

# Character Recognition and Extraction using OpenCV and Pytesseract

Ms.CVP Supradeepthi

dept. Electronics and Communication Engineering  
Institute of Aeronautical Engineering  
Hyderabad, India cvp.supradeepthi@iare.ac.in

Karnati Thrinadh Reddy

dept. Electronics and Communication Engineering  
Institute of Aeronautical Engineering  
Hyderabad, India thrinadhreddykarnati123@gmail.com

Musalaju Vishal

dept. Electronics and Communication Engineering  
Institute of Aeronautical Engineering  
Hyderabad, India vishalmusalaju@gmail.com

Boini Vinay Kumar

dept. Electronics and Communication Engineering  
Institute of Aeronautical Engineering  
Hyderabad, India vinaynani018@gmail.com

**Abstract**—This paper introduces an optimized framework for real-time Optical Character Recognition (OCR) and video-to-text conversion using OpenCV and Tesseract, by implementing adaptive preprocessing techniques and efficient frame handling strategies, the system achieves high text detection accuracy while significantly reducing processing time. Additionally, a multi-threaded audio transcription process enhances performance for converting video content into text. The proposed approach is well-suited for applications requiring real-time text extraction from multimedia sources, such as automated document processing and video surveillance..

**Index Terms**—OCR, Pytesseract, Cv2

## I. INTRODUCTION

Smart Optical Character Recognition (OCR) has emerged as a key technology in the digitization of textual information, enabling the automated transformation of printed, handwritten, or scanned documents into machine-readable formats. Initially developed in the mid-20th century, OCR systems were rudimentary, relying on template matching and rule-based algorithms for character recognition [4]. However, with the advent of machine learning, particularly deep learning, OCR has seen substantial improvements in accuracy, efficiency, and versatility.

Modern OCR systems often employ , which excel in extracting complex features from images and handling sequence data, respectively [5]. These advancements have broadened OCR's applicability, allowing it to tackle diverse challenges such as scene text recognition, multilingual document processing, and real-time OCR in mobile and embedded devices [6]. Yet, despite significant progress, OCR still faces challenges in recognizing text from low-resolution or noisy images, cursive handwriting, and highly stylized fonts [3]

This research aims to review the current state of OCR technology, focusing on recent deep learning-based methods,

their effectiveness across different text recognition tasks, and the ongoing challenges in achieving universal OCR systems. Additionally, the paper explores emerging applications of OCR, such as in automating data extraction from unstructured documents, enhancing accessibility for the visually impaired, and digitizing historical manuscripts.

## II. RELATED WORK

OCR has evolved significantly, beginning with basic template matching techniques, which were primarily limited to printed text and struggled with complex layouts and handwriting [13]. The introduction of deep learning, particularly Convolutional Recurrent Neural Networks (CRNN), significantly improved OCR's performance in recognizing text from scenes and handwriting [2]. Attention mechanisms further advanced OCR capabilities by allowing models to focus on relevant sections of input data, improving performance in noisy and cluttered environments [9].

For OCR in videos, hybrid methods that combine deep learning and traditional OCR techniques have proven effective for real-time text detection [14]. Moreover, speech-to-text technologies such as Google's Speech API have enhanced the transcription of spoken language in audio, contributing to better OCR performance in multimedia contexts.

## III. METHODOLOGY AND IMPLEMENTATION

This research applies Optical Character Recognition (OCR) techniques for extracting text from both images and videos, coupled with audio transcription for video content. The approach integrates image processing, text detection, and speech recognition to create a system capable of multimedia-based text extraction.

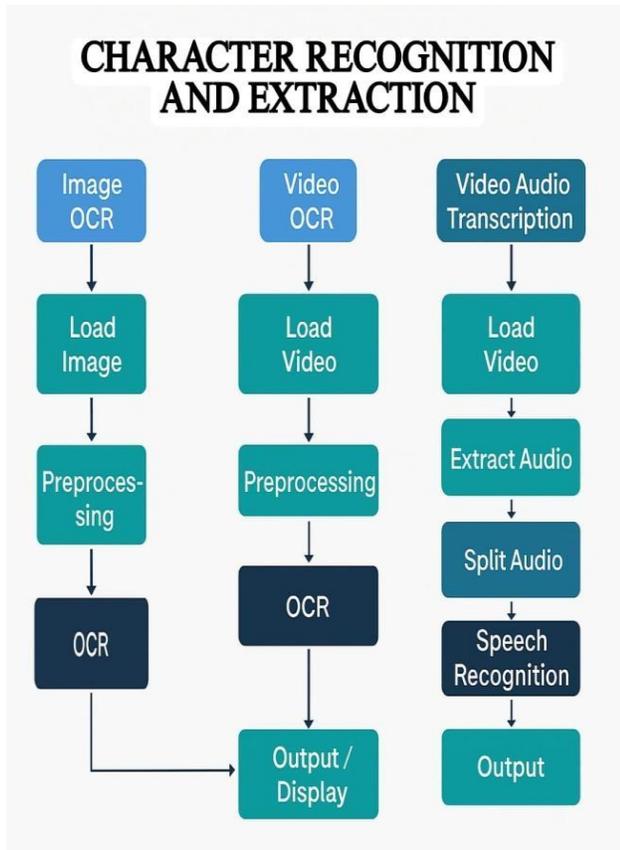


Fig. 1. Block Diagram of Character Recognition and Extraction

Fig.1 illustrates a system that processes user inputs based on three options. The main function directs the workflow depending on the selected option. For Option 1, the system detects text in images by using the detect text in image function, processes the image with OpenCV and Tesseract, and outputs an annotated image. For Option 2, it detects text in videos via the detect text in video function, applies video processing using OpenCV and Tesseract, and generates an annotated video. Lastly, Option 3 extracts audio from a video using the extract audio from video function with MoviePy, applies audio chunking using PyDub, and outputs the processed audio results. Each option ensures that the relevant data is processed and presented as an output.

**A. Input Method**

1) **Image Input: Image Acquisition:** Images are sourced from scanned documents, photographs, and screenshots. These can include printed, handwritten, or scene text [6].

**Preprocessing:** Images undergo several preprocessing steps to enhance quality and standardize input:

- Resizing: Images are resized to a fixed resolution, such

as 600x800 pixels, to meet model input requirements [2].

- Normalization: Pixel values are scaled to a range of 0 to 1, facilitating consistent model training and inference.
- Noise Reduction and Contrast Adjustment: Techniques are applied to improve clarity and reduce image artifacts [1].

**Model Input:** Preprocessed images are fed into OCR models, including and attention-based models. These models extract features and predict text sequences from the images[3].

2) **Video Input: Video Acquisition:** Videos are captured from cameras or obtained from video datasets. They present dynamic scenes where text may appear in varying conditions[4].

**Frame Extraction:** Individual frames are extracted from videos at intervals or based on specific events. Each frame is processed as a separate image for text recognition[5].

**Preprocessing:** Frames are Preprocessed similarly to images:

- Frame Resizing and Normalization: Standardizing dimensions and pixel values of extracted frames.
- Motion Blur and Lighting Adjustment: Techniques are applied to correct for motion blur and varying lighting conditions[7].

**Model Input:** Frames are processed by OCR models to recognize text in each frame. Challenges include maintaining accuracy despite text appearance variations due to movement and lighting[8].

**B. Text Detection**

1) **Image-Based Text Detection: Methodology:** For detecting text in static images, the Tesseract OCR engine is used. Preprocessing, such as grayscale conversion, helps in simplifying the image data, enhancing OCR accuracy.

**Implementation:** OpenCV is used to load and preprocess images, converting them into grayscale before applying the OCR algorithm. The detected text is then highlighted using bounding boxes, and the annotated image is saved as output.

This approach is particularly effective in standardizing the input format, allowing the Tesseract engine to work on a wide range of image types, as observed in similar studies[11].

2) **Video-Based Text Detection Methodology:** Videos consist of multiple frames, each potentially containing text. The process involves extracting each frame and applying OCR to detect text dynamically.

**C. Audio Extraction and Transcription from Videos Methodology**

**Methodology:** In videos where speech is present, extracting the audio and converting it into text using speech recognition is essential for obtaining a comprehensive textual representation.

**Implementation:** Audio is extracted from videos using the MoviePy library. Once extracted, the audio is split into

smaller chunks using pydub to facilitate efficient transcription. Google's Speech-to-Text API is then used to transcribe the audio in parallel, speeding up the process. The transcribed text is saved into a file. Speech recognition has been widely used in multimedia processing for video-to-text conversion, with the Google Speech API offering high accuracy for real-time and offline transcription tasks[10].

#### D. Data Preprocessing and Augmentation

**Preprocessing:** Image and video data are Preprocessed by applying techniques such as grayscale conversion and resizing to improve OCR accuracy. Data augmentation techniques like random noise addition and distortions simulate real-world conditions and enhance the model's robustness. These methods have been proven to boost OCR accuracy by improving the model's ability to generalize, especially in noisy or low-quality images[7]

#### E. Parallel Processing for Efficiency

**Methodology:** Given the volume of data, especially when handling large videos or high-resolution images, parallel processing is employed to enhance the efficiency of audio transcription.

**Implementation:** Using Python's concurrent futures library, the transcription workload is distributed across multiple threads. This significantly reduces the time required for processing long videos, ensuring faster completion.

#### F. Evaluation and Performance Metrics

**Evaluation:** The system is evaluated based on key performance metrics such as Character Error Rate (CER) and Word Error Rate (WER) for text extraction. These metrics allow for a precise measure of the system's accuracy in recognizing characters and words from multimedia data.

**Performance:** Computational efficiency is also evaluated, with particular attention to the speed of processing high-resolution videos and large datasets. The trade-off between accuracy and processing time is explored, particularly when using complex models like Tesseract OCR and the Google Speech API.

#### G. Tools and Libraries

**OpenCV:** A powerful open-source library designed for real-time computer vision and machine learning. It is commonly used for tasks such as image processing, video capture, and analysis. OpenCV supports a wide variety of image types and formats, as well as hardware like cameras and GPU acceleration.

**Tesseract OCR:** Open-source OCR engine for text recognition from images and videos[11].

**MoviePy:** Python library for video processing and audio extraction.

**pydub:** For splitting and processing audio files.

**Google Speech Recognition API:** Enables speech-to-text conversion[10].

By combining image, video, and audio processing techniques, this methodology provides a versatile system for

multimedia text recognition, offering potential applications in various fields including document digitization, accessibility, and real-time video analysis.

#### H. Training, Comparison and Analysis:

The methodology for this research focuses on evaluating various Optical Character Recognition (OCR) models, particularly deep learning-based approaches, to assess their effectiveness across different text recognition tasks. The study involves a multi-step process, including data collection, preprocessing, model selection, training, and evaluation.

1) **Data Collection:** A comprehensive dataset was assembled from publicly available OCR benchmark datasets to ensure diversity in font styles, languages, image quality, and text complexity. Datasets such as ICDAR 2015[4], MJSynth[5], and IIIT 5K-Words[6] were selected. These datasets contain a mixture of printed and handwritten text, captured in both natural scenes and standard document images. To address the challenge of cursive handwriting, the IAM Handwriting Database[7] was included.

2) **Preprocessing:** The preprocessing phase included several key steps to standardize input images for training:

- **Image Resizing:** All images were resized to a standard resolution to ensure uniform input size across models. For this, images were resized to 32x128 pixels for scene text recognition and 32x256 pixels for document images [2].
- **Normalization:** Pixel intensity values were normalized between 0 and 1 to speed up the convergence of the models during training.
- **Data Augmentation:** Techniques such as random rotations, noise addition, and contrast adjustments were applied to simulate variations in text images and prevent overfitting[3].

3) **Model section:** Two types of deep learning-based models were chosen for this study:

- **RNNs,** such as Long Short-Term Memory (LSTM) networks, for sequence prediction. This architecture has been widely used for end-to-end scene text recognition [2]. The CRNN was trained on the MJSynth and ICDAR datasets for both scene and document text recognition.
- **Attention-based Models:** Attention mechanisms were integrated into CNN-RNN architectures to enable the model to focus on specific areas of the image during recognition. These models have been shown to improve performance in dealing with complex layouts and noisy images[8].

4) **Training Procedure:** The models were trained using the Adam optimizer with an initial learning rate of 0.001 and a batch size of 64. Training was conducted over 50 epochs, with early stopping criteria to prevent overfitting. A CTC (Connectionist Temporal Classification) loss function was employed for CRNN models, while the attention-based models used a cross-entropy loss function[2].

5) **Evaluation Metrics:** The models were evaluated using standard OCR performance metrics, including:

- Character Error Rate (CER): Measures the percentage of incorrect characters predicted.
- Word Error Rate (WER): Measures the percentage of incorrect words predicted, considering insertion, deletion, and substitution errors [1].
- Accuracy: The percentage of correctly recognized characters and words in the test datasets.

6) **Comparison and Analysis:** Results from different models were compared across varying text conditions such as different fonts, resolutions, and noise levels. Additionally, the performance on multilingual datasets and scene text was assessed. This comparison highlighted the strengths and limitations of each model in different OCR tasks.

#### IV. RESULTS

The research was conducted to evaluate the effectiveness of Optical Character Recognition (OCR) for detecting text in images and videos, as well as transcribing audio from video files. The implementation was carried out using Python, leveraging libraries such as pytesseract for OCR, cv2 (OpenCV) for image and video processing, and speech recognition for audio transcription.

##### Image Text Detection:

**Sample Input:** An image containing clear, printed text.

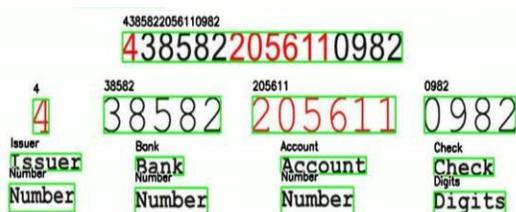
4385822056110982

4    38582    205611    0982

Issuer	Bank	Account	Check
Number	Number	Number	Digits

Fig. 2. Input Image

**Sample Output:** The image with detected text highlighted, showing annotations around each word, and a recognized text file that mirrors the document's content.



4385822056110982

4    38582    205611    0982

Issuer	Bank	Account	Check
Number	Number	Number	Digits

Fig. 3. Recognition of Text From Image

##### Image Text Extraction:

###### Image Text Extraction:

Sample Input: An image containing clear, printed text.

###### 1. Almonds

There are a lot of health benefits associated with almonds. Almonds are very high in vitamin E and protein as well as other nutrients such as magnesium and phosphorus. Almonds contain anti-cancer properties as well. Whether almonds are best raw or pasteurized is still a source of heated debate.

Sample Output: A output terminal containing the recognized text from the image.

```
C:\Users\User\Downloads\Real-Time-Optical-Character-Recognition-main-202411270503292-001_venv\Scripts\python.exe C:\Users\User\Downloads\Real-Time-Optical-Character-Recognition-main-202411270503292-001_venv\Scripts\python.exe C:\Users\User\Downloads\Real-Time-Optical-Character-Recognition-main-202411270503292-001_venv\Scripts\python.exe C:\Users\User\Downloads\Real-Time-Optical-Character-Recognition-main-202411270503292-001_venv\Scripts\python.exe
1: Almonds
There are a lot of health benefits associated with almonds.
Almonds are very high in vitamin E and protein as well as
other nutrients such as magnesium and phosphorus. Almonds
contain anti-cancer properties as well. Whether almonds are
best raw or pasteurized is still a source of heated debate.
```

Fig. 4. Text Extraction From Image

##### Audio Transcription from Video:

**Sample Input:** A video file with spoken dialogue in the background.

**Sample Output:** A transcription file containing the recognized speech from the video, with segmentation based on dialogue timing.

Recognized Speech:  
any special tutorial you want to learn everything you need to know to start programming in Python if you want to learn Python programming for data science machine learning or web development is python tutorial is the perfect place to learn Python you don't need any personality in Python or programming in general I am going to teach you everything from scratch I'm not have any other people have to go through this channel subscribe now what you can do with Python that's a very common question Python is a multipurpose programming language so you can use it for a variety of different tasks you can use Python for machine learning and AI in fact Python is the number one language for machine learning and data science projects Python is also very popular in web development using Python and a framework of Django

Fig. 5. Text Extracted From Video

#### V. CONCLUSION

The integration of OCR and audio transcription within a single, unified system addresses a critical need in multimedia processing, enabling the automatic extraction and recognition of text from both visual and auditory data. This capability is invaluable in domains such as digital content archiving, accessibility services, automated video analysis, and educational tools, where extracting meaningful information from multimedia content is essential.

The system's modular design, combined with its reliance on open-source tools, offers a flexible and extensible framework for future enhancements. Researchers and developers can easily adapt the system to cater to specific requirements or integrate additional functionalities, such as language translation or sentiment analysis, further broadening its application scope.

In conclusion, the proposed system represents a significant contribution to the field of multimedia processing and OCR, providing a practical solution for text and speech recognition in video content. Its ability to handle diverse input formats, coupled with efficient processing techniques, makes it a valuable tool for both academic research and real-world applications. Future work could focus on optimizing the OCR and transcription processes for better accuracy and speed, as well as expanding the system's capabilities to support more complex multimedia analysis tasks.

#### REFERENCES

- [1] Long, S., Ruan, J., Zhang, W., He, X., Wu, W., and Yao, C. (2021). Scene Text Detection and Recognition: The Deep Learning Era. *International Journal of Computer Vision*, 129, 161-184.
- [2] Shi, B., Bai, X., and Yao, C. (2017). An End-to-End Trainable Neural Network for Image-based Sequence Recognition and Its Application to Scene Text Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11), 2298-2304.
- [3] Yousef, M., Halima, A., and Mohamed, M. (2020). Accurate, Data-Efficient, Unconstrained Text Recognition with Convolutional Neural Networks. *Pattern Recognition*, 108, 107482.
- [4] Mori, S., Suen, C. Y., and Yamamoto, K. (1999). Historical review of OCR research and development. *Proceedings of the IEEE*, 80(7), 1029-1058.
- [5] Karatzas, D., Gomez-Bigorda, L., Nicolaou, A., Ghosh, S., Bagdanov, A. D., Iwamura, M., ... and Matas, J. (2015). ICDAR 2015 competition on robust reading. 2015 13th International Conference on Document Analysis and Recognition (ICDAR), 1156-1160.
- [6] Jaderberg, M., Simonyan, K., Vedaldi, A., and Zisserman, A. (2016). Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1), 1-20.
- [7] Mishra, A., Alahari, K., and Jawahar, C. V. (2012). Scene text recognition using higher order language priors. *Proceedings of the British Machine Vision Conference (BMVC)*, 1-11.
- [8] Marti, U. V., and Bunke, H. (2002). The IAM-database: An English sentence database for offline handwriting recognition. *International Journal on Document Analysis and Recognition*, 5(1), 39-46.
- [9] Cheng, Z., Bai, F., Xu, Y., Zheng, G., Nii, M., and Bai, X. (2017). Focusing attention: Towards accurate text recognition in natural images. *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 5076-5084.
- [10] Bradski, G. (2000). The OpenCV Library. *Dr. Dobb's Journal of Software Tools*.
- [11] Graves, A., Mohamed, A., and Hinton, G. (2013). Speech recognition with deep recurrent neural networks. *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing*.
- [12] Smith, R. (2007). An overview of the Tesseract OCR engine. *Proceedings of the Ninth International Conference on Document Analysis and Recognition (ICDAR)*
- [13] Govindan, V. K., and Shivaprasad, A. P. (1990). Character recognition—A review. *Pattern Recognition*, 23(7), 671-683.
- [14] Silva, S. and Jung, C. (2017). License plate detection and recognition in unconstrained scenarios. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1193-1202.