



OCR-MRD: performance analysis of different optical character recognition engines for medical report digitization

Pulkit Batra¹ · Nimish Phalnikar¹ · Deepesh Kurmi¹ · Jitendra Tembhurne¹ · Parul Sahare¹ · Tausif Diwan¹

Received: 15 March 2023 / Accepted: 25 October 2023 / Published online: 24 November 2023

© The Author(s), under exclusive licence to Bharati Vidyapeeth's Institute of Computer Applications and Management 2023

Abstract In the modern era, the necessity of digitization is increasing in a rapid manner day-to-day. The healthcare industries are working towards operating in a paperless environment. Digitizing the medical lab records help the patients in hassle-free management of their medical data. It may also prove beneficial for insurance companies for designing various medical insurance policies which can be patient-centric rather than being generalized. Optical Character Recognition (OCR) technology is demonstrated its usefulness for such cases and thus, to know the best possible solution for digitizing the medical lab records, there is a need to perform an extensive comparative study on the different OCR techniques available for this purpose. It is observed that the current research is focused mainly on the pre-processing image techniques for OCR development, however, their effects on OCR performance specially for medical report digitization yet not been studied. Herein this work, three OCR Engines viz Tesseract, EasyOCR and DocTR, and six pre-processing techniques: image binarization, brightness transformations,

gamma correction, sigmoid stretching, bilateral filtering and image sharpening are surveyed in detail. In addition, an extensive comparative study of the performance of the OCR Engines while applying the different combinations of the image pre-processing techniques, and their effect on the OCR accuracy is presented.

Keywords Classifier · Feature extraction · Optical character recognition · Image pre-processing · Computer vision · Medical reports · Digitization

1 Introduction

The age of digitization has made it mandatory to have digital records to make the data easily available. Philip et al. [1] observed that hospitals look after digitization of medical records in order to reduce costs and also improve accessibility to records. This medical record is extremely important from patient as well as from doctor's perspective, as it serves a basis for early diagnosis and tracking the history of the ailment. The hospitals which are equipped with the usage of digitized records showing better performance in terms of less time for searching the records and more time for patient care as compare to hospitals with paper records [1]. In such scenario, it is observed that there is a need for hospitals to move towards process of paperless environment, wherein a tool is required to convert the existing paper records into their digital form in a simple and efficient manner.

Usually, printed documents are first scanned and then saved in the form of image in order to store the information present in them. However, utilizing this available information by only reading the text within images is quite a difficult task. Optical Character Recognition (OCR) is one of the active research areas that involves number of

✉ Jitendra Tembhurne
jtembhurne@iiitn.ac.in

Pulkit Batra
bt18cse052@iiitn.ac.in

Nimish Phalnikar
bt18cse099@iiitn.ac.in

Deepesh Kurmi
bt18cse086@iiitn.ac.in

Parul Sahare
parulsahare@iiitn.ac.in

Tausif Diwan
tdiwan@iiitn.ac.in

¹ Indian Institute of Information Technology, Nagpur, Maharashtra, India

attempts to develop an automated system that extracts and processes text from the image files. The objective of OCR is to provide an editable digital format by modifying or converting any form of text from document images such as handwritten and printed scanned texts for deep and further processing [2]. The OCR mainly focuses on the recognition of two types of characters i.e. machine printed and handwritten. The varied characteristics makes them a challenging area of research [3]. Number of challenges such as variations in font size, font colour, background colour, image resolution, etc. which the OCR faces while performing the digitization, and thus, makes it vulnerable to inaccuracies [2]. This creates the pre-processing of images a crucial step wherein the characteristics of images are enhanced. Therefore, there is a need to carry out by taking in consideration the characteristics of the OCR, as every OCR Engine responds differently to the characteristics of the image e.g. resolution, colour correction, contrast etc. Subsequently, Braille system is developed for written communication to help visually impaired persons wherein language is converted into Braille characters. This optical Braille recognizer reads pattern from images and translate into text [21]. In [22], character recognition on Offline handwritten is investigated using machine-encoded text and support vector machine (SVM) is applied for classification in forensic applications. Here, local features from the image are extracted to improve the accuracy, moreover, 74.32% of accuracy is reported for character recognition.

Even though number of pre-processing techniques developed in the past, these techniques are dependent upon the category of images and which features should be extracted from these images. Most of the available comparative studies are mainly focused on a specific characteristic of the image, i.e., from where data needs to be extracted, however, a study which compares various set of tools and techniques based on the utilization is needed. Research in the digitization of medical lab reports is still one of the unexplored areas. Study relating to the aftermath or advantages of medical data digitization is found quite abundantly but investigation regarding optimal approaches to digitize the data relating to specific medical domains is quite insufficient. In the propose work, we focus on digitizing medical lab reports which are a major resource in understanding patients' health charts and his/her medical condition changes over a period of time in response to the treatment they are receiving.

The paper is organized as follows; Sect. 2 discuss on OCR engines and various image preprocessing techniques. Literature review based on the medical report generation is presented in Sect. 3. Section 4 highlights the proposed methodology for medical report digitization. Results and discuss is found in Sect. 5, finally the conclusion and future scope is presented in Sect. 6.

2 OCR engines and image pre-processing techniques

In this paper, we have mainly focused upon three open-source OCR Engines detailed in [10, 13, 17] and six pre-processing techniques, which is mentioned in [14] to draw necessary conclusions for this work. Our attention while choosing the OCR engines is to find out open-source, easily accessible and robust engines so that the target user can develop their applications without suffering from any licensing issues.

2.1 OCR engines

Here, three open-source OCR Engines are described and there technical details are presented as follows:

Tesseract OCR: An open-source text recognition (OCR) Engine, available freely under the Apache 2.0 license. It can be run both from the command line and through GUI Interface (using third Party compatibility) [18]. It consists of a fully-featured API and can be compiled for a variety of targets including Android and the iPhone. It is also available for Linux, Windows and Mac OS X platforms. The development of tesseract is focused on line finding, features/classification methods, and the adaptive classifier for achieving the best accuracy [9].

EasyOCR: It is ready-to-use, open-source OCR with support for 80+ languages and all popular writing scripts including Latin, Chinese, Arabic, Devanagari, Cyrillic, etc. EasyOCR is adaptable to any state-of-the-art model plug-in [5]. In Fig. 1 the flow of the development of EasyOCR is shown [19].

docTR: An open source OCR engine available under Apache 2.0 license. docTR is powered by TensorFlow 2 and PyTorch. The text detection and text recognition is performed using the following techniques [6].

Text detection

- Double-Butterfly Network (DBNet): Real-time scene text detection with differentiable binarization.
- LinkNet: Exploiting encoder representations for efficient semantic segmentation.

Text recognition

- Convolutional recurrent neural network (CRNN): The combination of two of the most prominent neural networks, CNN (Convolutional Neural Network) followed by the RNN (Recurrent Neural Networks). Optimal neural networks model for image-based sequence recognition and its application to scene text recognition.

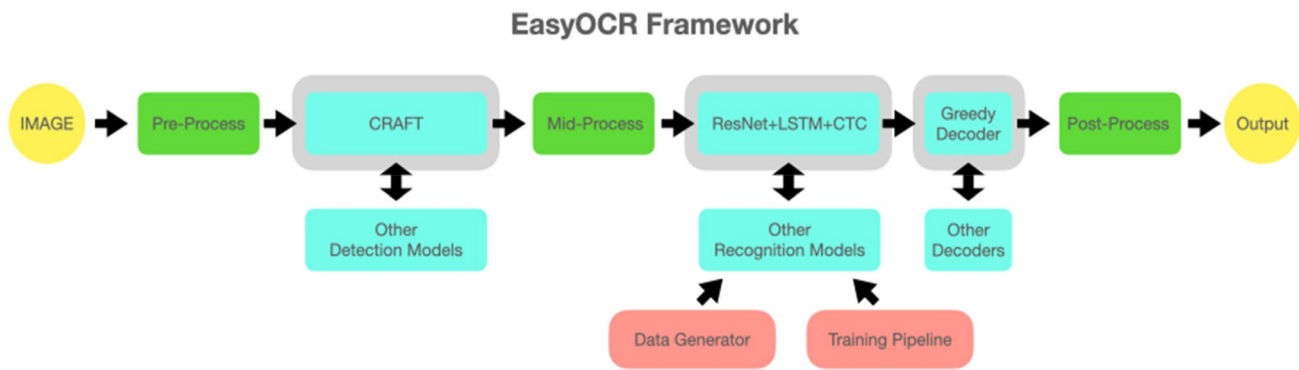


Fig. 1 Easy OCR framework

- Show, Attend and Read (SAR): A simple and strong baseline for irregular text recognition.
- Multi-Aspect Non-local Network for Scene Text Recognition (MASTER): A scene text recognizer targeting to identify characters in a sequence from image with straight text.

2.2 Pre-processing

Image binarization: Image binarization is the process of converting a document image from 3D pixels array format into 2D format, i.e. into a bi-level document image. Image pixels are separated into a dual collection of pixels, i.e. black and white. The main goal of image binarization is the segmentation of documents into foreground text and background.

Brightness transformation: Here, the brightness information is given by the Eq. (1):

$$g(x) = \alpha f(x) + \beta, \alpha > 0 \tag{1}$$

The parameters α = gain and β = bias parameters are said to control contrast and brightness, respectively. The image brightness and contrast vary as we change α and β .

Gamma correction or power law transform: Gamma correction carries out a non-linear operation on the source image pixels and can cause saturation of the image to be altered. Gamma correction simply can be defined by the following power-law expression:

$$V_{out} = AV_{in}^\gamma \tag{2}$$

where, the non-negative real input value V_{in} is raised to the power γ and multiplied by the constant A to get the output value V_{out} . In the common case of $A = 1$, inputs and outputs are typically ranging from 0–1.

Sigmoid stretching: Sigmoid function is a continuous nonlinear activation function. Statisticians call this function as the logistic function.

$$g(x, y) = \frac{1}{1 + e^{c(th - f_s(x, y))}} \tag{3}$$

Here, $g(x, y)$: enhanced pixel value, c : contrast factor, th : threshold value and $f_s(x, y)$: original image.

The amount of lightening and darkening can be varied to control the overall contrast enhancement by adjusting the contrast factor ‘ c ’ and threshold value [20].

Image smoothing (bilateral filtering): Image blurring is achieved by convolving the image with a low-pass filter kernel, which helps in removal of noises. Here, a bilateral filter is a Gaussian filter in space and is highly effective in noise removal while keeping edges sharp.

Image sharpening (Kernel method): High-pass filter is used to emphasize the fine details in the image to enhance the sharpness of the image. It is opposite to that of low-pass filter and uses a different convolution kernel than a low-pass filter.

$$Kernel = \begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

3 Literature review

It is noted from the literature review that research in the digitization of medical lab reports is still one of the unexplored areas. Most of the available comparative studies are focused on a specific area (such as image-smoothing) in the process of digitizing printed or handwritten data, however, a collective study for a complete process model is needed. It is found out that the study relating to the advantages of medical data digitization is quite abundant, but studies regarding optimal approaches to digitize data relating to specific medical domains are quite insufficient.

Scott et al. [1] discussed how hospitals with the usage of digitized records show a decrease in time to search for records and an increase in time for direct patient care, compared to hospitals with paper records. Suter et al. [2] projected the importance of big data analytics in the digitization of medical data to enable value-based healthcare. To further knowing the nuances of the OCR techniques, [4] and [7] are helpful in understanding all phases of OCR and understanding the problems that are faced during text recognition [3]. Mello et al. [8] detailed the best value of parameters for digitization (resolution, brightness, contrast, number of colors, etc.) that offers an overview of the impact on changing parameters and methods consists of OCR and their accuracy.

In OCR, image binarization helps in representation of textual data to ease the process of character recognition. The binarization in old documents as an image is processed in [23] and enhanced performance is reported. Authors, claimed that the proposed system is better for images with degraded documents. In [24], segmentation and recognition of document image is investigated using Dijkstra's algorithm and wavelet transform is applied for word segmentation. The accuracy of 98.1% is achieved for word segmentation and 97.6% of accuracy obtained for text-line segmentation, respectively. Subsequently, Tesseract-OCR is adopted for lot number detection in counting the stocks [25]. This system is proposed in hospital during COVID-19 to detect the stock quantity along with lot number of items that persist in hospital.

The novel pre-processing method is presented in [11] for scanned document wherein it detects and corrects the skew. The proposed method, minimizes the area of interest, independent of content and no limitation of skew angle. Subsequently, attentive generative network is employed to solve the problem of image denoising by embedding visual attention [12]. Due to visual attention, the noise region is attracted and forms the balance between removal of noise and preservation of texture. In [15], the detailed study is proposed on improving the output quality for the OCR system. Moreover, it very vital to digitized the medical information [16] to offer suitable techniques to process the medical data under different disease conditions. Table 1 present the detailed findings from the literature review performed.

4 Method and materials

It is seen that every OCR technique responds differently to the pre-processing techniques, which are applied to the input image. Herein, we have devised an approach to perform a combinatorial and permutation analysis using the set of pre-processing techniques that outputs the most valid sets with the best accuracy. The approach is divided into the following different sub-parts.

Lab reports: The input contains the lab reports which are in Portable Network Graphics (PNG) image format. A collection of lab reports is considered for performing this comparative study.

Image pre-processing techniques: The standalone pre-processing techniques on which the combinatorial and permutation study will be performed. These techniques will then be fed to the optimizer for further calculations.

Optimizer: The optimizer receives the pre-processing techniques as an input and uses the three OCR Engines. It then calculates the accuracy of text extraction after applying each pre-processing technique separately on all OCR engines. The technique whose accuracy is greater than the threshold accuracy (average of all accuracies) is supplied to the 'Combination Generator' for performing further analysis. The technique which performs below the threshold accuracy is discarded in this process.

Combination generator: The Combination Generator receives the input from Optimizer. Here, $2^{n+1} - 1$ combinations from the input set are generated, where 'n' is the number of pre-processing techniques. For each combination, $k!$ permutations are generated where 'k' is the length of a combination. For each permutation, pre-processing techniques are applied to the image in the order specified from the permutation. The image is then fed to each OCR and accuracy is calculated for the text extraction. A dictionary is maintained with key as permutation and value as accuracy which is then utilized to produce the final results. Figure 2 shows an example of how the combination generator gets its input and after that how it generates various combinations.

Accuracy calculator: To calculate accuracy, comparison the ground truth values of text with the extracted values of text from the lab reports is performed. The scoring is on the basis of proximity of extracted word to the original word, which is calculated using the Levenshtein distance. The scoring is explained in Table 2.

$$Accuracy = \frac{\sum(score)}{Total\ No.\ of\ words\ in\ original\ image} \times 100 \quad (4)$$

Optimized results: Using the dictionaries obtained from the 'Combination Generator', results are produced. The results is computed based on various parameters, presented as follows:

1. Accuracy without pre-processing vs accuracy of best combination of pre-processing.
2. Most suitable pre-processing techniques for each OCR.
3. Comparison of execution time of the OCR Engines.
4. Top five combinations of pre-processing for the best performing OCR.

Table 1 Summary of literature review

Refs.	Detailed description of the proposed model	Outcome	Dataset and model	Result
[1]	How digitizing medical records has been helpful for hospitals in reducing costs and improving accessibility? Clinical usability of digital records	Helpful in deciding the importance of features extracted according to need	Dataset was collected from two NHS hospitals using a work sampling approach. (size = 406) Two-way ANOVA models and Mann–Whitney <i>U</i> test	The hospital with the usage of digitized records showed a decrease in time to search records and an increase in time for direct patient care as compared to hospitals with paper records
[2]	How Big-Data Analytics helps in predicting outcomes through digitized data	Helpful in deciding the importance of features to determine better outcomes, which in turn will help in deciding optimal methods which will be more susceptible to our need	NA	As the digitization of medical data gets fully embraced and grows, Big-Data Analytics will help us enable value-based healthcare
[3]	A review of a comparative study of different feature extraction techniques used in OCR performed on handwritten text in Indian Scripts	How different feature extraction methods react to font style, text stroke, font ligatures, etc. and to choose the optimal one	NA	Based on a comparative study, we find that the feature extraction method that is best suited for one particular recognition application may not give optimum performance for the other application
[4]	A detailed overview of general phases of an OCR system such as pre-processing, segmentation, normalization, feature extraction, classification and post-processing	Helpful in understanding all phases of OCR and techniques in every phase. And helpful to understand common problems that we face in text recognition	NA	After performing all the phases of OCR we can get an accurate output of text recognition
[5]	Explained the various stages in text recognition and also explained the handwritten OCR systems classification according to the text type	Giving an idea to deal with different types of handwritten text like cursive, separate discrete text etc	NA	Explained all the steps of OCR for handwritten text and also discussed a study on some recent research papers like Electricity Billing using OCR
[6]	How to remove background images from the file without losing the quality of the text	Helpful in understanding the techniques of removing background images and thus to get a more accurate outcome	Some files with background images and with different background colors	Methods to remove background images and thus improve the accuracy of OCR are discussed
[7]	Provides information about the different OCR building steps, OCR engines available and about the techniques used in each of them and also about various OCR toolsets considering a variety of parameters	Gives a comprehensive understanding of the various processes in OCR building and the OCRs available currently	The dataset included OCR engines like Tesseract 4.0, GoOCR, OCRopus, ABBYY etc	OCRs were classified considering situations like distortion, handwritten, printed and other features of input to find suitable options
[8]	Provides insight into the difficulties of dealing with OCRs regarding the choice of the best method and in the best value of parameters for digitization (resolution, brightness, contrast, number of colors etc.)	We get an overview of the impact the changing parameters and methods have on the OCR and their accuracy	The documents have on average 3714.75 characters and 624.25 words. The documents were digitized with an HP Scanjet 4C scanner with a maximum resolution of 720 dpi. Brightness was set to 123% and contrast to 133%	Testing OCR engines on image types, rotation, maximum resolution, input file formats, font selection and multiple layouts
[9]	Provides information about the Tesseract OCR engine and its evolution	We get an overall review of the most popular OCR engine Tesseract and help in understanding the details of the OCR engine	The classifier was trained on a mere 20 samples of 94 characters from eight fonts in a single size, but with four attributes (normal, bold, italic, bold italic), making a total of 60,160 training samples	Results of current vs old Tesseract were computed

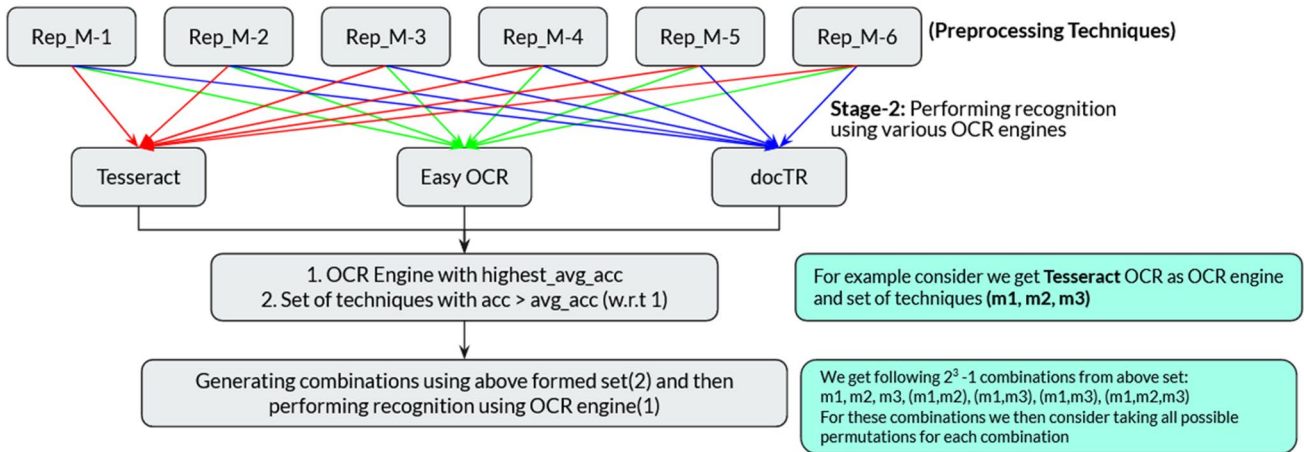


Fig. 2 Pre-processing techniques flow

Table 2 Levenshtein distance scores

Levenshtein condition	Score value
Exact match/Levenshtein distance 0	1.0
Inserting/deleting one character/Levenshtein distance 1	0.9
Substituting one character/inserting/deleting two characters/Levenshtein distance 2	0.8
Levenshtein distance > 2	0

Here, Fig. 3 shows how the aforementioned sub-parts form the flow of the proposed methodology.

5 Results and discussion

For the proposed methodology, we conducted a detailed comparative analysis of OCR engines and combinations of pre-processing techniques on a dataset of 39 images/lab-reports. The system configuration for performing the experimentation are—Intel i5 8th Generation CPU with 16 GB DDR4 RAM and 4 GB Nvidia 1050 GPU. Here,

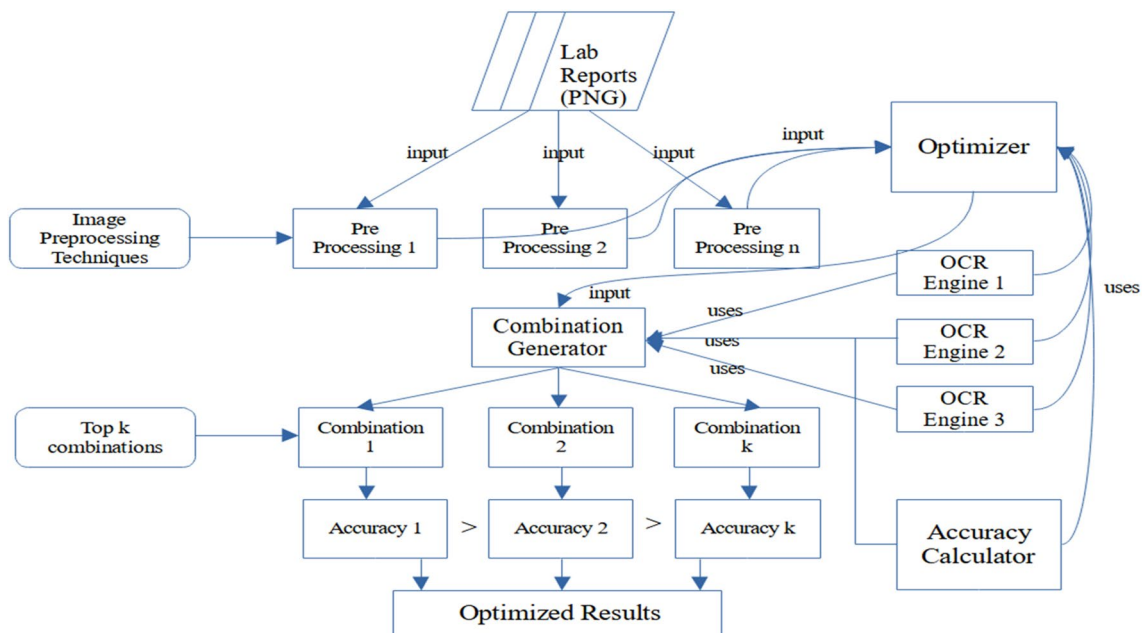


Fig. 3 Proposed system for identification of medical data

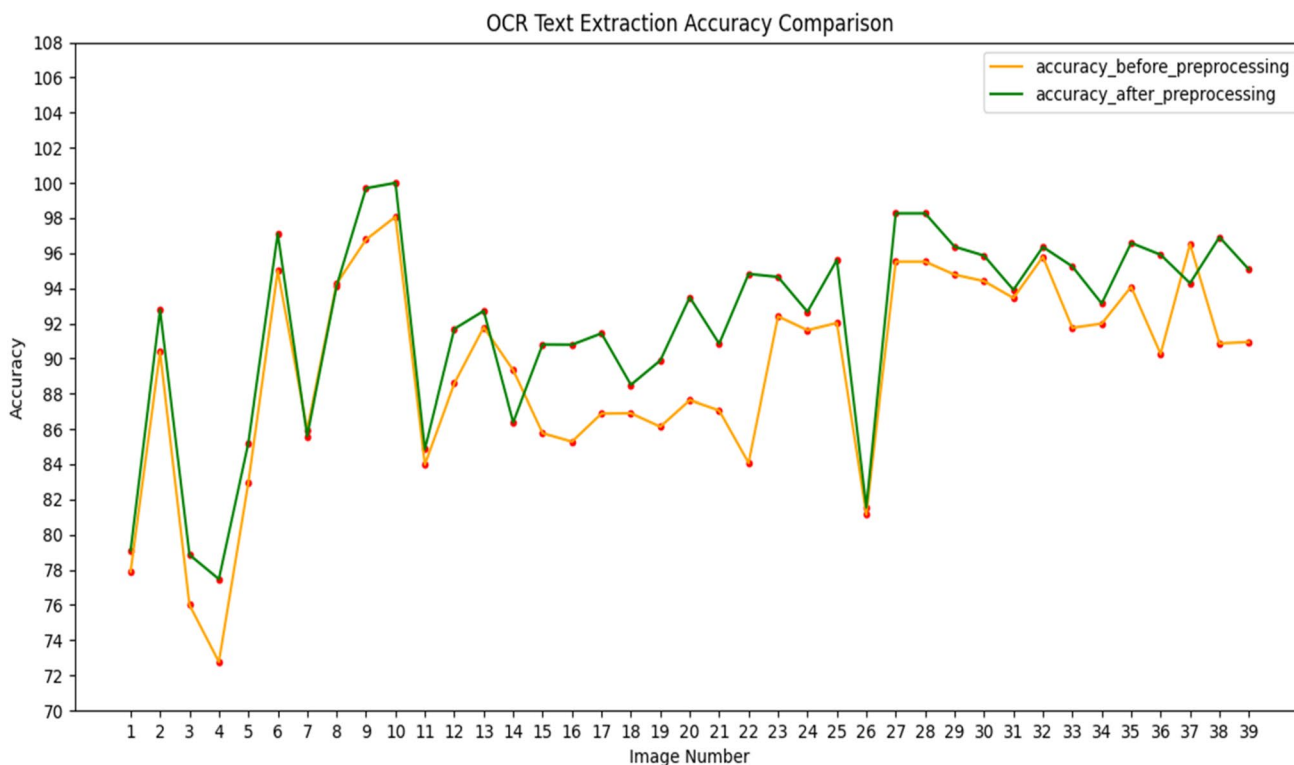


Fig. 4 Accuracy before pre-processing vs accuracy after pre-processing

Table 3 Top 2 pre-processing methods for each OCR

OCR engine	Top 2 pre-processing methods
Tesseract-OCR	Image smoothing Gamma correction
Easy-OCR	Sigmoid stretching Brightness transformation
docTR-OCR	Image smoothing Brightness transformation

Fig. 4 presents a comparison between the text extraction accuracy before any pre-processing techniques are applied and after the pre-processing techniques are adopted with the best-identified combination. Based on the OCR text extraction accuracy comparison graph, it is observed that the average accuracy of text extraction is improved by 2.88%. The average text extraction accuracy of 91.9% is reported, after applying the pre-processing techniques. Table 3 shows the most favorable pre-processing techniques for the OCRs. From this work, it is concluded that ‘Tesseract-OCR’ is 87% faster when compare to ‘Easy-OCR’ and 89% in comparison with ‘docTR-OCR’. After ‘Tesseract-OCR’, ‘Easy-OCR’ engine stands in second

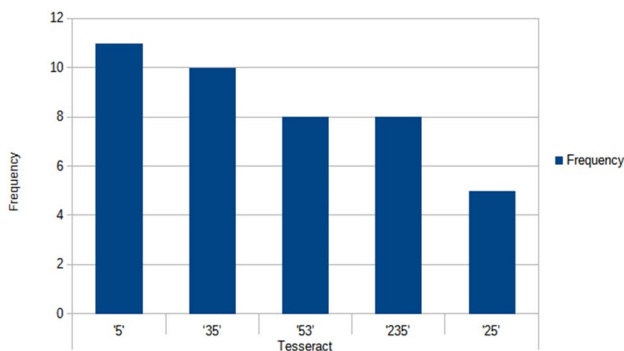


Fig. 5 Top five combinations w.r.t their frequency for Tesseract-OCR

Table 4 Method combination and the corresponding pre-processing techniques order

Method combination	Pre-processing name
'5'	Image smoothing
'35'	Gamma correction → Image smoothing
'53'	Image smoothing → Gamma correction
'235'	Brightness transformation → Gamma correction → Image smoothing
'25'	Brightness transformation → Image smoothing

Table 5 Overall comparison between the 3 OCRs

OE	A_WoP	A_WP	ET
Tesseract-OCR	85.67	92.46	5–10 s
Easy-OCR	75.86	78.84	3–5 min
docTR-OCR	81.29	86.24	5–10 min

OE OCR engine; *A_WoP* avg accuracy without pre-processing; *A_WP* avg accuracy with pre-processing; *ET* execution time (per-image)

position in terms searching, which is found to be 12% faster than ‘docTR-OCR’. In Fig. 5, top five combinations of pre-processing technique for Tesseract OCR are presented, whereas Table 4 shows the corresponding pre-processing techniques.

Table 5 depicted the overall comparison between the three OCRs based on average accuracy without pre-processing, with pre-processing and execution time (per-image). When the whole dataset is considered, it is found out that ‘Tesseract-OCR’ is 87% faster than ‘Easy-OCR’ and 89% faster in comparison with ‘docTR-OCR’. The ‘Easy-OCR’ is the second fastest engine, after ‘Tesseract-OCR’. Moreover, ‘Easy-OCR’ is 12% faster, compared with ‘docTR-OCR’. In Table 6, compatibility of a given pre-processing method with a given OCR is presented.

6 Conclusion and future work

In this paper, we present the detailed study, indicating the noticeable increase in the text extraction accuracy for all the OCR Engines, when the right set of pre-processing

techniques are applied before performing the text extraction. We have also identified the favorable pre-processing techniques for each OCR engine and the most optimal combinations amongst these pre-processing techniques to achieve the best text extraction accuracy. This work provides useful insights of a document with a structure similar to that of a medical lab report, where one of the pre-processing techniques can be chosen for optimal results, and that can help for further research in the field of text extraction. In addition, it supports in devising techniques for improved pre-processing, thus, improves the accuracy of text extraction. This work leverages, in optimizing the OCR engines for optimal lab report text recognition.

Due to the limited availability of processing power, experimentation is not performed on a larger dataset. However, experimenting on better computing infrastructure with larger dataset, we can observe better results which will be diverse and may offer the comparison of accuracies over the combinations between OCR engines, etc. Currently, our dataset consists of lab reports in PNG format, which can be extended to other formats. Moreover, apart from targeting on pre-processing combinations, future research could focus on the factors among the combinations that affects the accuracy for a certain OCR engine, if and only if detailed insights are available.

Funding No funding was received for conducting this study.

Data availability The data that support the findings of this study are available from the corresponding author [Jitendra Tembhurne], upon reasonable request.

Table 6 OCR engine and pre-processing technique compatibility

OE	CPT
Tesseract-OCR	BI: 56.41% BT: 74.36% GC: 79.49% SS: 51.28, ISm: 87.18, ISh: 12.82
Easy-OCR	BI: 2.57% BT: 92.31% GC: 61.54% SS: 97.44% ISm: 79.49% ISh: 28.21%
docTR-OCR	BI: 10.26% BT: 74.36% GC: 43.59% SS: 53.85% ISm: 74.36% ISh: 56.42%

OE OCR engine; *CPT* compatibility of pre-processing techniques (percentage of how many times a given method was chosen for the base-set); *BI* image binarization; *BT* brightness transformation; *GC* gamma correction; *SS* sigmoid stretching; *ISm* image smoothing; *ISh* image sharpening

Declarations

Conflict of interest The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Scott PJ, Curley PJ, Williams PB, Linehan IP, Shaha SH (2016) Measuring the operational impact of digitized hospital records: a mixed methods study. *BMC Med Inf Decis Mak* 16(1):1–13
2. Suter-Crazzolara C (2018) Better patient outcomes through mining of biomedical big data. *Front ICT* 5:30
3. Tawde GY, Kundargi J (2013) An overview of feature extraction techniques in OCR for Indian scripts focused on offline handwriting. *Int J Eng Res Appl* 3(1):919–926
4. Hamad K, Kaya M (2016) A detailed analysis of optical character recognition technology. *Int J Appl Math Electron Comput* 4:244–249
5. Karthick K, Ravindrakumar KB, Francis R, Ilankannan S (2019) Steps involved in text recognition and recent research in OCR; a study. *Int J Recent Technol Eng* 8(1):2277–3878
6. Shen M, Lei H (2015) Improving OCR performance with background image elimination. In: 2015 12th International conference on fuzzy systems and knowledge discovery (FSKD). IEEE, pp 1566–1570
7. Jain P, Taneja K, Taneja H (2021) Which OCR toolset is good and why: a comparative study. *Kuwait J Sci* 48(2)
8. de Mello CA, Lins RD (1999) A comparative study on OCR tools. In: *Vision interface*, vol 99, pp 224–231
9. Smith R (2007) An overview of the Tesseract OCR engine. In: *Ninth international conference on document analysis and recognition (ICDAR 2007)*, vol 2. IEEE, pp 629–633
10. Vitlani P, Kumbharana CK (2015) Comparative study of character recognition tools. *Int J Comput Appl* 118(9):31–36
11. Shafii M, Sid-Ahmed M (2015) Skew detection and correction based on an axes-parallel bounding box. *Int J Doc Anal Recogn (IJ DAR)* 18(1):59–71
12. Lin K, Li TH, Liu S, Li G (2019) Real photographs denoising with noise domain adaptation and attentive generative adversarial network. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*
13. List of Top 5 Open Source OCR Tools (2020). <https://www.hitec hnectar.com/blogs/open-source-ocr-tools/>. Accessed on 17th Oct 2022
14. Gupta B (2018). Improve accuracy of OCR using image preprocessing. <https://medium.com/cashify-engineering/improve-accuracy-of-ocr-using-image-preprocessing-8df29ec3a033>. Accessed on 17th Oct 2022
15. Improving the quality of the output (2021). <https://tesseract-ocr.github.io/tessdoc/ImproveQuality.html>. Accessed on 25th Oct 2022
16. Why is it important to digitize medical records? (2019). <https://www.managedoutsourcing.com/blog/why-is-it-important-to-digitize-medical-records/>. Accessed on 25th Oct 2022
17. Optical character recognition—OCR text recognition (2021). <https://www.v7labs.com/blog/ocr-guide>. Accessed on 30th Oct 2022
18. Devopedia (2019). Levenshtein distance. <https://devopedia.org/levenshtein-distance>. Accessed on 30th Oct 2022
19. EasyOCR (2021). <https://www.jaided.ai/easyocr/>. Accessed on 30th Oct 2022
20. Kannan P, Deepa S, Ramakrishnan R (2010) Contrast enhancement of sports images using modified sigmoid mapping function. In: 2010 International conference on communication control and computing technologies. IEEE, pp 651–656
21. Juneja K, Rana C (2020) Alignment and disruption robust binary mapper for optical Braille recognition. *Int J Inf Technol* 12(4):1291–1298
22. Joseph FJJ (2020) Effect of supervised learning methodologies in offline handwritten Thai character recognition. *Int J Inf Technol* 12(1):57–64
23. Rani U, Kaur A, Josan G (2019) A new binarization method for degraded document images. *Int J Inf Technol* 9(1):1–19
24. Sahare P, Tembhurne JV, Parate MR, Diwan T, Dhok SB (2023) Script independent text segmentation of document images using graph network based shortest path scheme. *Int J Inf Technol* 15(4):2247–2261
25. Lertsawatwicha P, Phathong P, Tantasanee N, Sarawutthinun K, Siriborvornratanakul T (2023) A novel stock counting system for detecting lot numbers using Tesseract OCR. *Int J Inf Technol* 15(1):393–398

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.