# Execution of OCR for Comparative Determination of Text from Images for Serif and Sans Serif Typefaces

Sukhpreet Kaur[1], Gagandeep Jagdev[2*]

*[1]Department of Computer Science, Guru Gobind Singh Khalsa College,
Bhagta Bhai Ka, Bathinda, Punjab, India.
[2*]Department of Computer Science, Punjabi University Guru Kashi College, Damdama Sahib, Punjab, India.
Email: [1]sukhpreetchanny50@gmail.com, [2*]drgagan137@pbi.ac.in*

*Abstract*−**Optical Character Recognition (OCR) pacts with the problem of identifying all types of different characters. OCR is capable of recognizing and converting both printed and handwritten characters into machine-readable digital formats. OCR is not only used to recognize proper words but can also read numbers and codes. The research paper focuses on performing the extraction of the text from the input image. The elaborated model has been proposed to conduct the extraction of the text from images. The implementation of the proposed model has been conducted on the MATLAB simulation tool. The three different cases with different font types (Times New Roman, Consolas, and Arial) have been tested on the proposed model and the readings for four different performance evaluation parameters (similarity achieved, difference, common symbols, and different symbols) have been achieved in all three cases. The paper elaborates the different steps of the conducted implementation in step-wise pictorial form for better understanding in all three implemented instances. The fonts have been taken from Serif (Arial and Times New Roman) and Sans Serif (Consolas) families. The conclusion part of the research paper details the obtained results in tabular as well as the graphical format and concludes that the accuracy achieved in extracting text from images is higher in the case of Serif fonts as compared to Sans Serif.**

**Keywords:***Arial, Consolas, Font, Text extraction, Times New Roman.*

## I. INTRODUCTION

The purpose of OCR is to conduct the mechanical and electrical conversion of scanned images of typewritten and handwritten text into machine text [1, 2]. It enables to digitize the printed texts which assist in electronic searches, enabling compact storage of text, and providing the text for machine translation, text mining, and text to speech conversion [3, 4]. In the recent past, the OCR found its widespread use in industries, banking, education, and research and brought a revolution in the world of the document management process [5, 6]. OCR has turned the scanned documents into fully searchable documents as the text is recognized via computers [7, 8]. OCR prevents the user from manually typing the documents if they are required to enter these into electronic databases. Instead, OCR excerpts pertinent information and enters it automatically [9, 10, 11]. Consider a code or a serial number comprising of alphabets and numerals that are supposed to be digitized. OCR can convert these codes into a digital form and provide the required output [12, 13]. But OCR does not consider the genuine nature of the object which is to be scanned [14, 15]. The job of the OCR is to take a look at the characters and convert them into digital format [16]. For instance, if one scans the word, OCR would learn and recognize the letters of the word, but have nothing to do with the meaning of the word [17]. The research work conducted in the paper primarily focuses on recognizing and extracting the alphabets and numerals from the given image based on the input image comprising alphabets and numerals to form the template. The four different performance evaluation parameters have been used to test the worth of the conducted research. These parameters are briefly mentioned below.

### A. Similarity Achieved

*"Similarity achieved"* refers to how accurately the text has been extracted from the image under consideration. The value of the *"similarity achieved"* is obtained in percentage. The greater is the value of *"similarity achieved"*, the better are the results obtained.

### B. Difference

*"Difference"* refers to the text that has not been accurately identified in the image under study. The lower is the value of the parameter "difference", the better are the results obtained. The formula for evaluating the value of the difference is mentioned as under.

Difference = 100 – Similarity achieved

## C. Common Symbols

*"Common symbols"* refers to the characters and numerals which have been accurately matched between the input image and the obtained output. Greater is the number of the *"common symbols"*, better are the results obtained.

## D. Different Symbols

"*Different Symbols*" refers to the characters and numerals which have been mismatched between the input image and the obtained output. The lower is the value of "different symbols", the better are the results obtained.

## II. Research Methodology

This section elaborates on the adopted research methodology to extract the text from the images. The proposed methodology is elaborated in the flowchart displayed in Fig. 1 followed by the relevant algorithm.



Fig. 1: Flowchart Depicts the Proposed Methodology for Extracting the Text from the Image

## Algorithm

1. Provide the input image *img1*.
2. Crop the image.
3. Convert the image *img1* to binary (1 for alphabet or numeral and 0 for space).
4. Upload the cropped image and generate a Template.
5. Generated Template should have created a .mat file.
6. Read the image *img2* from which the text is to be extracted.
7. Read one line per loop from the image *img2*.
8. Read the image *img2* from which the text is to be extracted.
9. Read one line per loop from the image *img2*.
10. Set the boundary for *img2*.
11. Create a text file to save results.
12. Load the image *img2* to the created template.
13. Analyze the number of words in line through spaces between letters.
14. Read letters in their original dimensions.
15. Estimate the space occupied by the letters.
16. Remove the line that has been read.
17. If there more lines left in the image *img2*
18. Goto 13
19. Else
20. Goto 18
21. End

## III. Implementation and Results

This section practically demonstrates the proposed research work via three instances. MATLAB has been used as a simulation tool to extract the text from the images.

## A. Case 1

Input images
First input image — *img1* (Fig. 2)
Second input image — *img2* (Fig. 3)
Font — Times New Roman
Font size — 18

Fig. 2 represents the image *img1* showing the lowercase characters from 'a' to 'z', uppercase characters from 'A' to 'Z', and numerals from '0' to '9' written in Times New Roman font having font size 18.

a b c d e f g h i j k l m n o p q r s t u v w x y z

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z

0 1 2 3 4 5 6 7 8 9

Fig. 2: Figure Depicts the Input Image *img1* with Written Characters and Numerals Using Times New Roman Font

Fig. 3 shows the second input *img2* which would be provided as input to the proposed model after successful creation of the template. The text has to be extracted from *img2*.

The purpose of OCR is to conduct the mechanical and electrical conversion of scanned images of typewritten and handwritten text into machine text. It enables to digitize the printed texts which assists in electronic searches, enabling compact storage of text, and providing the text for machine translation, text mining, and text to speech conversion. In recent past, the OCR found its widespread use in industries, banking, education, and research and brought a revolution in the world of document management process. OCR has turned the scanned documents into fully searchable documents as the text is recognized via computers. OCR prevents the user from manually typing the documents if they are required to enter these into electronic databases. Instead, OCR excerpts pertinent information and enters it automatically.

Fig. 3. Figure shows the second input *img2* which would be provided as input to the proposed model

Fig. 4 represents the GUI (Graphical User Interface) created to execute the proposed methodology. The GUI comprises three push buttons on the left side of the form titled *"ocr_start"*. The first push-button *"Click to Create Template"* is dedicated to the creation of the template via uploading the input image file *img1*. The second push button titled *"Click to initiate with text extraction"* is intended to initiate the process of extracting the text from the second input image file *img2*. The third push-button titled *"Exit"* terminates the GUI window.
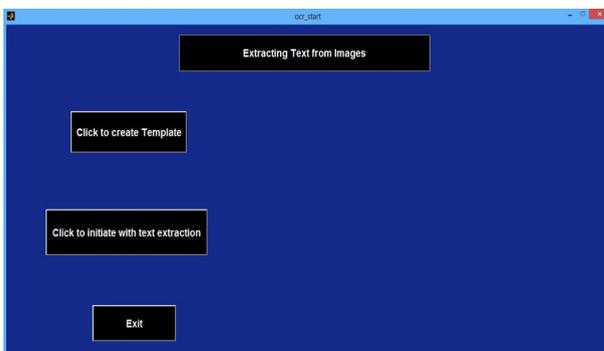


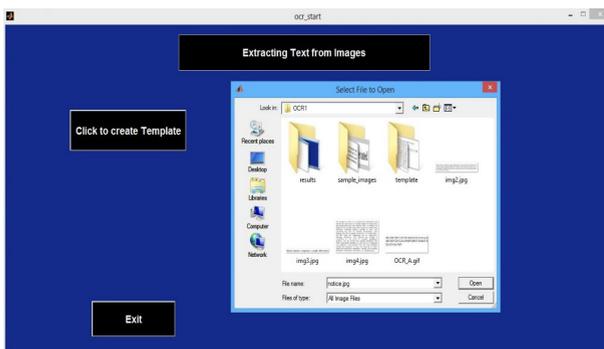Fig. 4: Figure depicts the designed GUI for extracting the text from the image



Fig. 5: Figure depicts the browsing and uploading of the image *img1* into the proposed system

Fig. 5 depicts the browsing and uploading of the image *img1* into the proposed system.

Fig. 6 shows the uploaded image *img1*. A menu with two options of *"Back"* and *"Continue"* appears on the screen. The *"Back"* button redirects to the previous stage and the *"Continue"* button leads to the creation of the template.

Fig. 7 pops the message confirming the successful creation of the template.
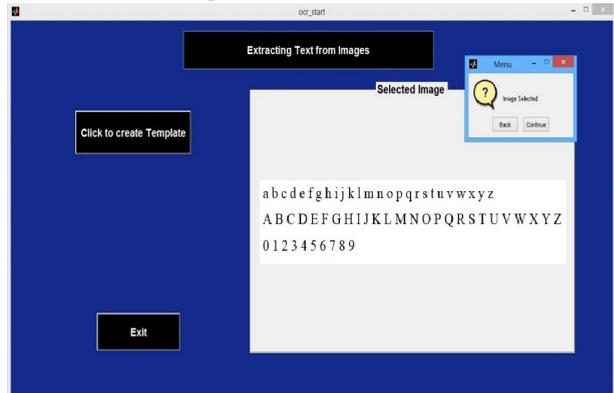


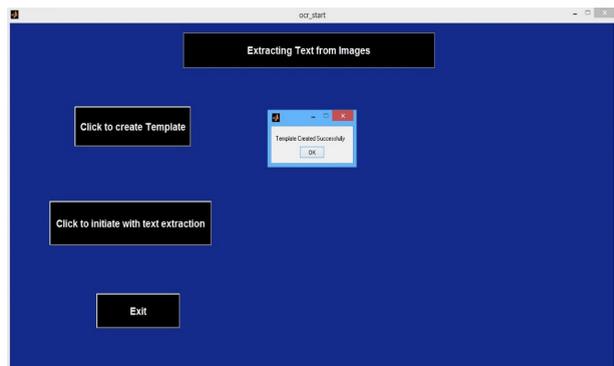Fig. 6: Figure shows the uploaded image *img1* on GUI



Fig. 7: Figure displays the popped message confirming the successful creation of the template

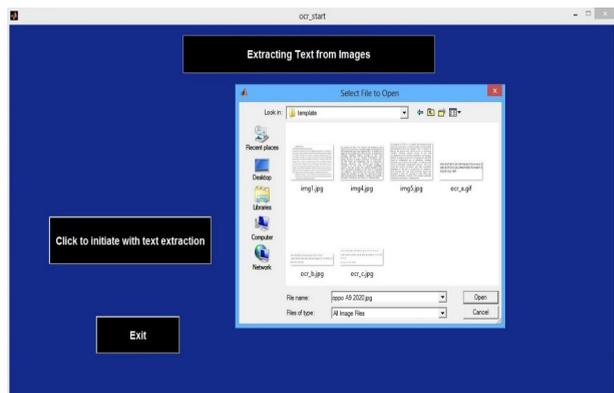Fig. 8 depicts the browsing and uploading of the image *img2* from which the text is to be extracted.



Fig. 8: Figure depicts the browsing and uploading of the image *img2*

Fig. 9 shows the image *img2* successfully uploaded for the extraction of text.
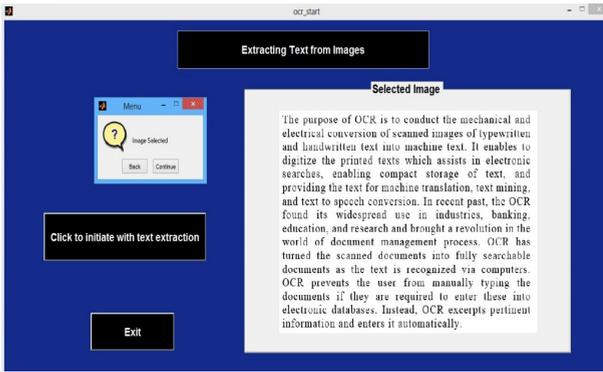


Fig. 9: Figure depicts the image *img2* been successfully uploaded

Fig. 10 shows the confirmation of the text been successfully extracted from the image *img2*.
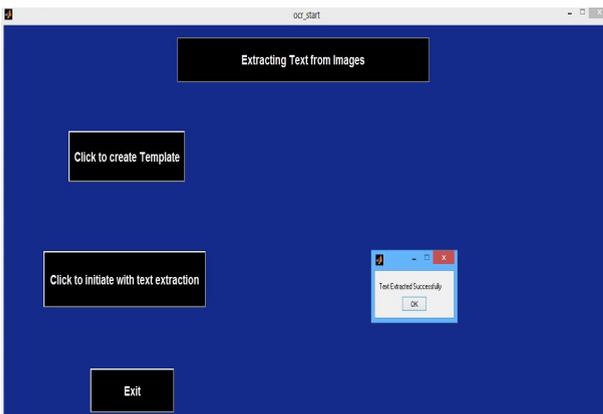


Fig. 10: Figure shows the confirmation of the text been successfully extracted

Fig. 11 shows the finally extracted text from the image *img2* as per the proposed methodology.
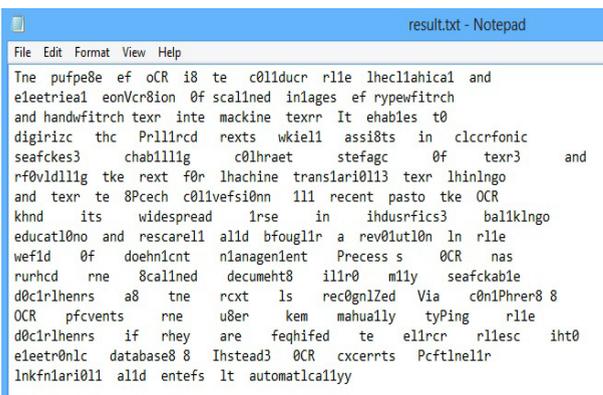


Fig. 11: Figure depicts the extracted text from the image *img2*

*Results obtained:*
Similarity achieved – 58.98%
Difference – 41.025%

Common symbols – 611
Different symbols – 425

### B. Case 2

Input images

First input image       – *img1* (Fig. 12)
Second input image    – *img2* (Fig. 13)
Font                 – Consolas
Font size            – 18

Fig. 12 represents the image *img1* showing the lowercase characters from 'a' to 'z', uppercase characters from 'A' to 'Z', and numerals from '0' to '9' written in Consolas font having font size 18.

Fig. 13 shows the second input *img2* which would be provided as input to the proposed model after successful creation of the template. The text has to be extracted from *img2*.



Fig. 12: Figure depicts the input image *img1* with written characters and numerals using Times New Roman Font



Fig. 13: Figure shows the second input *img2* which would be provided as input to the proposed model

Fig. 14 shows the uploaded image *img1*. A menu with two options of *"Back"* and *"Continue"* appears on the screen. The *"Back"* button redirects to the previous stage and the *"Continue"* button leads to the creation of the template.

Fig. 15 pops the message confirming the successful creation of the template.

Fig. 16 shows the image *img2* successfully uploaded for the extraction of text.

Fig. 17 shows the finally extracted text from the image *img2* as per the proposed methodology.

*Results obtained:*

| | |
|---|---|
| Similarity achieved | – 49.22% |
| Difference | – 50.78% |
| Common symbols | – 603 |
| Different symbols | – 622 |

## C. Case 3

Input images
First input image – *img1* (Fig. 18)
Second input image – *img2* (Fig. 19)
Font – Arial
Font size – 18

Fig. 18 represents the image *img1* showing the lowercase characters from 'a' to 'z', uppercase characters from 'A' to 'Z', and numerals from '0' to '9' written in Arial font having font size 18.

Fig. 19 shows the second input *img2* which would be provided as input to the proposed model after successful creation of the template. The text has to be extracted from *img2*.

Fig. 20 shows the uploaded image *img1*. A menu with two options of *"Back"* and *"Continue"* appears on the screen. The *"Back"* button redirects to the previous stage and the *"Continue"* button leads to the creation of the template.
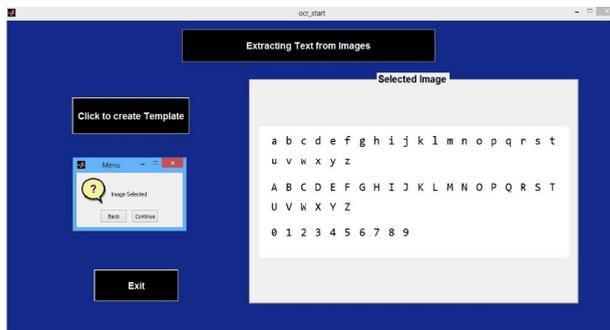


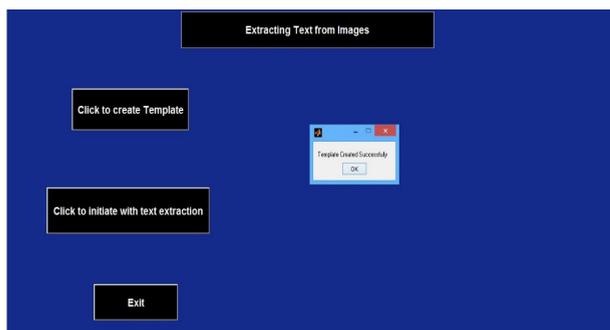Fig. 14: Figure shows the uploaded image *img1* on GUI



Fig. 15: Figure displays the popped message confirming the successful creation of the template
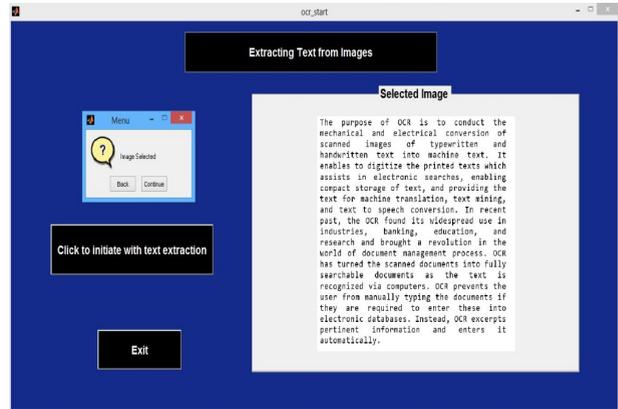


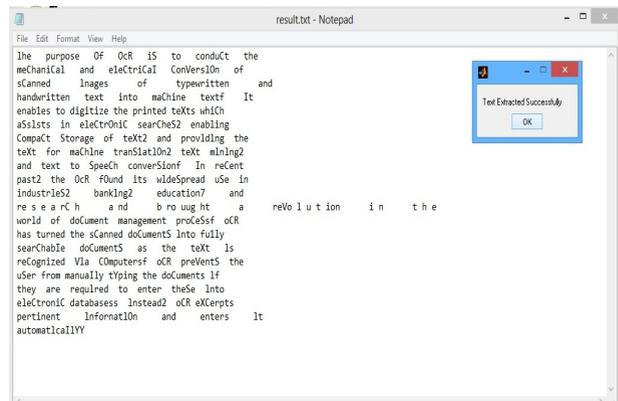Fig. 16: Figure depicts the image *img2* been successfully uploaded



Fig. 17: Figure depicts the extracted text from the image *img2*



Fig. 18: Figure depicts the input image *img1* with written characters and numerals using Arial

The purpose of OCR is to conduct the mechanical and electrical conversion of scanned images of typewritten and handwritten text into machine text. It enables to digitize the printed texts which assists in electronic searches, enabling compact storage of text, and providing the text for machine translation, text mining, and text to speech conversion. In recent past, the OCR found its widespread use in industries, banking, education, and research and brought a revolution in the world of document management process. OCR has turned the scanned documents into fully searchable documents as the text is recognized via computers. OCR prevents the user from manually typing the documents if they are required to enter these into electronic databases. Instead, OCR excerpts pertinent information and enters it automatically.

Fig. 19: Figure shows the second input *img2* which would be provided as input to the proposed model
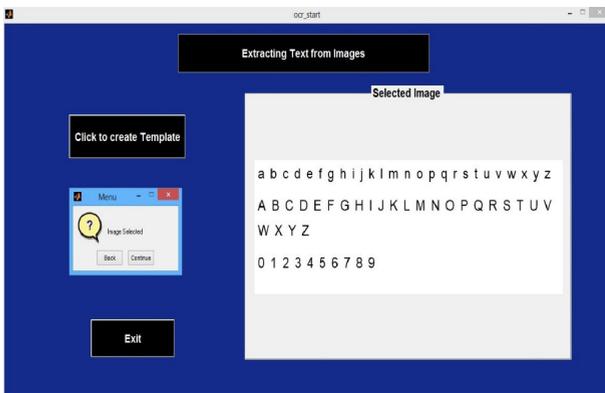
Fig. 20: Figure shows the uploaded image *img1* on GUI

Fig. 21 pops the message confirming the successful creation of the template.

Fig. 22 shows the image *img2* successfully uploaded for the extraction of text.

Fig. 23 shows the finally extracted text from the image *img2* as per the proposed methodology.

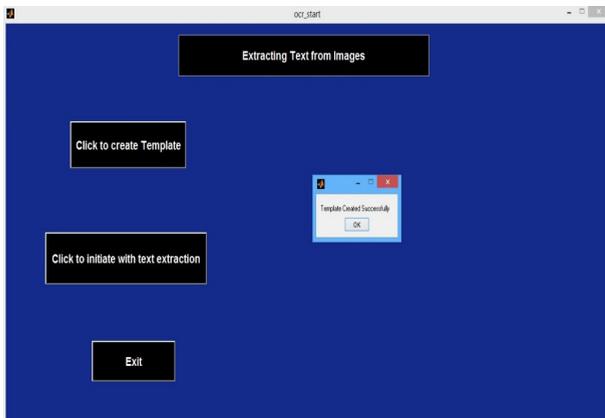Fig. 24 depicts the extracted text from the image *img2*



Fig. 21: Figure displays the popped message confirming the successful creation of the template
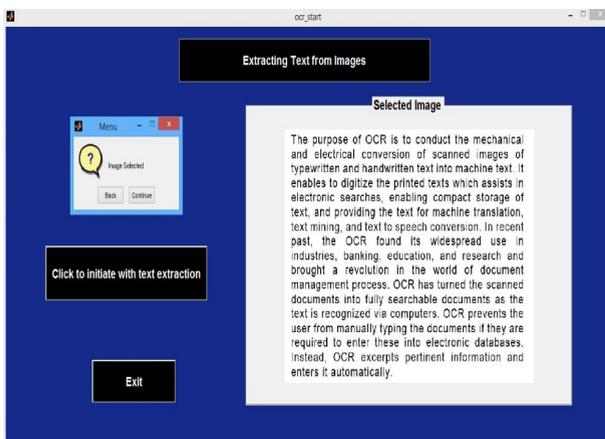


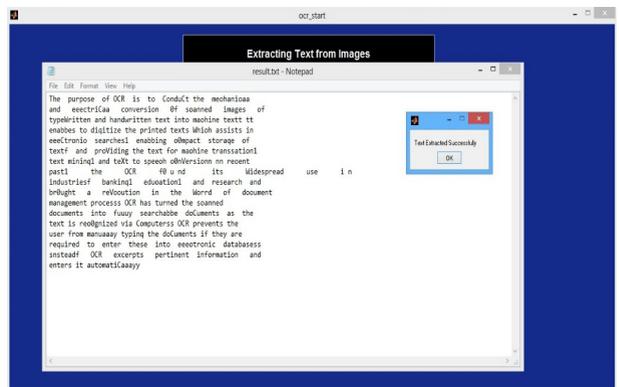Fig. 22: Figure depicts the image *img2* been successfully uploaded



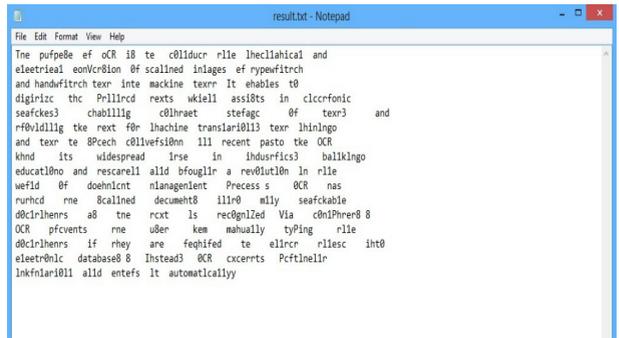Fig. 23: Figure shows the finally extracted text from the image *img2* as per the proposed methodology



Fig. 24: Figure depicts the extracted text from the image *img2*

*Results Obtained:*

| | |
|---|---|
| Similarity achieved | – 61.35% |
| Difference | – 38.65% |
| Common symbols | – 665 |
| Different symbols | – 419 |

## IV.   CONCLUSION

The proposed model for extracting the text from the images has been tested on three different cases in Section III. The values for four performance evaluation parameters have been obtained in all three cases and have been summarized in Table 1 below.

Table 1. Table denotes the obtained values of different performance evaluation parameters from three tested cases

| Font Name | Similarity achieved (%) | Difference (%) | Common symbols | Different symbols |
|---|---|---|---|---|
| Times New Roman | 58.98 | 41.025 | 611 | 425 |
| Consolas | 49.22 | 50.78 | 603 | 622 |
| Arial | 61.35 | 38.65 | 665 | 419 |

The obtained results indicate that the maximum similarity in percentage has been obtained in with font Arial followed by Times New Roman and finally followed by Consolas. The greater value of the "similarity achieved" parameter represents the better extraction of the text from the images.

Fig. 25 denotes the comparative graphical representation of the parameter "*similarity achieved*" in percentage among the three tested cases.

Fig. 26 denotes the comparative graphical representation of the parameter "*difference*" in percentage among the three tested cases.
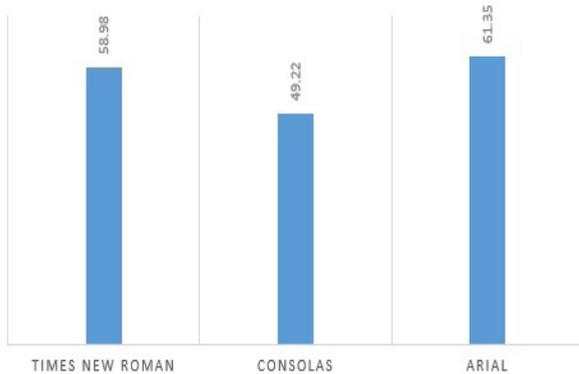
**SIMILARITY ACHIEVED (%)**

Fig. 25: Figure depicts the proportional graphical representation of the parameter "*similarity achieved*" among the three tested cases

Fig. 27 denotes the comparative graphical representation of the parameter "*common symbols*" among the three tested cases.

Fig. 28 denotes the comparative graphical representation of the parameter "*different symbols*" among the three tested cases.
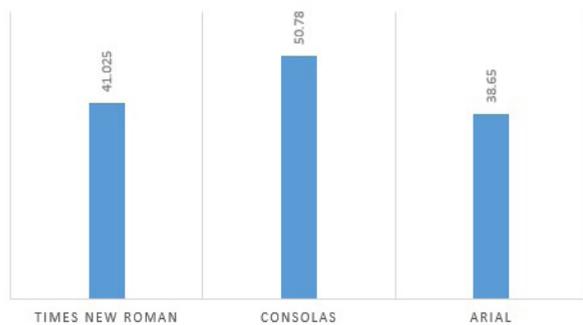
**DIFFERENCE (%)**

Fig. 26: Figure denotes the proportional graphical representation of the parameter "*difference*" in percentage among the three tested cases
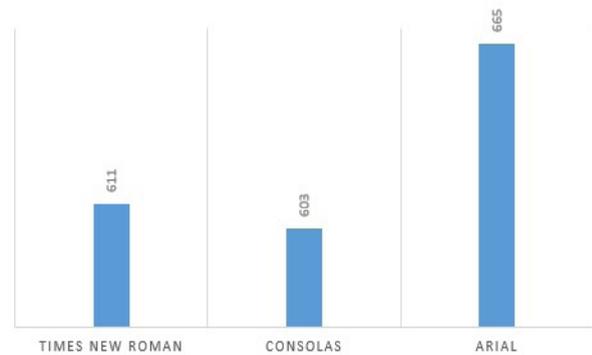
**COMMON SYMBOLS**

Fig. 27: Figure denotes the proportional graphical representation of the parameter "*common symbols*" among the three tested cases
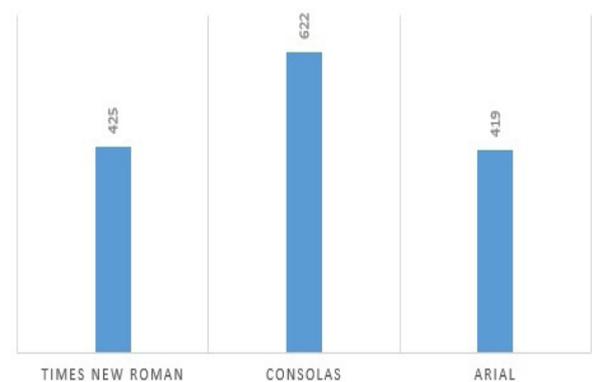
**DIFFERENT SYMBOLS**

Fig. 28: Figure denotes the comparative graphical representation of the parameter "*common symbols*" among the three tested cases

So, based on the conducted research work it can be concluded that the accuracy in terms of extraction of text from the images is best conducted with the fonts falling under the Serif category (Arial and Times New Roman) as compared to Sans Serif (Consolas). In the future, many more fonts with varying font sizes can be used to obtain check out the best among them in terms of different performance evaluation parameters. Tesseract with Python can be used to future enhance the values of participating performance evaluation parameters.

REFERENCES

[01]  Ali, A.A.A., Suresha M.: A novel approach to correction of a skew at document level using an arabic script. Int J Comput Sci Inf Technol. 8(5), 569–573 (2017).

[02]  Abdullah M.A., Al-Harigy L.M., Al-Fraidi H.H.: Off-line arabic handwriting character recognition using word segmentation. J Comput. 4,40–44 (2012).

[03] Elaiwat S., Abu-zanona M.A.: A three stages segmentation model for a higher accurate off-line Arabic handwriting recognition. World Comput Sci Inf Technol J. 2(3), 98–104 (2012).

[04] Adiguzel H., Sahin E., Duygulu P.: A hybrid for line segmentation in handwritten documents. In: International conference on frontiers in handwriting recognition (ICFHR), pp. 503–508, (2012).

[05] Banerjee S., Mullick K., Bhattacharya U.: A robust approach to extraction of texts from camera captured images. In: International Workshop on Camera-Based Document Analysis and Recognition, pp. 30 -46, SpringerWashington, (2013).

[06] Smith R.: An overview of the Tesseract OCR engine. In: Document Analysis and Recognition, 2007. ICDAR 2007. pp. 629–633, IEEE,Paraná(2007).

[07] Qiaoyang, Y., Doermann, D.: Text detection and recognition in imagery: A survey. IEEE Trans. Patt. Anal. Mach. Intell.37, 1480–1500 (2015).

[08] Wang S., Fu C., Li Q.: Text detection in natural scene image: a survey. In: International Conference on Machine Learning and Intelligent Communications. pp. 257–264, Springer, Shanghai (2016).

[09] Zhang, H., Zhao, K., Song, Y.Z., Guo, J.: Text extraction from natural scene image: a survey. Neurocomputing. 122, 310–323 (2013).

[10] Bai, X., Yao, C., Liu, W.: Strokelets: a learned multi-scale mid-level representation for scene text recognition. IEEE Trans. Image Process.25, 2789–2802 (2016).

[11] Girshick, R.: Fast R-CNN. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1440–1448. IEEE, Las Condes (2015).

[12] Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: unified, real-time object detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.pp. 779–788. IEEE, Las Vegas (2016).

[13] Zhang, Z., Zhang, C., Shen, W., Yao, C., Liu, W., Bai, X.: Multi-oriented text detection with fully convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4159–4167. IEEE, Las Vegas (2016).

[14] Shi, B., Bai, X., Belongie, S.: Detecting oriented text in natural images by linking segments. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2550–2558. IEEE (2017).

[15] Wang, Y., Shi, C., Xiao, B., Wang, C., Qi, C.: CRF based text detection for natural scene images using convolutional neural network and context information. Neurocomputing. 295, 46–58 (2018).

[16] Ubul, K., Tursun, G., Aysa, A., Impedovo, D., Pirlo, G., Yibulayin, T.: Script identification of multi-script documents: a survey. IEEE Access. 5, 6546–6559 (2017).

[17] Gomez, L., Nicolaou, A., Karatzas, D.: Improving patch-based scene text script identification with ensembles of conjoined networks. Patt. Recogn.67, 85–96 (2017).