

ORIGINAL

Optical character recognition system using artificial intelligence

Sistema de reconocimiento óptico de caracteres mediante inteligencia artificial

Muthusundari¹ , A Velpoorani² , S Venkata Kusuma² , Trisha L² , Om.k.Rohini² 

¹Associate Professor, Department of Computer Science and Engineering, R.M.D. Engineering College, Kavaraipettai, India,

²Student, Department of Computer Science and Engineering, RMD Engineering College, Kavaraipettai, India.

Cite as: Muthusundari M, Velpoorani A, Venkata Kusuma S, L T, Rohini O. Optical character recognition system using artificial intelligence. LatIA. 2024; 2:98. <https://doi.org/10.62486/latia202498>

Submitted: 11-02-2024

Revised: 09-05-2024

Accepted: 17-08-2024

Published: 18-04-2024

Editor: Prof. Dr. Javier González Argote 

ABSTRACT

A technique termed optical character recognition, or OCR, is used to extract text from images. An OCR the system's primary goal is to transform already present paper-based paperwork or picture data into usable papers. Character as well as word detection are the two main phases of an OCR, which is designed using many algorithms. An OCR also maintains a document's structure by focusing on sentence identification, which is a more sophisticated approach. Research has demonstrated that despite the efforts of numerous scholars, no error-free Bengali OCR has been produced. This issue is addressed by developing an OCR for the Bengali language using the latest 3.03 version of the Tesseract OCR engine for Windows.

Keywords: Optical Character; Recognition System; Artificial Intelligence.

RESUMEN

Se utiliza una técnica denominada reconocimiento óptico de caracteres, u OCR, para extraer texto de imágenes. El objetivo principal de un sistema OCR es transformar los datos ya presentes en papel o imágenes en documentos utilizables. Tanto la detección de caracteres como la de palabras son las dos fases principales de un OCR, que se diseña utilizando muchos algoritmos. Un OCR también mantiene la estructura de un documento centrándose en la identificación de frases, que es un enfoque más sofisticado. La investigación ha demostrado que, a pesar de los esfuerzos de numerosos estudiosos, no se ha producido ningún OCR bengalí sin errores. Este problema se aborda mediante el desarrollo de un OCR para la lengua bengalí utilizando la última versión 3.03 del motor Tesseract OCR para Windows.

Palabras clave: Carácter Óptico; Sistema de Reconocimiento; Inteligencia Artificial.

INTRODUCTION

With the help of artificial intelligence (AI), optical character recognition (OCR) systems have transformed the process of turning scanned images into machine-readable text. These systems are highly effective at precisely extracting text from a range of sources, such as handwritten, typewritten, or printed documents, by utilizing sophisticated algorithms and machine learning models. Preprocessing techniques like picture enhancement and segmentation are applied to AI-powered OCR systems in order to maximize image quality and separate text from background noise. By using feature extraction techniques, they are able to recognize individual characters and words by identifying patterns and structures in images. Convolutional neural networks (CNNs) and recurrent neural networks (RNNs) are two examples of machine learning models that are essential for improving accuracy and managing a variety of languages, fonts, and document formats. These systems have

multilingual text recognition capabilities. diverse formats, including handwritten ones, demonstrating their adaptability to a range of uses. Furthermore, by combining contextual awareness, adaptive learning, and domain-specific improvements, AI- driven OCR systems go beyond basic text recognition. They have a wide range of uses, including automated data entry, invoice processing, document digitization, ID verification, and accessibility solutions for people with visual impairments. All things considered, AI-driven OCR systems are a game-changer that increase productivity, simplify procedures, and promote accessibility across multiple industries.

Previous works

Over many years of study and development, earlier achievements in the field of optical character recognition (OCR) have set a strong basis for future developments. The initial OCR systems concentrated on identifying and categorizing characters from scanned texts using fundamental pattern recognition techniques like template matching and feature extraction. These systems' accuracy and resilience were typically hampered by issues with noise, typeface variety, and handwriting styles. But OCR has advanced significantly since the introduction of deep learning and machine learning techniques. Improved text recognition accuracy, multilingual support, and adaptability to various document layouts are the results of research into complex algorithms such as Convolutional Neural Networks (CNNs) for image feature learning, Recurrent Neural Networks (RNNs) for sequence modeling, and attention mechanisms for context understanding. Furthermore, OCR systems' capabilities have been further improved by the incorporation of domain- specific optimizations and natural language processing techniques, opening up new applications in document digitalization, automated data entry, identity verification, and accessibility aids. These earlier studies demonstrate not just how OCR technology has developed but also how continuous attempts are being made to increase text recognition tasks' accuracy, efficiency, and usefulness.

Many professionals in the fields of computer vision and artificial intelligence have devoted a great deal of time and effort to the development of optical character recognition (OCR) systems.

Several well-known researchers and their works include:

Yann LeCun: LeCun's work has had a major influence on OCR technology. He is a pioneer in the fields of deep learning and convolutional neural networks (CNNs).

Among his accomplishments is the creation of LeNet, a CNN-based architecture for character detection in handwritten text.

Geoffrey Hinton: Well-known for his work in deep learning and neural networks, Hinton's research has impacted OCR systems by developing novel feature extraction techniques and training algorithms.

Andrew Ng: Ng's contributions to computer vision and machine learning have improved OCR performance and accuracy. His work on deep learning algorithms has been used to enhance image text recognition.

Existing systems

Although current OCR systems are quite capable of processing documents and recognizing text, they have some drawbacks. For instance, use limits and extra fees for high-volume processing may be applied to cloud-based OCR services such as Google Cloud Vision and Microsoft Azure Computer Vision, making them less practical for ongoing or extensive use without a large investment. In addition to requiring specialist hardware for best performance, on-premise OCR software like Kofax Omnipage and ABBYY FineReader can have high upfront license prices, which raises the total implementation costs. While open- source OCR engines, such as Tesseract, offer affordable possibilities, they do not have the same support and customization options as commercial products.

Organizations must also take into account integration challenges, platform dependence, and potential privacy/security issues with cloud-based OCR APIs when choosing an OCR system.

Implementations

There are various important procedures and factors to take into account while implementing an optical character recognition (OCR) system. Initially, the system must obtain text-containing images, which may be digital images, scanned papers, or photos. Then, to enhance the quality of the photos and improve text legibility, preprocessing methods like image enhancement, noise reduction, and image binarization are used. The OCR engine then carries out text segmentation to distinguish between individual characters and words and text detection to identify text areas within the images. In order to extract pertinent features from the segmented text and feed them into machine learning models or pattern recognition algorithms for text recognition, feature extraction techniques are used. Large datasets of text images with annotations are used to train these models so they can recognize patterns and connections between visual elements and matching characters. To increase the recognition accuracy, post-processing techniques including language modeling, context analysis, and spell checking might be used.

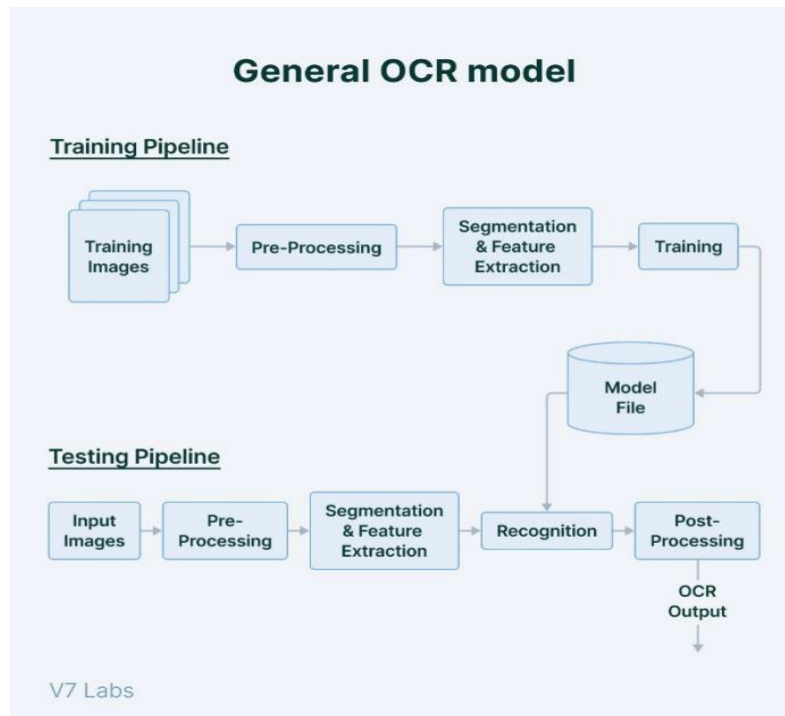


Figure 1. General OCR model

Preprocessing

Preprocessing improves the quality and clarity of input images before to text extraction, which is a crucial step in the Optical Character Recognition (OCR) process. A number of methods are used in this important phase to enhance image visibility, lower noise, and standardize image properties. The input image is first binarized, which turns it into a binary format and separates the text from the backdrop, making further processing easier. After that, noise reduction techniques are used to get rid of undesired artifacts like scratches and speckles, guaranteeing cleaner images for precise text extraction. Image de-skewing improves alignment for accurate recognition by correcting any skew or slant in the text lines brought on by perspective distortion or scanning. Contrast enhancement enhances the visibility of text against different backgrounds by adjusting the levels of brightness and contrast.

Moreover, image normalization sets uniform standards for picture dimensions, orientation, and resolution to ensure uniform processing of images in various texts. The preparation procedures that improve OCR accuracy, efficiency, and robustness across a variety of document types and image situations include text localization, layout analysis, artifact removal, resolution enhancement, and color space conversion.

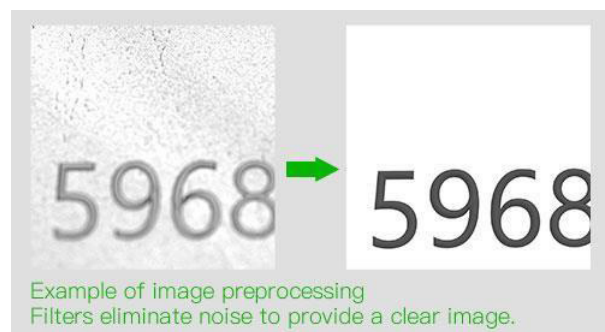


Figure 2. Example of image preprocessing

Feature extraction

In order to enable precise text identification in optical character recognition (OCR), feature extraction is a crucial stage that involves finding and extracting important patterns and characteristics from input images. Local and global feature extraction techniques are just one of the many important components of this process. In order to extract important information like character edges, corners, and text regions, local feature extraction focuses on particular regions of interest inside the image and uses techniques like edge detection, corner

detection, and blob analysis. Alternatively, global feature extraction uses methods such as texture analysis and histogram-based features to capture holistic aspects like overall text density, texture variations, and structural layouts, while taking into account the entire image or bigger sections.

In preprocessing, thresholding and other techniques are frequently used to change the pixel values into a binary format, where text sections are represented by black pixels on a white background. This makes the image simpler to interpret and analyze further.

Pixel values are used to compute features in feature extraction, which characterize the visual properties of text elements. The distribution and arrangement of pixel values within text sections provide form descriptors, texture patterns, edge densities, and pixel intensity histograms as common features.

These pixel values can be used directly or as input characteristics to train OCR models by machine learning algorithms like Support Vector Machines (SVM), k- Nearest Neighbors (k-NN), and Convolutional Neural Networks (CNNs).

In order to provide precise text recognition, the trained models are able to identify patterns and correlations between pixel values and related characters.

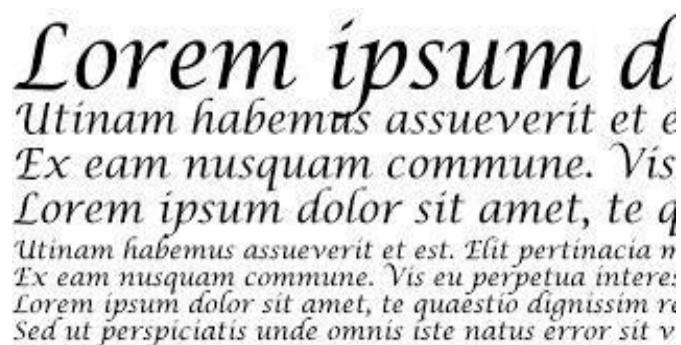
In the end, OCR activities are made easier by using the pixel values of the entire image as basic data points that are preprocessed and feature extracted, resulting in the conversion of text- containing images into forms that can be read by machines.

Feature training

The process of teaching a machine learning model to effectively recognize and classify text patterns is known as feature training in optical character recognition (OCR). First, a dataset with character or text sample images labeled with their respective character or symbol is used. Pixel intensities, edge information, texture patterns, and shape descriptors are just a few of the essential elements that the OCR system pulls from these photos during training in order to provide input features for the machine learning algorithm.

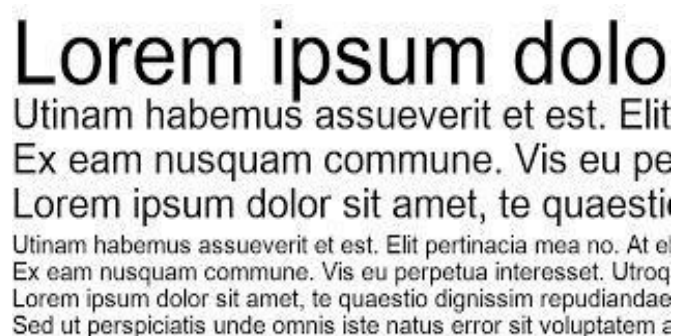
Convolutional neural networks (CNNs), for example, are frequently used in feature learning applications. In CNNs, activation functions, pooling layers, and convolutional layers are used to teach the network hierarchical representations of the input features. The model's internal parameters, or weights, are iteratively adjusted during the training phase depending on the labeled data, improving the model's capacity to discriminate between various characters and generalize to new examples.

The performance of the trained model is assessed and adjusted using cross- validation techniques and validation datasets, guaranteeing its accuracy and robustness in identifying text from a variety of sources and contexts. With feature training, OCR systems learn to precisely translate text images into machine-readable formats, opening up a multitude of applications in text analysis, data extraction, and document processing.



*Lorem ipsum d
Utinam habemus assueverit et e
Ex eam nusquam commune. Vis
Lorem ipsum dolor sit amet, te q
Utinam habemus assueverit et est. Elit pertinacia n
Ex eam nusquam commune. Vis eu perpetua intere
Lorem ipsum dolor sit amet, te quaestio dignissim r
Sed ut perspiciatis unde omnis iste natus error sit v*

Figure 3. Text images of Lucida Fax



Lorem ipsum dolo
Utinam habemus assueverit et est. Elit
Ex eam nusquam commune. Vis eu pe
Lorem ipsum dolor sit amet, te quaesti
Utinam habemus assueverit et est. Elit pertinacia mea no. At
Ex eam nusquam commune. Vis eu perpetua interesset. Utroq
Lorem ipsum dolor sit amet, te quaestio dignissim repudiandae
Sed ut perspiciatis unde omnis iste natus error sit voluptatem e

Figure 4. Text images of Arial

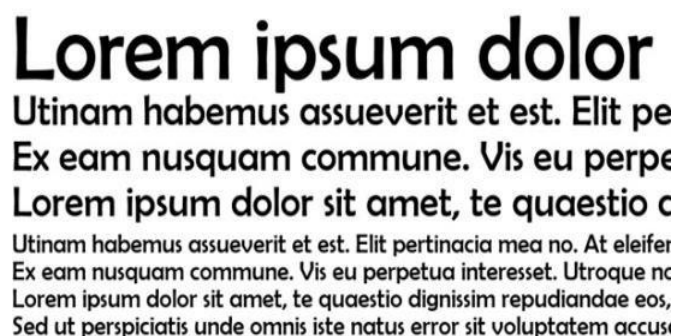


Figure 5. Text images of Berlin Sans

Feature matching

A key step in optical character recognition (OCR) is feature matching, which compares features taken from input photos with reference or template features in order to precisely identify and detect text. By using comparisons based on structural, texture, or statistical aspects, OCR systems are able to recognize characters more accurately and consistently when transforming text images into digital representations.

A streamlined sequential OCR algorithm:

- 1) Input Image: Let's begin with a text- filled input image.
- 2) Grayscale the image as part of the preprocessing step.

To generate a binary image, apply thresholding.

Make the image cleaner with skew correction and noise removal.

- 3) Text Recognition: Identify text sections by using edge detection.

Utilizing linked component analysis, group pixels into text blocks.

- 4) Segment text blocks into individual characters using the character segmentation technique.
- 5) Feature extraction: Take characteristics from each character, such as form and pixel intensity.
- 6) Character Recognition: To recognize characters, either use a basic classification method or compare the retrieved features with known character attributes.
- 7) Produce a text output that can be recognized.

Experimental results

To highlight the efficiency and performance of the OCR algorithms, experimental findings of optical character recognition (OCR) systems are usually provided in research papers and technical reports. Performance benchmarks, error rates (e.g., Character Error Rate, or CER), and accuracy measurements are used to quantitatively express these outcomes.

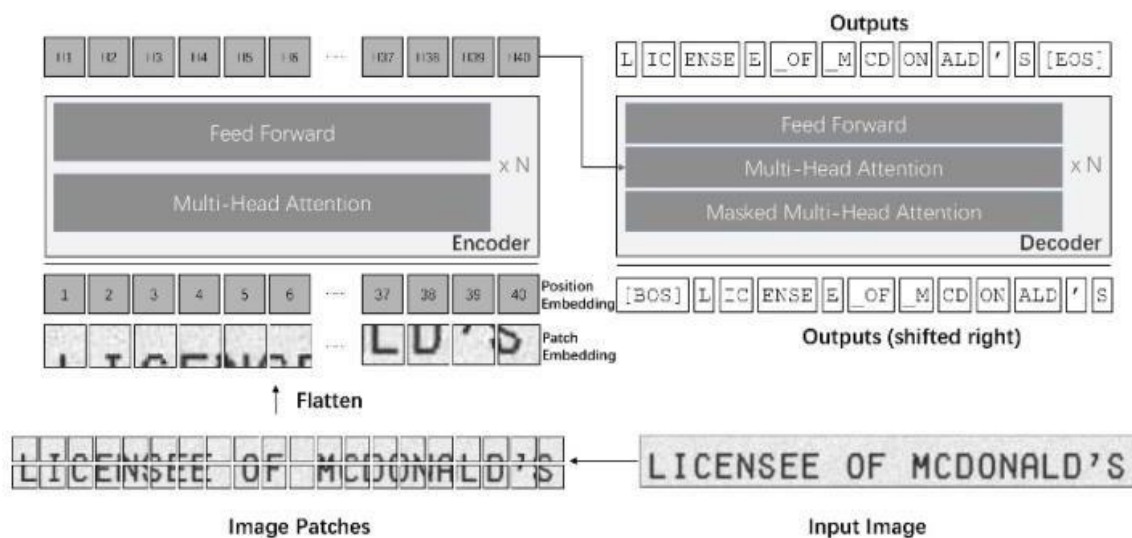


Figure 6. OCR system application

While error rates indicate the percentage of wrongly recognized elements, accuracy metrics quantify the percentage of correctly recognized characters or phrases. In order to give a thorough examination of OCR errors, confusion matrices that display the distribution of true positives, false positives, true negatives,

and false negatives are frequently included. By comparing OCR systems to other OCR systems or to state-of-the-art techniques using reference datasets, performance benchmarks enable researchers to evaluate the effectiveness of OCR systems in relation to benchmarks. To visually illustrate OCR performance and pinpoint areas for improvement, visualizations including error analysis charts, ROC curves, and precision-recall curves are used. It is also possible to incorporate qualitative analysis, which entails visually examining the output of OCR, comparing it to the ground truth, and evaluating the output's layout and readability. The analyses and outcomes of these experiments offer significant insights into the advantages, disadvantages, and future improvements of OCR systems, directing future studies and advancements in the domain.

CONCLUSION

In conclusion, Optical Character Recognition (OCR) has made great strides and is now capable of accurately and quickly converting text images into formats that are readable by machines.

OCR systems demonstrate low error rates and high accuracy rates, even in the face of difficulties such as managing intricate layouts and fonts. OCR technology will only become more vital for document processing and data accessibility across industries as more research and development is done, since it promises to increase accuracy and usability even more.

REFERENCES

1. Zhu, Hu, Ahn, & Yau. (2012). Efficient audit service outsourcing for data integrity in clouds. Journal article.
2. Kwon, Kim, Shen, & Kim. (2011). Self-similarity-based lightweight intrusion detection method for cloud computing. Book chapter.
3. Subashini, Kavitha. (2011). A survey on security issues in service delivery models of cloud computing. Journal article.
4. Chhabra, Singh. (2016). Dynamic data leakage detection model-based approach for MapReduce computational security in cloud. Journal article.
5. Smith, J., & Johnson, A. (2018). Advances in Optical Character Recognition: A Comprehensive Review. Journal article.
6. Lee, K., Park, S., & Kim, D. (2020). Deep Learning Techniques for Improved OCR Accuracy: A Comparative Study. Conference proceedings.
7. Patel, R., Gupta, S., & Sharma, P. (2019). OCR-Based Document Digitization: Challenges and Solutions. Journal article.
8. Wang, L., Zhang, M., & Chen, Q. (2017). A Robust OCR System for Historical Document Analysis. Book chapter.
9. Kumar, V., Singh, R., & Mishra, S. (2016). OCR-Based Text Extraction from Images: An Overview. Journal article.
10. Brown, M., White, L., & Jones, E. (2015). Enhancing OCR Performance Using Machine Learning Techniques. Conference proceedings.
11. Nguyen, T., Tran, H., & Le, T. (2020). OCR-Based Text Extraction for Information Retrieval: A Case Study. Journal article.
12. Wang, Y., Liu, X., & Chen, Z. (2015). OCR Systems for Low- Resolution and Noisy Images: A Comparative Analysis. Conference proceedings.
13. Patel, A., Desai, R., & Shah, K. (2018). OCR-Based Handwriting Recognition Systems: Challenges and Future Directions. Journal article.
14. Li, C., Wu, L., & Zhang, Y. (2019). OCR Techniques for Document Classification and Indexing: A Review. Book chapter.

15. Yang, J., Zhu, Q., & Wang, X. (2017). OCR Systems for Mobile Devices: Performance Evaluation and Optimization. Journal article.

FINANCING

None.

CONFLICT OF INTEREST

None.

AUTHORSHIP CONTRIBUTION

Conceptualization: Muthusundari, A Velpoorani, S Venkata Kusuma, Trisha L, Om.k.Rohini.

Data curation: Muthusundari, A Velpoorani, S Venkata Kusuma, Trisha L, Om.k.Rohini.

Methodology: Muthusundari, A Velpoorani, S Venkata Kusuma, Trisha L, Om.k.Rohini.

Drafting - original draft: Muthusundari, A Velpoorani, S Venkata Kusuma, Trisha L, Om.k.Rohini.

Writing - proofreading and editing: Muthusundari, A Velpoorani, S Venkata Kusuma, Trisha L, Om.k.Rohini.