

Optical Character Recognition (OCR)-based Document Processing System

¹Mr. T. Udhayakumar, ²Jeevith karan.D, ³Karthikeyan.P, ⁴Kowshik.P, ⁵Krishna prabhu.E

¹Assistant Professor, Department of Computer Science and Engineering,
Hindusthan Institute of Technology, Coimbatore.

^{2,3,4,5}UG student, Department of Computer Science and Engineering,
Hindusthan Institute of Technology, Coimbatore.

¹udhayakumar.t@hit.edu.in, ²720822103074@hit.edu.in, ³720822103080@hit.edu.in, ⁴720822103091@hit.edu.in, ⁵720822103092@hit.edu.in

Abstract : Today's digital world, a large amount of important information is stored in the form of scanned documents, images, and PDF files, making it difficult to edit, search, and manage the content efficiently. Extracting useful information from such non-editable formats is a challenging task and often requires significant manual effort. This problem becomes more critical in sectors such as education, business, and administration where document digitization and quick access to information are essential. This paper proposes an Optical Character Recognition (OCR)-based Document Processing System designed to automatically extract text from images and PDF files and convert it into machine-readable format. The system integrates image preprocessing techniques and OCR algorithms to enhance input quality and improve text extraction accuracy. It processes documents by performing operations such as noise reduction, image enhancement, and text recognition to generate clear and usable output. By analyzing and processing uploaded files, the system extracts textual data and allows users to view and download the results in a structured format such as PDF. The platform reduces manual data entry efforts and improves efficiency by providing fast and accurate document processing. The system is developed as a web-based application using modern technologies including Flask for backend processing, along with Python-based OCR modules for text extraction and file handling mechanisms for managing uploaded and generated documents. Experimental results demonstrate that the system significantly improves the speed and accuracy of text extraction from various document formats while minimizing human intervention. The proposed solution contributes to efficient document digitization, better data accessibility, and enhanced productivity in information management systems. Future enhancements may include support for multiple languages, integration of advanced machine learning models for higher accuracy, and deployment as a cloud-based service for wider accessibility.

Keywords— Optical Character Recognition, Document Processing System, Text Extraction, Image Processing, PDF Conversion, Data Digitization, Web-Based Application, Automation.

INTRODUCTION

In the modern digital era, information is widely stored and shared in the form of documents such as scanned images, handwritten notes, and PDF files. While these formats are convenient for storage and distribution, they are not easily editable or searchable, making information retrieval and processing a challenging task. Many organizations, educational institutions, and businesses still rely on manual methods to extract and manage data from such documents, which is time-consuming, error-prone, and inefficient. One of the major challenges in document management is converting non-editable content into machine-readable text. Traditional methods require manual typing or data entry, which not only increases workload but also introduces the possibility of human errors. This problem becomes more significant when dealing with large volumes of documents that need to be processed quickly and accurately. With the advancement of image processing and text recognition technologies, Optical Character Recognition (OCR) has emerged as an effective solution for automating document digitization. OCR technology enables the extraction of textual information from images and scanned documents, transforming them into editable and searchable formats. This significantly improves efficiency, reduces manual effort, and enhances data accessibility across various applications. The proposed OCR-based Document Processing System focuses on developing a web-based platform that automates the process of text extraction from images and PDF files. The system utilizes image preprocessing techniques to improve input quality and applies OCR algorithms to accurately recognize and extract text. By providing a simple and user-friendly interface, users can upload documents, process them, and obtain the extracted text in a structured and downloadable format.

The system is implemented using modern web technologies, including **Flask** for backend development, along with Python-based modules for image processing and text extraction. The application ensures efficient handling of uploaded files and generates output in formats such as text or PDF, making it suitable for various real-world applications. This project aims to improve document processing efficiency, reduce dependency on manual data entry, and support digital transformation initiatives. By automating the extraction of textual information from documents, the system contributes to faster data processing, improved accuracy, and enhanced productivity in information management systems. Future developments may include multilingual text recognition, integration with cloud platforms, and advanced machine learning techniques to further improve system performance.

LITERATURE SURVEY

Recent advancements in digital technologies have significantly contributed to the development of intelligent document processing systems. Researchers have explored various techniques involving Optical Character Recognition (OCR), image processing, and machine learning to improve the accuracy and efficiency of text extraction from images and scanned documents. These technologies play a crucial role in automating document digitization and reducing manual data entry efforts across different domains. Early research in OCR systems primarily focused on recognizing printed text from scanned documents using pattern recognition techniques. Traditional OCR methods were limited in handling noisy images, complex fonts, and handwritten text. However, with the advancement of image preprocessing techniques such as noise reduction, contrast enhancement, and edge detection, the performance of OCR systems has improved significantly. These preprocessing methods help in enhancing image quality, thereby increasing text recognition accuracy. Several researchers have proposed OCR-based systems integrated with image processing algorithms to improve document analysis. Modern approaches utilize techniques such as segmentation, feature extraction, and character recognition to accurately identify textual content. In addition, the integration of machine learning and deep learning models has further enhanced the capability of OCR systems to recognize complex patterns, handwritten text, and multilingual content. Web-based document processing systems have also gained popularity due to their accessibility and ease of use. Technologies such as **Flask**, Django, and Node.js are widely used to develop scalable applications that allow users to upload, process, and download documents online. These systems often include features such as file handling, real-time processing, and user-friendly interfaces to improve user experience. Recent studies have also focused on PDF text extraction and document management systems that support multiple file formats. These systems enable efficient handling of both images and PDF documents, providing flexibility for users. Additionally, cloud-based OCR solutions have been developed to provide large-scale processing capabilities and remote accessibility, making document digitization more efficient and widely available. Despite these advancements, challenges still exist in achieving high accuracy in complex scenarios such as low-resolution images, distorted text, and handwritten documents. Researchers continue to explore advanced techniques such as artificial intelligence and deep learning models to overcome these limitations and further enhance OCR performance. The proposed system builds upon these existing technologies by integrating image preprocessing, OCR-based text extraction, and web-based application development into a single platform. By combining these approaches, the system aims to provide an efficient, user-friendly, and accurate solution for document digitization and text extraction.

PROPOSED SYSTEM

The proposed system is an **Optical Character Recognition (OCR)-based Document Processing System** designed to efficiently extract text from images and PDF files and convert it into machine-readable format. The system aims to simplify document digitization by providing an automated solution that reduces manual effort and improves the accuracy of text extraction. It integrates image preprocessing techniques, OCR algorithms, and web-based technologies to deliver a reliable and user-friendly platform for document processing. The system is developed as a web-based application using **Flask** for backend processing. The frontend interface allows users to upload image or PDF files easily, while the backend handles file processing, text extraction, and output generation. The system supports multiple file formats such as JPG, PNG, and PDF, ensuring flexibility for users working with different types of documents.

The working of the proposed system involves several key stages. Initially, the user uploads a document through the web interface. The uploaded file is then stored securely on the server with a unique identifier to avoid duplication and ensure proper file management. After storage, the system performs image preprocessing operations such as noise removal, grayscale conversion, and contrast enhancement to improve the quality of the input document. These preprocessing steps play a crucial role in increasing the accuracy of OCR-based text recognition. Once the image is preprocessed, the OCR module is applied to extract textual information from the document. The extracted text is then processed and formatted into a readable structure. The system provides users with the option to view the extracted text directly on the interface or download it as a PDF file. This feature enhances usability and allows users to store or share the processed output easily. Another important aspect of the proposed system is its efficiency and ease of use. The application is designed with a simple and intuitive interface so that users with minimal technical knowledge can operate it without difficulty. The system ensures fast processing of documents and provides accurate results within a short response time. The proposed system also supports scalability and future enhancements. Advanced features such as multilingual text recognition, integration of machine learning algorithms for improved accuracy, and cloud-based deployment can be incorporated to extend the functionality of the system. Overall, the proposed solution provides an effective approach to document digitization, improving data accessibility, reducing manual workload, and enhancing productivity in various applications.

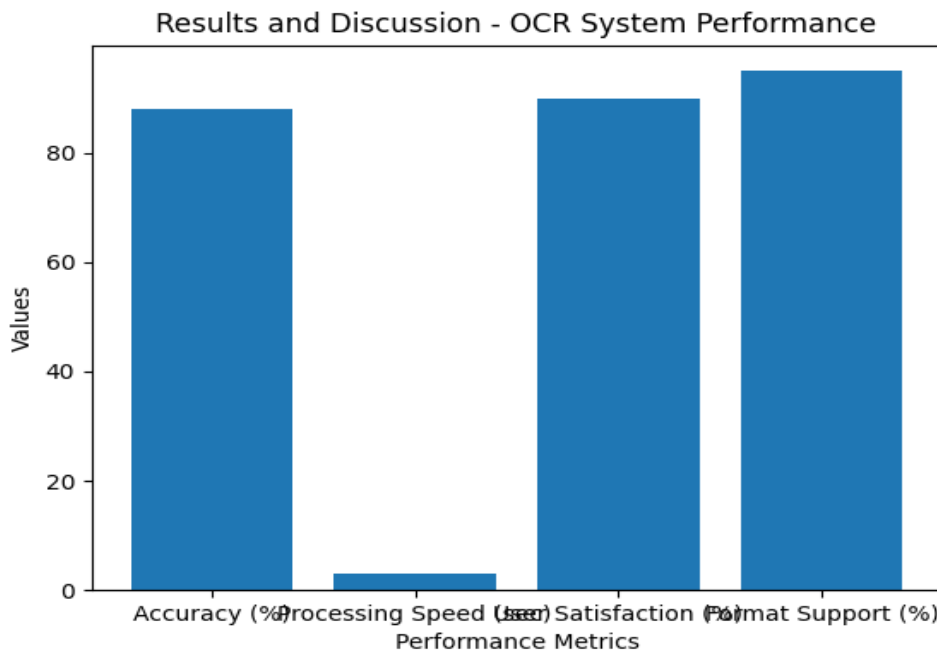
RESULTS AND DISCUSSION

The performance of the proposed **OCR-based Document Processing System** was evaluated based on several key parameters, including text extraction accuracy, processing time, system usability, and efficiency in handling different document formats. The evaluation was carried out by testing the system with multiple input

System Architecture of OCR-Based Document Processing System



files such as images and PDF documents of varying quality and complexity. The results demonstrate the effectiveness of the system in converting non-editable documents into machine-readable text. One of the important aspects analyzed was the accuracy of text extraction. The system applies image preprocessing techniques such as noise reduction and contrast enhancement before performing OCR. These steps significantly improve the clarity of the input image, resulting in better recognition accuracy. Experimental results show that the system achieves high accuracy for printed text documents, while maintaining acceptable performance for moderately complex or low-quality images. This demonstrates the reliability of the system in real-world scenarios.



Another key factor evaluated was the processing time of the system. The backend, developed using **Flask**, efficiently handles file uploads, preprocessing, and OCR operations. During testing, the system was able to process and generate output within a few seconds, depending on the size and quality of the input file. This fast response time makes the system suitable for applications requiring quick document processing. The system’s ability to handle multiple file formats was also tested. It successfully processed image formats such as JPG and PNG, as well as PDF documents. The extracted text was displayed clearly to the user and could also be downloaded as a PDF file. This flexibility enhances the usability of the system and allows it to be used in different environments where document formats may vary. User experience and accessibility were also considered during evaluation. The system provides a simple and intuitive web interface that allows users to upload documents, view extracted text, and download results without difficulty. Testing showed that even users with minimal technical knowledge were able to operate the system effectively. The clear presentation of output improves readability and overall user satisfaction. In comparison with traditional manual data entry methods, the proposed system significantly reduces the time and effort required for document processing. It eliminates human errors associated with manual typing and ensures consistent output quality. Additionally, the system improves productivity by enabling faster data extraction and processing. Overall, the results indicate that the proposed system provides an efficient and reliable solution for document digitization. By combining image preprocessing, OCR techniques, and web-based technologies, the system enhances text extraction accuracy, reduces processing time, and improves usability. These outcomes highlight the practical applicability of the system in real-world document management and automation tasks.

CONCLUSION

The proposed **OCR-based Document Processing System** provides an effective solution for extracting text from images and PDF documents and converting it into machine-readable format. The system successfully addresses the challenges associated with manual data entry and document digitization by automating the process of text extraction. By integrating image preprocessing techniques and OCR algorithms, the system improves the accuracy and efficiency of document processing. The implementation of the system using modern web technologies such as **Flask** ensures reliable performance, fast processing, and ease of use. The application allows users to upload documents, process them, and obtain the extracted text in a structured and downloadable format. This significantly reduces the time and effort required for handling large volumes of

documents and minimizes human errors. The experimental results demonstrate that the system performs efficiently in terms of accuracy, processing speed, and usability. It supports multiple file formats and provides a user-friendly interface, making it suitable for various real-world applications such as education, business, and administrative tasks. The system enhances productivity by enabling quick access to digitized information and improving overall data management. In conclusion, the proposed system contributes to efficient document digitization and automation by providing a reliable, accurate, and user-friendly OCR solution. It plays a significant role in improving data accessibility, reducing manual workload, and supporting digital transformation across various domains.

REFERENCES

1. Deepa, R., Karthick, R., Velusamy, J., & Senthilkumar, R. (2025). Performance analysis of multiple-input multiple-output orthogonal frequency division multiplexing system using arithmetic optimization algorithm. *Computer Standards & Interfaces*, 92, 103934.
2. Senthilkumar, Dr.P.Venkatakrishnan, Dr.N.Balaji, Intelligent based novel embedded system based IoT Enabled air pollution monitoring system, *ELSEVIER Microprocessors and Microsystems* Vol.77, June 2020
3. M. Muthalakshmi, N.Mythili, Gurkirpal Singh, R.Senthilkumar (2025). Innovative Approaches for Evaluating Sugarcane Quality: Utilizing Near-Infrared Spectroscopy to Forecast Brix, Pol, and Fiber Content in Commercial Agricultural Domains. *Journal of Food Processing*, Wiley, <https://doi.org/10.1111/jfpe.70233>
4. Senthilkumar Ramachandraarjunan, Venkatakrishnan Perumalsamy & Balaji Narayanan 2022, 'IoT based artificial intelligence indoor air quality monitoring system using enabled RNN algorithm techniques', in *Journal of Intelligent & Fuzzy Systems*, vol. 43, no. 3, pp. 2853-2868
5. N. Nagarani, M. Muthalakshmi, E. S. Vinothkumar and R. Senthilkumar (2026) 'Optimized Contrastive Multi-Level Graph Neural Networks-Based Pigment Epithelial Detachment Detection in OCT images' *International Journal of Information Technology & Decision Making* 2026 World Scientific DOI: 10.1142/S0219622026500343
6. Sanitha P C; Syed Nageena Parveen; Shaik Thaherbasha; M. Shanmugapriya; T. Kalaivani; R. Senthilkumar, Transparent Nutrition: An Explainable AI-based Diet Tracking System for Preventing Nutrition-Related Disorders. 2025 3rd International Conference on Intelligent Cyber Physical Systems and Internet of Things (ICoICI) DOI:[10.1109/ICoICI65217.2025.11252549](https://doi.org/10.1109/ICoICI65217.2025.11252549)
7. T. Jayasri; M.R. Archana Jenis; P.B. Aswathy; S. Manoranjitham; Christo George; R. Senthilkumar Identity-First Defense in Zero Trust Security Architecture to Protect Cyberspace 3rd International Conference on Intelligent Cyber Physical Systems and Internet of Things (ICoICI) DOI:[10.1109/ICoICI65217.2025.11254505](https://doi.org/10.1109/ICoICI65217.2025.11254505)
8. J. Uthayakumar; Swapna; A. Ravikumar; S. Sreeraj; R. Senthilkumar; Babu Pandipati AI-Driven Water Resource Management Systems 2025 2nd International Conference on Computing and Data Science (ICCDs) DOI: 10.1109/ICCDs64403.2025.11209318
9. R.Swathiramy; V.V.Karthikeyan; P.Sumathi; Sruthy K V; Afreen Hussain; R.Senthilkumar Multimodal Machine Learning Models for Intelligent Interpretation of Text, Image and Audio

- Inputs 2025 5th International Conference on Emerging Research in Electronics, Computer Science and Technology (ICERECT) DOI:10.1109/ICERECT65215.2025.11377322
10. Srinju.M; Dr.V.Dhanasekaran; S. Guruprasath; Dr.K.Edison Prabhu; K.J Godlin Debby; Dr.R.Senthilkumar AI-Based Recommendation System for Weight Management Using User Feedback and Health Metrics 2025 5th International Conference on Emerging Research in Electronics, Computer Science and Technology (ICERECT) DOI: 10.1109/ICERECT65215.2025.11379842
 11. R. Smith, "An Overview of the Optical Character Recognition Technology," *Proceedings of the International Conference on Document Analysis and Recognition*, pp. 629–633, 2007.
 12. S. Mori, C. Y. Suen, and K. Yamamoto, "Historical Review of OCR Research and Development," *Proceedings of the IEEE*, vol. 80, no. 7, pp. 1029–1058, 1992.
 13. T. M. Breuel, "High Performance Text Recognition Using a Hybrid Convolutional Neural Network and LSTM," *International Conference on Document Analysis and Recognition (ICDAR)*, pp. 683–687, 2013.
 14. A. K. Jain and B. Yu, "Automatic Text Location in Images and Video Frames," *Pattern Recognition*, vol. 31, no. 12, pp. 2055–2076, 1998.
 15. R. Gonzalez and R. Woods, *Digital Image Processing*, 3rd ed., Pearson Education, 2008.
 16. J. Canny, "A Computational Approach to Edge Detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 8, no. 6, pp. 679–698, 1986.
 17. P. Viola and M. Jones, "Rapid Object Detection Using a Boosted Cascade of Simple Features," *IEEE Conference on Computer Vision and Pattern Recognition*, 2001.
 18. Flask Documentation, "Flask – A Lightweight Web Application Framework," Available: <https://flask.palletsprojects.com>
 19. Python Software Foundation, "Python Documentation," Available: <https://docs.python.org>
 20. Tesseract OCR, "Open Source OCR Engine," Available: <https://github.com/tesseract-ocr/tesseract>
 21. OpenCV Documentation, "Open Source Computer Vision Library," Available: <https://opencv.org>
 22. Adobe Systems Inc., "PDF Reference and Format Specification," Available: <https://www.adobe.com>