# Text extraction using OCR: A Systematic Review

Rishabh Mittal

Department of Computer Science & Engineering
Amity School for Engineering and Technology
Amity University Uttar Pradesh, Noida (UP), India
rm.23mittal@gmail.com

Anchal Garg

Department of Computer Science & Engineering
Amity School for Engineering and Technology
Amity University Uttar Pradesh, Noida (UP), India
agarg@amity.edu

*Abstract*— In the digital era, almost everything is automated, and information is stored and communicated in digital forms. However, there are several situations where the data is not digitized, and it might become essential to extract text from those to store in digitized form. The latest technology such as Text recognition software has completely revolutionized the process of text extraction using Optical Character Recognition. Therefore, this paper introduces the concept, explains the process of extraction, presents the latest techniques, technologies, and current research in the area. Such a review will help other researchers in the field to get an overview of the technology.

*Keywords—Optical Character Recognition (OCR), Digital Image Processing(DIP), Text recognition, Pre-processing, Feature extraction*

## I. INTRODUCTION

Optical Character Recognition (OCR) is the computerized conversion of text or making a digital copy of the text through sources like handwritten documents, printed text, or from natural images [1]. It comes under the wide umbrella of Digital Image Processing (DIP) [2]. DIP is the process in which digital images are processed through a computer algorithm. DIP is a field that has applications in pretty much every other field like in Healthcare division for PET sweeps, Banking segment, Robotics, and so forth and is as yet developing with time. One of its major applications is pattern recognition which includes computer-aided diagnoses, handwriting recognition, and image recognition.

The need for text recognition software came because the amount of data in the world is growing at an exponential rate. All this data cannot be stored physically and hence need to preserve it digitally. Thus, it is done using Automatic Character Recognition, which utilizes OCR. OCR frameworks are these days normally used to extract text content from any computerized picture or natural image.

OCR encourages checked archives to turn out to be something other than picture records, transforming them into completely accessible reports. Therefore, it is generally utilized for different purposes like information extraction. With OCR, people do not need to retype huge reports when making a digital copy of them. OCR isolates material information and enters it normally. It empowers a machine to perceive the content from a picture as seen it when perusing a composed record and store it in a manner that is simpler to process upon later[3].

Based on the kind of information, OCR frameworks may be characterized into 3:

1. Handwritten- Those systems that work only on written text.

2. Machine printed- Those that work only on the text that is typed and then taken a hard copy of.

3. Specific type- There are a lot of factors in the working of the OCR like the language, font, etc. Thus, there are a lot of OCR systems for specific languages like Urdu for example.

OCR these days are used majorly for text recognition, but it was created to help blind or specially-abled people back in 1914. It allowed the blind people to have written text read to them out loud by a machine. Also, since back in the 1950s, technology was not so advance and there was not sufficient computing power, the advancement of OCR confronted numerous difficulties like speed and accuracy [1]. Early optical character versions required training on individual character images and worked on one word at a time. In its initial stage, the best OCR framework could just perceive 1 word per minute.

This paper summarizes the research in OCR. The paper is structured as follows: Section II covers modern OCR Systems,

## II. MODERN OCR SYSTEM

There are a lot of OCR engines that are used these days like the Google Drive OCR, Tesseract, Transym, OmniPage, etc. Many of them are paid, however, some are accessible for nothing.

[4] Tesseract is one of the most popular and commonly used engines in OCR frameworks that is available on the internet. It works on a step by step process involving 4 steps. It focuses on covering a wide range of languages and fonts instead of accuracy. When compared to Transym OCR (another open-source OCR engine), it is found that tesseract offers higher accuracy on average than Transym but it is not necessarily faster every time

Program/software uses these engines to convert digital images into a form that can be used to edit them. One of the most popular OCR programs is the Adobe Acrobat Pro which can perform a lot of OCR related functions. Abby Fine Reader is also one such program offering similar functions.

With the help of modern Wireless facilities, a device does not have to have the entire OCR system built into it, instead, it can simply send the data over the network where it will be worked upon and results are sent back to the device extremely quick. Online OCR is one such Web-based OCR converter that can function as a complete system over the internet.

## III.    PROCESS OF TEXT RECOGNITION USING OCR

As mentioned in [5] Creating a fully functional OCR framework requires many steps but they can be majorly grouped into 6 steps as listed in figure 1.
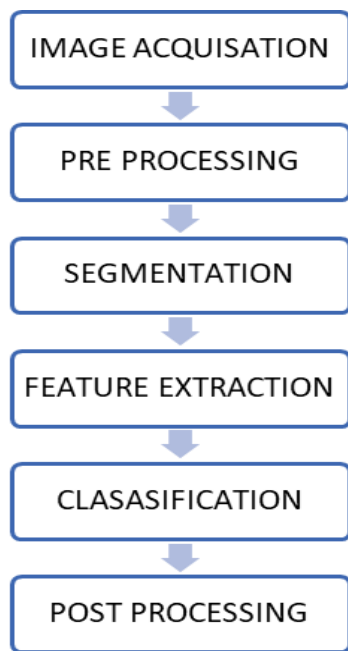
IMAGE ACQUISATION

↓

PRE PROCESSING

↓

SEGMENTATION

↓

FEATURE EXTRACTION

↓

CLASSIFICATION

↓

POST PROCESSING

Figure 1: Steps towards Recognition in OCR

### A.    The image acquisition:

Image acquisition is the underlaying advance of OCR that includes taking an advanced picture and changing it in an appropriate structure that can be viably processed upon later. This includes the quantization of picture and image compression.

### B.    Preprocessing

Pre-processing is critical to accomplish higher recognition rates as utilizing it makes the OCR increasingly strong. It starts with binarization [6] (converting multi-tone image into black and white) and then, includes the following image enhancement procedures [7] that improve the detail of the picture.

1. Spatial image filtering operations- these are commonly used to either smoothen the picture or make edges increasingly obvious. It is further divided into two parts -point processing and mask processing. Point processing is done on individual pixels while mask processing is done on a group of pixels.

2. Thresholding- it isolates the data from an unnecessary part of the image. It is normally applied to greyscale images. It can be divided into two parts: global and local. Global strategies lean toward a solitary worth cut off for the whole archive while the local methods essentially apply to explicit applications and more often than not, they don't function admirably on isolated applications.

3. Noise removal- Various kinds of noises are identified with devices that are used to capture images like a photon, lighting in the surroundings, electronics in the device, and so forth. Although thanks to modern technology it is possible to reduce the noise while photographing to almost insignificant levels.

4. Screw detection/correction- Natural images can often be rotated making it tough for an OCR framework to work upon it. So, these images ought to be rectified before processing further. There are various methods for it like Hough transform.

### C.    Segmentation

Segmentation, as the name suggests, is used to isolate the required substance from the rest of the image.[8] It involves various steps:

1. Page segmentation- it is used to isolate content from the rest of the image, resulting in a text-only image. It can be portrayed into three general characterizations: top-down, base up, and a half and half strategies.

2. Character segmentation- it is seen as one of the chief strides in preprocessing especially in cursive contents where symbols are related together. Henceforth, there are various systems created for character division, and most by far of them are content unequivocal and may not work with various substances.

3. Image size normalization- The outcome from the character division stage gives isolated characters that make way to next process; in this way, the detached parts get standardized in a particular size tentatively relying upon the utilization later and the feature extraction or grouping techniques utilized, at that point highlights are removed from all characters with a steady size to give information consistently.

4. Morphological processing- Sometimes, few pixels may be evacuated creating openings to certain pieces of the pictures. It is like having a few gaps where a portion of the pixels was expelled during thresholding. Then again, the inverse can likewise be valid, it may overwrite separate items which result in difficulty to isolate characters; these strong articles take after masses and are difficult to decipher. The answer to these issues is Morphological Filtering. Helpful methods incorporate disintegration and expansion, opening and shutting, delineating, and diminishing and Skeletonization

## D. *Feature extraction*

In this step, each character is allocated a vector which at that point represents it. Its goal is to extract a group of features, that will increase the recognition for a small number of items and produce the same feature used in different instances of the same attribute.[9] This is done using different techniques:

1. Zoning- A character is typically separated into zones of predefined size. These predefined or matrix sizes are ordinarily of the request 2x2, 4x4, and so forth. At that point, the densities of pixels or features are broken down in various zones to form the representations.

2. Projection histogram features- They measure the number of pixels in a very specific manner. There are three kinds of projection histograms – flat, vertical, and diagonal.

3. Distance profile features- They measure the distance as the number of pixels from the merge box of the image to the edge of the symbol. They can be taken in any direction such as up, down, left, or right.

4. Background directional distribution (BDD)- It is used to calculate the amount of background distribution of each front pixel. Several masks are used in different directions where the average grey matter value is calculated using a specific mask.

5. The combination of various features- more than one method can be used together for even higher accuracy, as can be seen in figure 2.

| Feature Vector | Included features |
|---|---|
| FV1 | Zoning + profiles |
| FV2 | BDD + zoning |
| FV3 | BDD + histograms |
| FV4 | Profiles + horizontal and vertical histograms(HVH) |
| FV5 | BDD + HVH |
| FV6 | BDD + Profiles |
| FV7 | BDD + Diagonal (both) histograms |

Figure 2: Combination of feature extraction methods

## E. *Classification*

The feature vector that is acquired in the previous step is utilized for classification. To do the classification the information and many element vectors should be contrasted. [10] A classifier is utilized to analyze the element vector of info and the component vector of the information bank. The choice of classifier relies on the application, preparing the set, and the number of free parameters. There are numerous strategies utilized for classification, 3 of them, which are most normally utilized are-

1. Probabilistic Neural Network (PNN) classifier- It is a classifier that employs a multi-layered feed-forward neural network to classify unknown patterns using probability density function.

2. Support Vector Machines (SVM) classifier- These are supervised learning techniques that may be enforced for classification or regression. It takes a collection of input and predicts to classify them within the best 2 categories. SVM classifier is trained using a set of coaching data and to categorize test statistics a model is prepared based totally upon this.

3. K- Nearest Neighbor (K-NN) classifier- It uses a model-based lesson in accordance with an unknown pattern known for a particular distance or other similar activity. It divides the thing by the votes of its neighbors. Because it looks only at a neighboring object up to a certain level, it uses the spatial correlation of the distance function.

## F. *Post-Processing*

It is the last advance after the arrangement. As the outcomes are not 100% right, particularly for complex dialects, Post handling procedures can be performed to improve the precision of OCR frameworks. These procedures use characteristic language preparing, geometric, and etymological setting to address blunders in OCR results. Postprocessor ought not to take a lot of time and cause new blunders.

## IV. NEW TECHNIQUES

Other than the regular techniques referenced before, various strategies proposed by various researchers in the field have been investigated and are summarized below:

- A new methodology [11] to extract the text from natural images having 7 stages was proposed. Initially, the filtering process is utilized for pre-processing to improve the image. Lateral separation is done using the Thresholding method to separate the background from the required content. Then, the MSER (Maximally Stable External Region) is detected and the part which is not required is deleted. Then the stroke width calculation is applied by the stroke width variance algorithm and lastly, CNN (Convolutional neural network) is used to get the features required to spot the characters and that will be provided to the OCR to obtain the text.

- Another [12] proposed a neural network-based framework that operates based on BLSTM-Bidirectional Long Short-Term Memory that allows OCR to work at the word level. It leads to over 20% better results when compared to a regular OCR framework. It uses a method that does not require segmentation, that is one amongst the foremost common reasons for the error. Also, it found an over 9% decrease in character error compared to the more widely available OCR framework.

- [13] This technique proposed utilizing a two-advance iterative Conditional Random Field (CRF) calculation with Belief Propagation obstruction to isolate content

and non-content parts and afterwards utilizing OCR on the content part to give the ideal outcome. In the case of multiple text lines, two relational graphs are used to extract different text lines and the OCR confidence is utilized as a guide while finding text containing parts.

There are a lot more techniques other than these in research. These techniques attempt to improve the overall system of OCR by increasing accuracy or working on areas that are the most common error source.

## V. SOME MAJOR RESEARCH

There is a lot of application-based researches going on using OCR. Some of them are mentioned below:

- *Text Extraction from Historical Documents*: OCR's Complete Method for Historical Texts without font knowledge. It consists of 3 steps - initially, a step that incorporates binarization and enhancement. In the second step, a high-level separation method is used to distinguish line parts, words, and letters. The KNN integration theme is adopted to incorporate the symbols of the same group. Lastly, within the third step, in each image of the new document, the same previous classification method is utilized whereas the recognition is predicated on the information extracted from the previous step. It results in high throughput rates of up to 95.44% [14].

- *Text Extraction from Television:* This program uses several steps. First, the hosting service is used to get the data processed. Subsequently, the OCR algorithm is employed to separate the text from the given data. Finally, output presented in the previous step is compared to expectations, and a decision is made whether it is correct or not before issuing it. This program can be used as a practical test for TV sets [15].

- *Vehicle Identification Using Number Plate Recognition:* –OCR is used in many areas, security is one of them. This system can be used to keep track of traffic at the security entrance. First, a photograph of the automotive is taken and then, the number plate is separated using parts of the pictures. Afterwards, the OCR is employed to get the text from it and lastly, the details are compared to the inventory dataset to find the automobile owner's details. It makes a strong system with high reliability for security purposes [16] [17].

## VI. CURRENT WORK AND CHALLENGES IN IT

Several research papers suggest potential research advancements in this area. For instance, as per [18] currently improving components like Scan goals, filtered picture quality,

type of printer utilized whether ink-jet or laser, the nature of the paper, phonetic complexities, the lopsided brightening, and watermarks can impact the precision of OCR. Hence work can be done on improving the precision of OCR.

There are different issues looked by an OCR framework [19], particularly Chinese and Arabic character acknowledgment uniquely as these dialects contain muddled structures, lopsided text styles, and so on. There are various issues in the event of handwritten record pictures like the nearness of slanted, contacting or covering content lines, etc. [20] [21] therefore, new algorithms like Spiral Run Length Smearing Algorithm (SRLSA) have been under research. Additionally, since it isn't yet 100% precise, the extracted text despite everything must be crosschecked for errors, and for limited text, it is simply not worth utilizing it as it tends to be quicker to do physically compared to modern OCR systems.

Since every language and font is different from one another, there is no single method that can be applied to all to get the desired result with high accuracy [22]. In this manner, there is a great deal of research proceeding to improve the precision of the OCR framework for a specific language or font, and even for the same language, more than one method can be applied. Some examples are mentioned below:

1. Gurmukhi script- [23] There are different issues in Gurmukhi content OCR like covering characters, variable composing styles, comparability of certain characters, the unavoidable nearness of foundation commotion, and different sorts of mutilations. Likewise, when various strategies are utilized for manually written Gurmukhi content contrasting in the feature extraction strategy and classifier utilized. The distinct outcome is obtained with different strategies like with zone Density and BDD as feature extraction technique, and SVM with the RBF bit as the classifier, the highest accuracy was achieved which is 95.04%.

2. Devanagari script- [24] Zoning feature extraction methods are used for Devanagari script OCR. In the zoning technique, several methods can be used. When the 4x4 grid was used for zoning, there was a significant improvement, but when 2x2, 8x8, or 16x16 grids were used for zoning, the performance was equal to the original feature value. Thus, there is a significant improvement in yield when a 4x4 grid size is used for Devanagari.

3. Arabic script- [25] It's written from right to left and many characters may have different shapes depending on its contexts. Also, semi-cursive nature and discontinuities create more problems for the OCR framework. Generalized Hough Transformation can be used to improve accuracy in character segmentation and different fonts, the result obtained

was 86% accurate. Whereas, for cursive, a success rate of up to 97% could be obtained.

4. Bengali script- [26] Bengali is the fifth most communicated in language on the planet. It has no upper or lower case like the ones in English but Bengali characters have a part called "Matra" that remains connected making the process of recognition slower. Self-Organizing Map (SOM) otherwise referred to as Kohonen Neural Network (KNN) is used to group different characters. This results in saving of 33% recognition time as compared to standard methods.

5. Rashi font- [27] An algorithm supported by the utilization of fuzzy logic-based rules, depending upon the factual information of the investigated text style is utilized. This new methodology consolidates letter measurements and relationship coefficients in a lot of fuzzy-based principles, empowering the extraction of contorted letters that might not be recognized by some other method. It centers around Rashi text styles related to critiques of the book of the bible that is transcribed calligraphy.

6. Multilingual Indian document- [28] A thorough examination of the various databases including written and computer-generated records is done. The best result was found to be around 99% accurate. Also, algorithms that bolster multilingual Indian report pictures containing mixed writings of Devanagari and Latin contents are proposed. The primary two components are generated using the pixel of the character center, while the third element is calculated using the neighbor's pixel information.

## VII. OTHER APPLICATIONS OF OCR

The OCR has become an important part of modern-day technologies such as text to speech software, barcode scanners, etc. It is also used to assist the Augmented Reality or AR as it centers around Object Recognition. OCRs can be used not only to recognize text but other objects as well.

OCR was first created to help blind people but since then the technology has evolved a lot so it has become easier to create a text to speech engine as mentioned in [29].

With the help of other fields like Artificial Intelligence, it might also be used to help reduce manual work in many areas like in the banking sector or airport where a person has to fill out a form from which data has to be entered in the online database as OCR may extract the text, but it needs support in filling that data into an online database. For example, some banks in China use OCR combined with AI to provide 2-step authentication at ATMs.

Also, to cross-check or verify any information like matching a person's name from a very big database can be done using OCR much faster. Similarly, creating a digital copy of any document is extremely easy. It enables digital images to be used as a fully searchable document like Google Books.

## VIII. CONCLUSION

Optical Character Recognition has a variety of applications. It is being used to extract the text from ancient books and scripts, images, etc. There is ongoing research in this area to improve the precision of the extracted text and increasing the range on which a single OCR system can work. Since it came up there have been challenges that researchers have faced in this field, some of them that have been solved like the slow speed of early OCR systems but some are yet to solved like achieving a 100% accuracy and limited range for a single OCR. Thus, albeit there's an enormous amount of analysis occurring in this field, there's still a lot of scope for a lot more. In the future, new techniques or algorithms can be found to improve the accuracy and help create an OCR system that can work on any kind of dataset regardless of language, font, etc. Also, with the help of other fields like Artificial Intelligence and Augmented Reality, it can be used for a lot more applications and its possibilities become endless.

## REFERENCES

[1] N. Islam, Z. Islam and N. Noor," A Survey on Optical Character Recognition System", Journal of Information & Communication Technology-JICT Vol. 10 Issue. 2, December 2016.

[2] S. Muthuselvi and P. Prabhu," Digital Image Processing Techniques-A Survey", International Journal of Open Information Technologies vol.5 Issue.11, May-2016.

[3] P. M. Manwatkar and K. R. Singh, "A technical review on text recognition from images," 2015 IEEE 9th International Conference on Intelligent Systems and Control (ISCO), Coimbatore, 2015, pp. 1-5, doi: 10.1109/ISCO.2015.7282362.

[4] C. Patel, A. Patel, and D. Patel," Optical Character Recognition by Open Source OCR Tool Tesseract: A Case Study", International Journal of Computer Applications (0975 – 8887) Volume 55– No.10, October 2012.

[5] K. Hamad and M. Kaya," A Detailed Analysis of Optical Character Recognition Technology ", International Journal of Applied Mathematics, Electronics and Computers, 3rd September 2016.

[6] Poovizhi P," A Study on Preprocessing Techniques for the Character Recognition", International Journal of Open Information Technologies ISSN: 2307-8162 vol. 2, no. 12, 2014.

[7] R. Maini and H. Aggarwal," A Comprehensive Review of Image Enhancement Techniques" journal of computing, volume 2, issue 3, March 2010, issn 2151-9617.

[8] A. Shinde and D.G.Chougule," Text Pre-processing and Text Segmentation", IJCSET |January 2012| Vol 2, Issue 1,810-812.

[9] V. Prasad and Y. Singh," A study on structural method of feature extraction for Handwritten Character Recognition", Indian Journal of Science and Technology · March 2013.

[10] A. Nadarajan and Thamizharasi A, "A Survey on Text Detection in Natural Images", International Journal of Engineering Development and Research (IJEDR), ISSN: 2321-9939, Volume.6, Issue 1, pp.60-66, January 2018.

[11] C.P. Chaithanya, N. Manohar and Ajay Bazil Issac," Automatic Text Detection and Classification in Natural Images", International Journal of Recent Technology and Engineering (IJRTE) ISSN: 2277-3878, Volume-7, Issue-5S3, February 2019.

[12] N. Sankaran and C. V. Jawahar, "Recognition of printed Devanagari text using BLSTM Neural Network," Proceedings of the 21st International

Conference on Pattern Recognition (ICPR2012), Tsukuba, 2012, pp. 322-325.

[13] H. Zhang, C. Liu, C. Yang, X. Ding and K. Wang, "An Improved Scene Text Extraction Method Using Conditional Random Field and Optical Character Recognition," 2011 International Conference on Document Analysis and Recognition, Beijing, 2011, pp. 708-712, doi: 10.1109/ICDAR.2011.148.

[14] G. Vamvakas, B. Gatos, N. Stamatopoulos, and S. J. Perantonis, "A Complete Optical Character Recognition Methodology for Historical Documents," 2008 The Eighth IAPR International Workshop on Document Analysis Systems, Nara, 2008, pp. 525-532, doi: 10.1109/DAS.2008.73.

[15] I. Kastelan, S. Kukolj, V. Pekovic, V. Marinkovic, and Z. Marceta, "Extraction of text on TV screen using optical character recognition," 2012 IEEE 10th Jubilee International Symposium on Intelligent Systems and Informatics, Subotica, 2012, pp. 153-156, doi: 10.1109/SISY.2012.6339505.

[16] M. T. Qadri and M. Asif, "Automatic Number Plate Recognition System for Vehicle Identification Using Optical Character Recognition," 2009 International Conference on Education Technology and Computer, Singapore, 2009, pp. 335-338, doi: 10.1109/ICETC.2009.54.

[17] Vaishnav, A., & Mandot, M. (2020). Template Matching for Automatic Number Plate Recognition System with Optical Character Recognition. *Information and Communication Technology for Sustainable Development* (pp. 683-694). Springer, Singapore.

[18] Ayatullah Faruk Mollah, Nabamita Majumder, Subhadip Basu, and Mita Nasipuri," Design of an Optical Character Recognition system for Camera-based Handheld Devices ", IJCSI International Journal of Computer Science Issues, Vol. 8, Issue 4, No 1, July 2011.

[19] .Karthick, K.B.Ravindrakumar, R.Francis, and S.Ilankannan," Steps Involved in Text Recognition and Recent Research in OCR; A Study", International Journal of Recent Technology and Engineering (IJRTE) ISSN: 2277-3878, Volume-8, Issue-1, May 2019.

[20] S. Malakar, S. Halder, R. Sarkar, N. Das, S. Basu, and M. Nasipuri, "Text line extraction from handwritten document pages using spiral run-length smearing algorithm," 2012 International Conference on Communications, Devices and Intelligent Systems (CODIS), Kolkata, 2012, pp. 616-619, doi: 10.1109/CODIS.2012.6422278.

[21] M. N, S. R, Y. G and V. J, "Recognition of Character from Handwritten," *2020 6th International Conference on Advanced Computing and Communication Systems (ICACCS)*, Coimbatore, India, 2020, pp. 1417-1419, doi: 10.1109/ICACCS48705.2020.9074424.

[22] Gaurav Y. Tawde and Mrs. Jayashree M. Kundargi," An Overview of Feature Extraction Techniques in OCR for Indian Scripts Focused on Offline Handwriting", International Journal of Engineering Research and Applications (IJERA) ISSN: 2248-9622 www.ijera.com Vol. 3, Issue 1, January -February 2013, pp.919-926

[23] Pritpal Singh and Sumit Budhiraja," Feature Extraction and Classification Techniques in O.C.R. systems for Handwritten Gurmukhi Script – A Survey", / International Journal of Engineering Research and Applications (IJERA) ISSN: 2248-9622 www.ijera.com Vol. 1, Issue 4, pp. 1736-1739

[24] O. V. Ramana Murthy and M. Hanmandlu," Zoning based Devanagari Character Recognition", International Journal of Computer Applications (0975 – 8887) Volume 27– No.4, August 2011.

[25] Sofien Touj, Najoua Ben Amara, and Hamid Amiri," Generalized Hough Transform for Arabic Printed Optical Character Recognition", The International Arab Journal of Information Technology, Vol. 2, No. 4, October 2005.

[26] M. G. Kibria and Al-Imtiaz, "Bengali Optical Character Recognition using self-organizing map," 2012 International Conference on Informatics, Electronics & Vision (ICIEV), Dhaka, 2012, pp. 764-769, doi: 10.1109/ICIEV.2012.6317479.

[27] E. Gur and Z. Zelavsky, "Retrieval of Rashi Semi-cursive Handwriting via Fuzzy Logic," 2012 International Conference on Frontiers in Handwriting Recognition, Bari, 2012, pp. 354-359, doi: 10.1109/ICFHR.2012.262.

[28] Sahare P. and Dhok S. B., "Multilingual Character Segmentation and Recognition Schemes for Indian Document Images," IEEE Access, vol. 6, pp. 10603-10617, January 18, 2018.

[29] Dr. S. Manoharan," A Smart Image Processing Algorithm For Text Recognition, Information Extraction And Vocalization For The Visually Challenged", Journal of Innovative Image Processing (JIIP) (2019) Vol.01/ No. 01 Pages: 31-38.