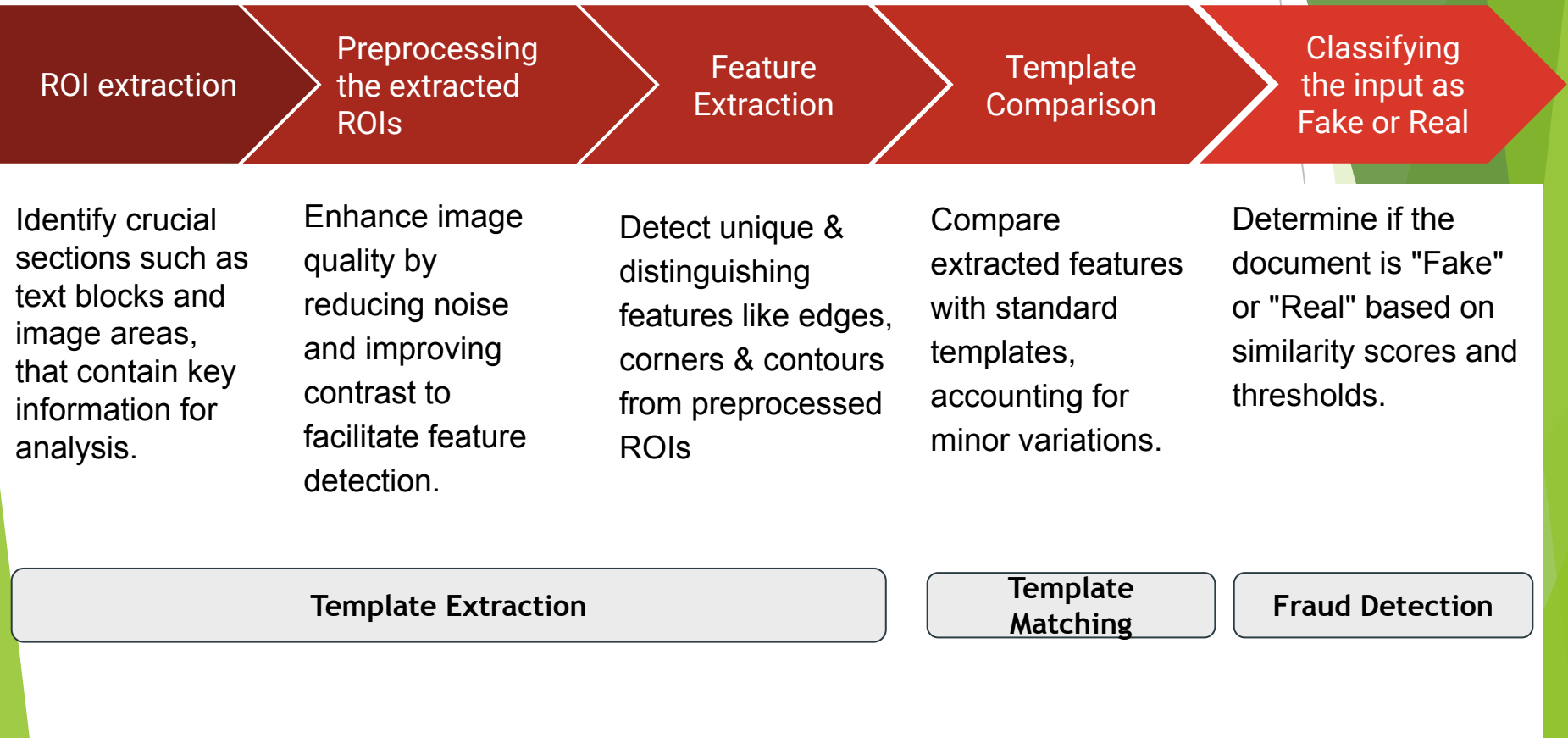
System Architecture of Interactive Query Answering

# EXPERIMENTAL RESULTS
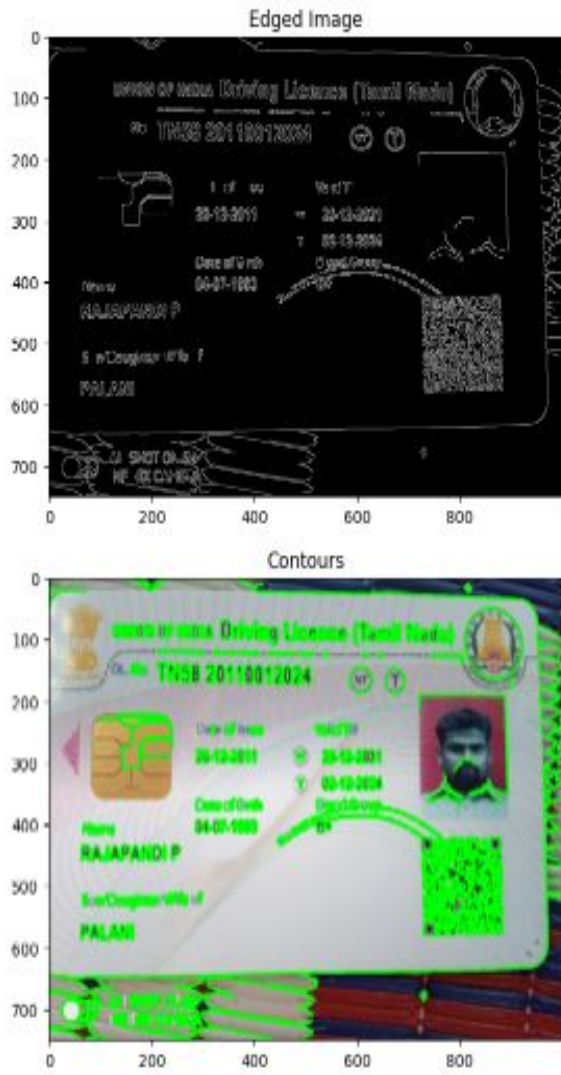
Comparison between LLM and BERT

| Length of pdf (in pages) | Confidence Score (LLM Model) | Confidence Score (BERT Transformer) |
|---|---|---|
| 1 | 98.43 | 80.04 |
| 4 | 94.46 | 56.51 |
| 12 | 92.85 | 55.46 |

# II) SIMILAR TEMPLATE MATCHING ALGORITHM (Code)

| ROI extraction | Preprocessing the extracted ROIs | Feature Extraction | Template Comparison | Classifying the input as Fake or Real |
|---|---|---|---|---|
| Identify crucial sections such as text blocks and image areas, that contain key information for analysis. | Enhance image quality by reducing noise and improving contrast to facilitate feature detection. | Detect unique & distinguishing features like edges, corners & contours from preprocessed ROIs | Compare extracted features with standard templates, accounting for minor variations. | Determine if the document is "Fake" or "Real" based on similarity scores and thresholds. |

**Template Extraction**

**Template Matching**

**Fraud Detection**

# ROI extraction



Edge and Contour detection

Document IMAGE

↓

Transform into grayscale — Using OpenCV

↓

Edge Detection — Using dynamic thresholding

↓

Contour Detection — for separating out specific areas

↓

List of extracted ROIs — **Output**

↓

Each ROI is passed for preprocessing

# Preprocessing ROIs

Original ROI

Processed ROI

Loop through each ROI

for (x, y, w, h) in rois:

Crop ROI from Document Image

roi_image = document_image[y:y+h, x:x+w]

preprocess_function(roi_image)

Noise Reduction

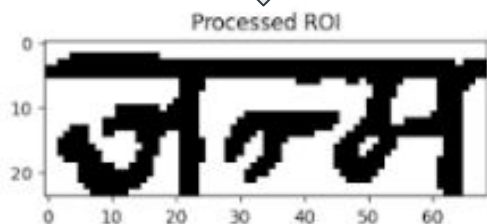denoised = cv2.fastNlMeansDenoising(gray, None, 30, 7, 21)
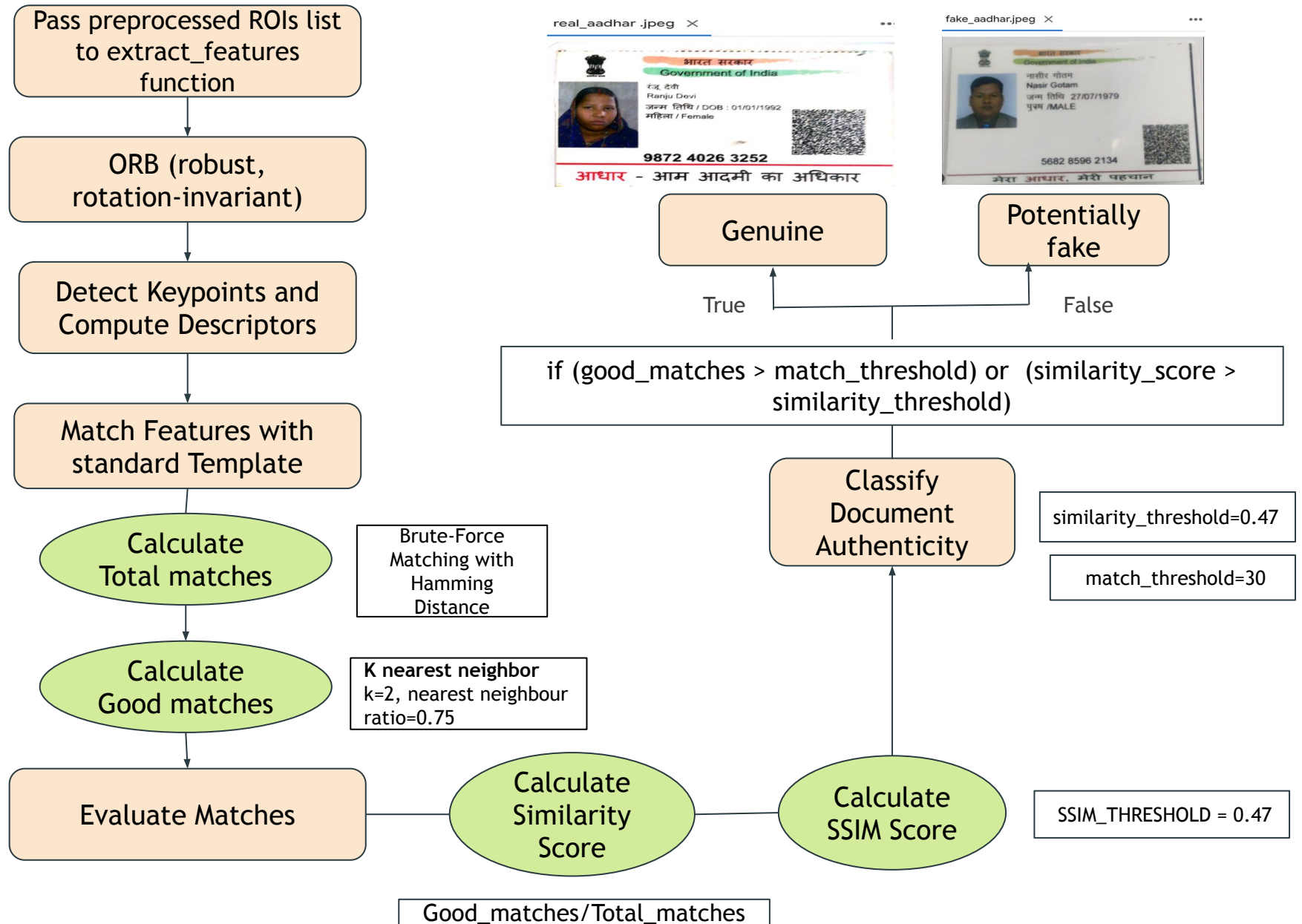
Contrast Enhancement

CLAHE

Binary Thresholding

Otsu

Morphological Operations

Store Preprocessed ROIs

OUTPUT

# Feature Extraction, Template Matching & Fraud Detection

Pass preprocessed ROIs list to extract_features function

↓

ORB (robust, rotation-invariant)

↓

Detect Keypoints and Compute Descriptors

↓

Match Features with standard Template

↓

Calculate Total matches

Brute-Force Matching with Hamming Distance

↓

Calculate Good matches

**K nearest neighbor**
k=2, nearest neighbour ratio=0.75

↓

Evaluate Matches

Calculate Similarity Score

Good_matches/Total_matches

Calculate SSIM Score

SSIM_THRESHOLD = 0.47

Classify Document Authenticity

similarity_threshold=0.47

match_threshold=30

if (good_matches > match_threshold) or (similarity_score > similarity_threshold)

True → Genuine

False → Potentially fake

real_aadhar .jpeg ×

fake_aadhar.jpeg ×

# EXPERIMENTAL RESULTS

Analysis of Document Authenticity through Feature Matching and SSIM Evaluation

| Image | Key features | SSIM | Decision Obtained | Expected Decision | Observation |
|---|---|---|---|---|---|
| aadhar_img1 | With display picture variation, no qr code, missing name in hindi | 0.29 | Potential fraud detected | Fraud | Significant changes detected very well |
| driver_Img2 | Black and white, rotated | 0.54 | Document is likely genuine | Real | Color invariant |
| passport_img3 | Face hidden | 0.20 | Fraud | Fraud | Recognises inconsistencies with Display image in the id |
| aadhar_img4 | Date font size changes and logo missing | 0.20 | Potential fraud detected | Fraud | Worked on subtle variation Like text font and size |
| aadhar_img5 | Colored, well aligned and illuminated | 0.21 | Potential fraud detected | Real | False positives are high. Here additional checks can be employed like re-uploading etc |
| pan_Img6 | Colored, well aligned and illuminated | 0.51 | Document is likely genuine | Real | Ideal document |

# III) DOCUMENT CLASSIFICATION

Our model employs a custom CNN architecture particularly designed with convolutional layers progressively increasing in filter count, dropout layers to reduce overfitting, and dense layers to finalize classification. This specialized CNN architecture provides powerful feature extraction, allowing the model to recognize detailed patterns unique to each document format. In parallel, the use of the pre-trained VGG16 model improves feature extraction capabilities, resulting in a more detailed representation of document features.
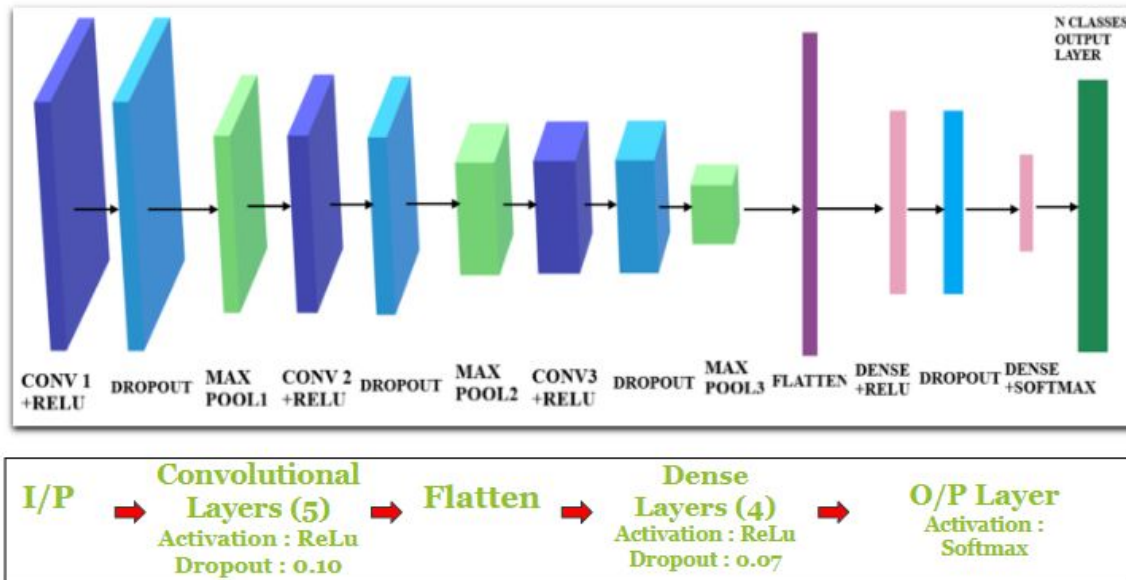
**Dataset Overview :**

| Document Type | Aadhar Card | Pan Card | Driving License | Voter Id | Passport |
|---|---|---|---|---|---|
| No. of Documents | 129 | 45 | 64 | 76 | 37 |



Few images from Training set

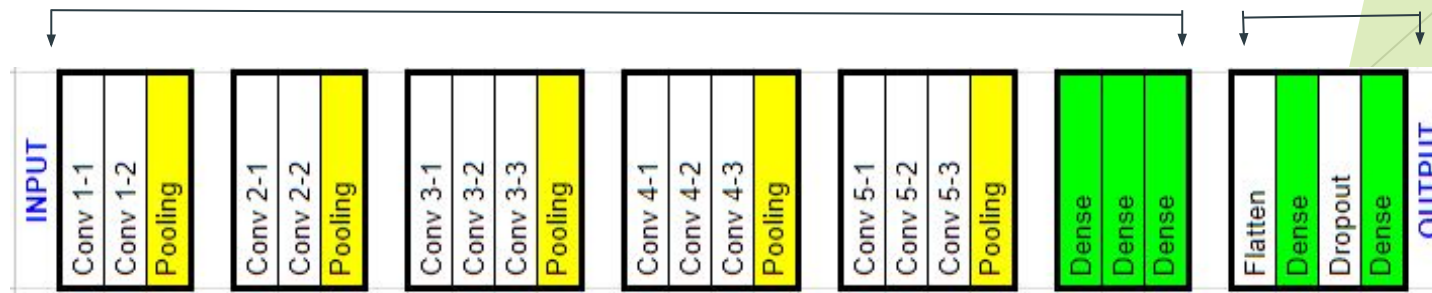# Proposed Classification Model

**CNN Model Layers**



**VGG16 Model Layers**      **VGG16**      **Additional layers**

# EXPERIMENTAL RESULTS
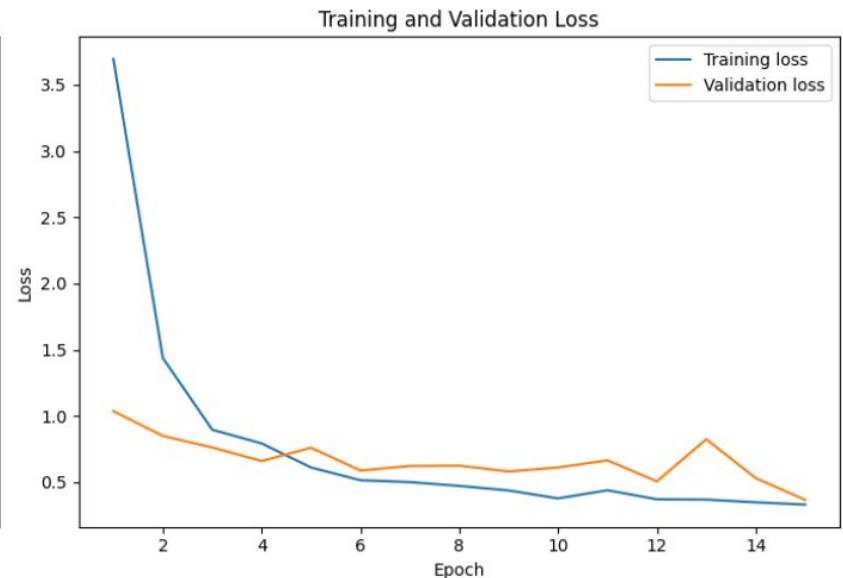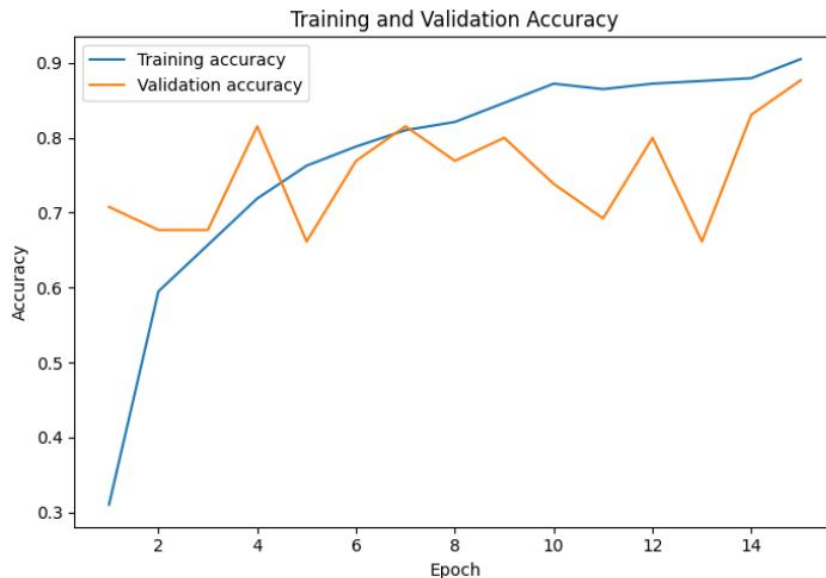
## Custom CNN model results with hyperparameter tuning

| S.No. | Model name | number of layers | Epochs | Other parameters | Validation Accuracy | Observation |
|-------|-----------|------------------|--------|------------------|---------------------|-------------|
| 1. | Custom CNN link | 9 | 40 | Dropout = 0.07 | 77.36 | limited data augmentation and a smaller network architecture |
| 2. | VGG 16 | Additional 4 | 10 | | 80.00 | Better generalization on our dataset due to pre-learned features |
| | link | | 15 | | 87.69 | additional training time |
| | | | 25 | | 70.77 | Model may have learnt noise in the training data |
| | | | 40 | | 87.69 | better generalization, leading to improved accuracy |
| | | | 10 | Learning rate = 1e-5<br><br>Patience = 5<br><br>Unfrozen layers = 4<br><br>Class importance applied= yes | 76.92 | Overfitting due to higher unfreezed layers |
| | | | 15 | Unfrozen layers= 2<br><br>patience=6 | 86.15 | Tackling of unbalanced dataset and lower learning rate for fine tuning purposes |
| 3. | VGG 19 | Additional 4 | 20 | Learning rate=1e-5 | 78.46% | |

# EXPERIMENTAL RESULTS

**Comparative Analysis of CNN and VGG16 model**

| Model | Number of layers | Epochs | Parameter | Train Accuracy (in %) | Validation Accuracy (in %) | Observation |
|-------|------------------|--------|-----------|----------------------|---------------------------|-------------|
| CNN | 9 | 40 | Dropout = 0.07 | 69.44 | 83.33 | Limited data augmentation and a smaller network architecture |
| VGG16 | Additional 5 | 40 | Learning rate = 1e-5 | 90.51 | 87.69 | Better generalisation leading to improved accuracy |

**Accuracy and loss graph for finalised VGG16 model**

# Information Extraction Model (OCR)

- Applied OCR for text extraction from scanned documents.
- Conducted image processing, text extraction, and cleaning.
- Extracted information using keywords from the cleaned text.
- Saved extracted data into a downloadable Excel file.

**Text extraction using Tesseract OCR engine.** link

| | Name | Father Name | Document Type | DOB | Address | Document Number | Sex |
|---|---|---|---|---|---|---|---|
| 0 | KUSUM LATA | DHANI RAM | PAN | 17-10-1992 | | AQSPL9772C | |
| 1 | NAVNEET NAYAL | | Adhaar Card | 16-10-1997 | | 472672990081 | MALE |
| 2 | CHAMDRKANT YADAV | MAHADEV YADAV | Voter ID Card | 01-06-1963 | E-SECTOR, BQLINE ROOM 2 MUMBAI | WIC7896681 | MALE |
| 3 | AHMED ALI SHAIKH | MOHM ALI SHAIKH | Driving License | 21-06-1992 | 13 Kisoli Village , Gulaothi Block , Bulandsha... | MH032014001542 | |
| 4 | PREM SINGH BOHRA | | Adhaar Card | 15-12-1988 | | 603313250609 | MALE |
| 5 | RAJESH BALKRISHNA MISHRA | BR MISHRA | PAN | 01-01-1990 | | AUUPM6954D | |

# IV) READ ALOUD

Converts uploaded text into speech using the gTTS library.

**Code Snippet**

```python
def text_to_audio():
    print("Please upload a PDF file.")
    menu_option = 1
    extracted_text = upload_pdf_and_convert(menu_option)

    # Clean the extracted text by removing newline characters.
    cleaned_text = [text.replace('\n', ' ') for text in extracted_text.values()]
    print("Extracted text : ", cleaned_text)

    # Combine the cleaned text into a single string.
    combined_text = ' '.join(cleaned_text)

    # Using gTTS library to convert the combined text to speech.
    tts = gTTS(text=combined_text, lang='en')

    # Saved generated speech as MP3 file.
    tts.save('output.mp3')

    print("\n\n Audio has been created!\n\n")

    # Return the created audio file and autoplay it.
    return Audio('output.mp3', autoplay=True)
```
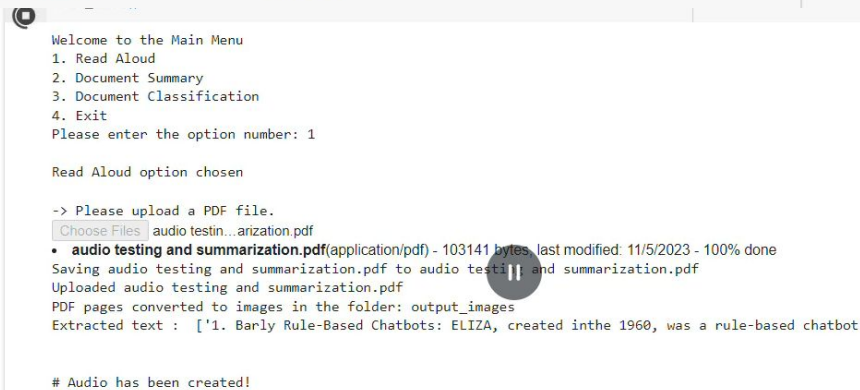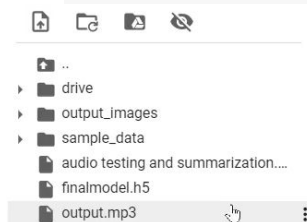
**Output**

```
drive
output_images
sample_data
audio testing and summarization....
finalmodel.h5
output.mp3
```

```
Welcome to the Main Menu
1. Read Aloud
2. Document Summary
3. Document Classification
4. Exit
Please enter the option number: 1

Read Aloud option chosen

-> Please upload a PDF file.
  Choose Files  audio testin...arization.pdf
  • audio testing and summarization.pdf(application/pdf) - 103141 bytes, last modified: 11/5/2023 - 100% done
Saving audio testing and summarization.pdf to audio testing and summarization.pdf
Uploaded audio testing and summarization.pdf
PDF pages converted to images in the folder: output_images
Extracted text :  ['1. Barly Rule-Based Chatbots: ELIZA, created in the 1960, was a rule-based chatbot

# Audio has been created!
```

# V) PARAPHRASING

- The work leverages the BART model from the transformers library to paraphrase text extracted from a PDF. It begins by importing necessary components and utilizing a function that splits a given text into smaller, manageable chunks.
- This is followed by the loading a pre-trained BART model and tokenizer, chunks the input text, and then paraphrases each chunk individually. The chunks are then recombined into a single paraphrased text.

**Actual text**

Values in a Python dictionary can be accessed by placing the key within square brackets next to the dictionary. Values can be written by placing key within square brackets next to the dictionary and using the assignment operator . If the key already exists, the old value will be overwritten. Attempting to access a value with a key that does not exist will cause a `KeyError`.

**Paraphrased text** `BLEU Score: 0.6408`

```
Values in a Python dictionary can be accessed by placing the key within square brackets next to the dictionary
 If the key already exists, the old value will be overwritten
 'Attempting to access a value with a key that does not exist will cause a KeyError'
```

# VI) SUMMARISATION

Document summarization condenses text into concise summaries, facilitating quick understanding and information extraction.

**Code Snippet**

```python
from transformers import pipeline

# Function to summarize a given text
def summarize_text(text):
    # Initialize the summarization pipeline
    summarizer = pipeline("summarization")

    # Generate the summary with specified length constraints
    summary = summarizer(text, max_length=150, min_length=30, do_sample=False)

    # Return the summarized text from the result
    return summary[0]['summary_text']

# Function to summarize text extracted from multiple files
def summarize_extracted_text(extracted_text):
    # Dictionary to hold the summarized text for each file
    summarized_text = {}

    # Iterate over each file and its extracted text
    for file_name, text in extracted_text.items():
        summarized_text[file_name] = summarize_text(text)

    return summarized_text
```

**Output**

```
Extracted text :  ['/content/output_images/page_2.png', '/content/output_images/page_1.png']
Summary of /content/output_images/page_2.png:
 Rule-based, scripted and neural network-based chatbot chatbots are among the most realistic challenges to AI chatbots . Google's
-------------------------------------------------
Summary of /content/output_images/page_1.png:
 Evolution of Conversational AI represents the journey from rudimentary rule-based chatbots to advanced systems capable of understa
-------------------------------------------------
```

# VII) GRAMMAR CHECK

- The Grammar Check feature leverages Natural Language Processing (NLP) techniques to ensure accurate and context-aware text corrections. It begins with extracting text from PDF documents using the `upload_pdf_and_convert` function.
- Next, a transformer-based model from Hugging Face's Transformers library, specifically the `pszemraj/flan-t5-large-grammar-synthesis` model, is used for grammar correction. The model, initialized via the `pipeline` function, processes the input text and generates a grammatically corrected version, which is retrieved through `results[0]['generated_text']`.

```
Grammar Check option chosen

-> Please upload a PDF file.
 Choose Files   Grammar check (3).pdf
• Grammar check (3).pdf(application/pdf) - 13274 bytes, last modified: 2/17/2024 - 100% done
Saving Grammar check (3).pdf to Grammar check (3).pdf
Uploaded Grammar check (3).pdf
PDF pages converted to images in the folder: output_images

Corrected text :  [{'generated_text': 'He is dancing. It is raining.'}]
```

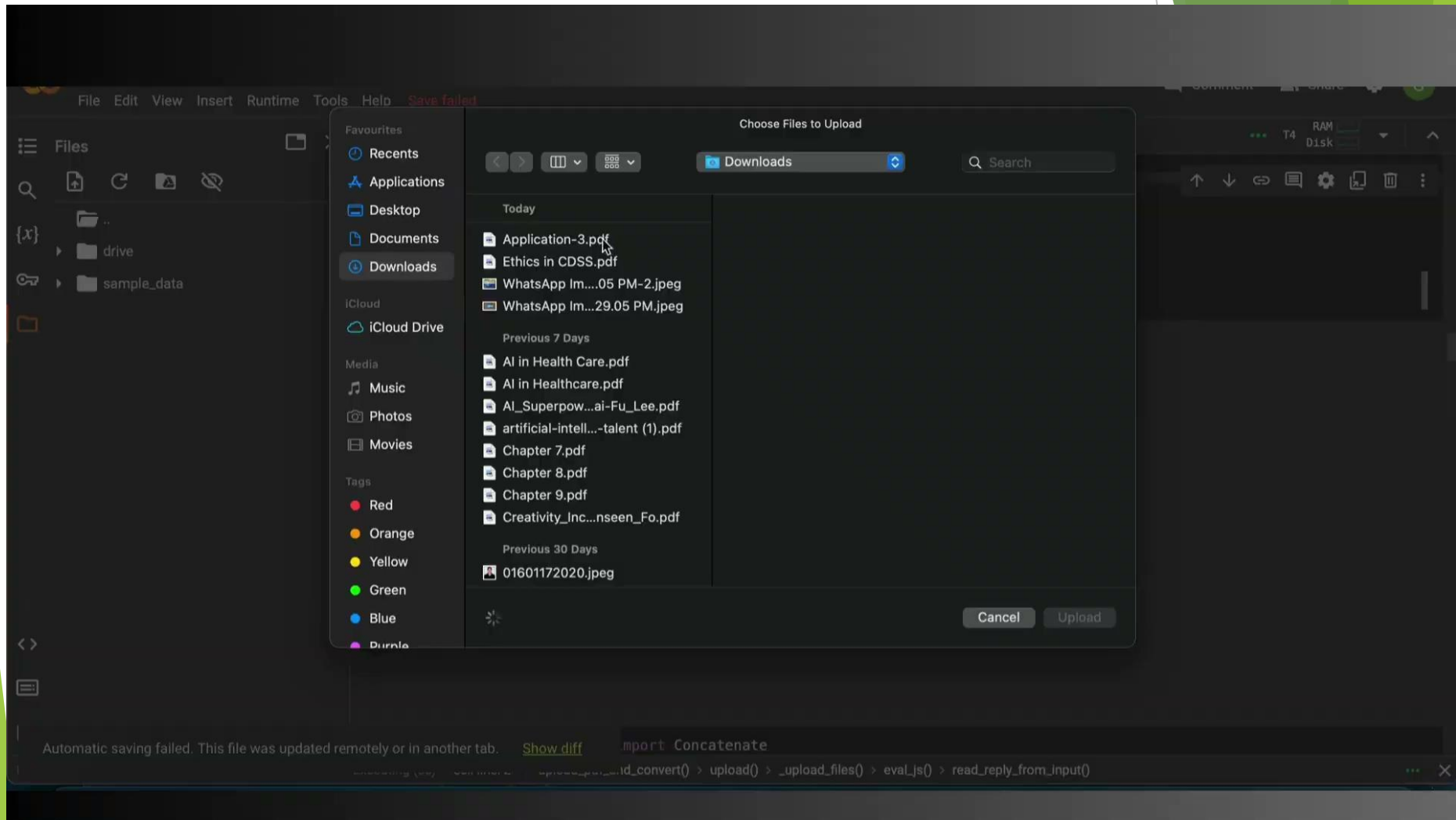Original: I is dancing. It are raining.

# EXPERIMENTAL RESULTS

**Comparative Analysis of Grammar Check Libraries**

| Library | LanguageTool | Hugging Face's Transformers Pipeline |
|---|---|---|
| Accuracy | 0.93 | 0.96 |

# CONCLUSION AND FUTURE WORK

- Different CNN designs have been evaluated and each displayed unique performance traits. Our custom CNN, with only minimal data enhancement, struck a compromise between efficiency and intricacy. Additional layers improved the ability to generalize performance of VGG16.

- Furthermore, our interactive Q&A system yielded reliable outcomes with high confidence levels.

- To ascertain document genuineness with most possible accuracy, we have fused up-to-date image processing techniques with the established techniques of comparing template shapes. Our project to develop AI-powered document intelligence has made certain that its templates are matched, answering interactive questions and document classification.

- Subsequent development stages of AI-driven document intelligence will major on enhancing the accuracy levels of our models by training them on wider data sets as well as applying more advanced classifiers.

# DEMO VIDEO

# REFERENCES

[1] Sammed S. Admuthe , Hemlata P. Channe, 2021, Document Image Classification using Visual and Textual Features, INTERNATIONAL JOURNAL OF ENGINEERING RESEARCH & TECHNOLOGY (IJERT) Volume 10, Issue 09 (September 2021),

[2] P. Dhakal, M. Munikar and B. Dahal, "One-Shot Template Matching for Automatic Document Data Capture," 2019 Artificial Intelligence for Transforming Business and Society (AITB), Kathmandu, Nepal, 2019, pp. 1-6, doi: 10.1109/AITB48515.2019.8947440.

[3] Yenigalla, Harshitha, Bommareddy Revanth Srinivasa Reddy, Batta Venkata Rahul, and Nannapuraju Hemanth Raju. "Similar Document Template Matching Algorithm." arXiv preprint arXiv:2311.12663 (2023).

[4] Li, Xiaoxi, Yujia Zhou, and Zhicheng Dou. "UniGen: A Unified Generative Framework for Retrieval and Question Answering with Large Language Models." In Proceedings of the AAAI Conference on Artificial Intelligence, vol. 38, no. 8, pp. 8688-8696. 2024.

[5] Zhuang, Yuchen, Yue Yu, Kuan Wang, Haotian Sun, and Chao Zhang. "Toolqa: A dataset for llm question answering with external tools." Advances in Neural Information Processing Systems 36 (2024).

# ACCEPTANCE LETTER

## ICAMC- 2024 Acceptance Notification for Paper ID-227

**Microsoft CMT** <email@msr-cmt.org>                          Fri, Apr 26, 2:38 PM
to me ▾

Dear Kanika Kanojia

Paper ID: 227
Title: AI-Driven Document Intelligence: A Comprehensive Approach to Classification, Q&A Handling, and Fraud Detection

We are glad to inform you that your manuscript has been accepted for the presentation in conference ICAMC 2024 and for publication in proceeding/journal.

Also, we would like to inform you that your manuscript has been reviewed and the comments from reviewers are at the bottom of this e-mail. Please incorporate the reviewers stated concerns and update the revised paper (in Camera ready paper, both word and pdf)

You are requested to send the following documents in attachment at ehsan.asgar@hmritm.ac.in in a Zipped file (Paper ID Number as a file name).
1.      Camera Ready Paper after incorporating the reviewer comments (pdf and word, both files).  File Name:  Paper ID Number_camera.doc and Paper ID Number_camera.pdf
2.      Copyright Form (Paper ID Number_CTP.pdf).
        (Download from https://docs.google.com/document/d/1wxuXXoF96fJT_nVmchjqRXzyZmUeBJJp/edit)
3.      Receipt of Registration Fee paid. File Name: (Paper ID Number _Fee receipt.pdf).
4.      Response to Reviewers (word file). File Name:  Paper ID Number _response.doc


Please follow the given MS Word template for Camera Ready Paper.
https://icamc-2024.vercel.app/guidelines

# RESEARCH PAPER [(link)](link)

## AI-Driven Document Intelligence: A Comprehensive Approach to Classification, Q&A Handling, and Fraud Detection

Sonanshi Goel[1], Shambhavi Rai[1], Princy Singhal[1], Mehak Aggarwal[1], Kanika Kanojia[1], *Ritika Kumari[1,2], Poonam Bansal[1]

[1]Department Of Artificial Intelligence and Data Sciences, IGDTUW, Delhi, 110006, India.
[2]USICT, Guru Gobind Singh Indraprastha University, Dwarka, New Delhi, India.

Contributing authors: sonanshig02@gmail.com; shambhavi.rai1604@gmail.com; princysingla11@gmail.com; mehakagg1313@gmail.com; kanikakj07@gmail.com; *ritikakumari@igdtuw.ac.in; poonambansal@igdtuw.ac.in;

### Abstract

With the growth in dependency of electronic verification systems used both by government and private organizations, there has been an increase in usage of digitized manuscripts. This demands for user friendly and efficient document interaction tools. Optimisation of the document management process has become AI driven with features like identity document classification, authenticity checking algorithms and interactive question answering. If we effectively identify the document type, we can streamline workflows and enable secure verification processes.The project employs Deep Learning techniques to classify documents like PAN cards, Aadhar cards etc. Also we employ a Similar Document Template Matching Algorithm which can seamlessly extract features, match the template with a real document and thereby detect fraud. Additionally LLMs have been leveraged to enable interactive question answering over documents. Experimental findings in our study showed the effectiveness of our work in improving efficiency, accuracy and usability in document management.

**Keywords:** CNNs for Document classification, LLMs for Question Answering, Langchain, Similar Document Template Matching Algorithm, VGG16, Fraud Detection,ORB

# AWARD FOR BEST PAPER PRESENTATION

# THANK YOU