

Extraction of information from bill receipts using optical character recognition

Vedant Kumar, Pratyush Kaware,
Pradhuman Singh, Dr. Reena Sonkusare
Electronics & Telecommunications Department
Sardar Patel Institute of Technology
Mumbai, India

Siddhant Kumar
Electronics Department
Sardar Patel Institute of Technology
Mumbai, India

Abstract—This paper presents an application of optical character recognition (OCR) which can extract information from images of bills and receipts; store it as machine-processable text; in an organized manner for ease of access. It can do this efficiently even in the presence of watermarks on the bills or any shadows in the images of the bills. In developing this application, OpenCV has been used for the processing of the images and the Tesseract OCR engine has been used for optical character recognition and text extraction. The image is first processed using OpenCV for the removal of any shadows or watermarks present in it. For longer invoices, by employing the image bifurcation process, the data can be easily extracted which was not possible earlier. Furthermore, the processed image is passed on to the Tesseract OCR engine for the retrieving of text present in it. The text is then searched for important information, such as the total amount spent and the date on the receipt, using string processing.

Index Terms—Image Processing, Optical Character Recognition, OpenCV, Tesseract OCR Engine

I. INTRODUCTION

Paper bills have been traditionally used to store information about financial transactions. Although this is a reliable way of storing information, it is very time-consuming to search through them and looking for a specific transaction is tedious. These types of tasks can be done easily by a computer, but the paper bills cannot be stored digitally without manually typing their contents into a computer, which in itself is a very time-consuming task. So there is a need for extracting and storing this information from the bills automatically. This can be achieved by using optical character recognition. In this day and age, optical character recognition has reached a level where it can give a very accurate output from images of printed text, to say, in the same order and with the spacing of the words as found in the images. So an optical character recognition can be used to store this type of printed information as machine-processable text which would make it easily accessible using a computer. Tools such as OpenCV and Tesseract OCR engine have been used to develop the system mentioned above, and have finally implemented it as an Android application. Any android phone user can use this app to convert images of their bills to the text that would then be stored in an organized manner. The app works offline so that it does not depend on the availability of the internet. This further reduces the time required to retrieve text from

an image of a bill. It is simple to use as the user has to click a picture of the bill from a reasonable distance and the information in the bill would be stored. The information on the bills can be easily acquired according to their date.

II. LITERATURE SURVEY

OCR (Optical Character Recognition) is an engine that enables the recognition and retrieval of text information of printed or handwritten character units from images or scanned documents that contain texts. OCR as a process consists of various sub-processes, namely, text localization, character segmentation, and character recognition. Although some methods work without segmentation [1]. Such systems are developed using various techniques, such as using least square support vector machines to classify the characters [2]. K-Nearest neighbor models have also been used [3]. But mostly to recognize the presence of a single character in an image, Convolutional Neural Networks (CNN) are being used [4]. Text is an arbitrary sequence of characters, and detecting these characters with higher accuracy is a problem that can be solved using recurrent neural networks (RNNs) [5] and long short term memory (LSTM) that is a popular form of RNN. Words are obtained by organizing text lines into blobs. These lines and regions are further analyzed for proportional text. Text lines are separated into words individually according to the kind of spacing between them. The process of recognition is divided into two steps. In the first step, each word is recognized. Each perfect word is further passed to an adaptive classifier as training data. The input image is analyzed and processed in parts line by line feeding into the LSTM model. Tesseract is an optical character recognition engine that is available for a variety of operating systems. It uses CNN and LSTM to detect and extract texts from images accurately [6].

Now, for the images of a text with noise or shadow, OCR engines like Tesseract have a noticeable drop in performance [7]. Training the models on blurred text images help improve the accuracy and reduce the error rate in such cases [8]. Pre-processing on the image can be done using OpenCV to minimize the noise. Open CV is an open-source software [9] library that holds applications in the domains of Machine learning and computer vision. The processing of images can be done

by incorporating the OpenCV library. With the processing step, the images will be more clear, and the relevant data can be extracted efficiently. This pre-processing step can involve determining the region of interest and cropping the image accordingly, noise removal, thresholding, dilation, erosion, and contour detection. Once these steps are performed, now OCR engines can be used to read the image and extract the text from it accurately.

III. TOOLS USED

A. OpenCV

OpenCV was a library originally for C/C++ that has been widely used for processing image data. There are many pre-defined useful functions in this library for doing some basic transformations on an image. A ported version of this library is also available for android development in Java which has been used for the initial processing of the image before passing it to the extraction stage.

B. Tesseract OCR Engine

It is an open-source library maintained by Google for text recognition. It can detect text in many languages and gives the text in formatted form i.e. in the order that it appears in the input image. It provides the output text quickly after giving it the input of an image of a reasonable length. This library is also available for android development in the form of an SDK which can be used as a dependency in an android app.

IV. METHODOLOGY

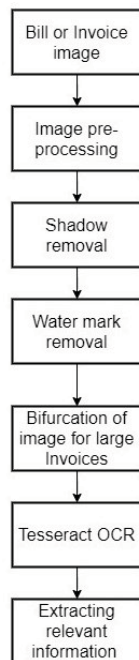


Fig. 1: Process flow of methodology

As shown in Fig. 1, our proposed system uses the following steps for getting the data from a bill or receipt:

- (1) Initially, basic processing steps are being done like contour detection, noise removal, and functions like erosion and dilation.
- (2) The watermark and shadows are then being removed from the bill.
- (3) Then segmenting the bill into parts.
- (4) Passing it through the Tesseract OCR engine to get the whole text.
- (5) Finally getting the important information like the total amount and the date of the bill using regular expressions.

A. Watermark and shadow removal

Bills generally have watermarks as shown in Fig. 10 marked in red circles. These watermarks hinder the process of text extraction i.e. causes text on the watermark to be missed. Sometimes these watermarks contain text themselves that can act as noise in the image of a bill as Tesseract detects text of all sizes in a line that may cause it to miss the overlapping text on the watermark.

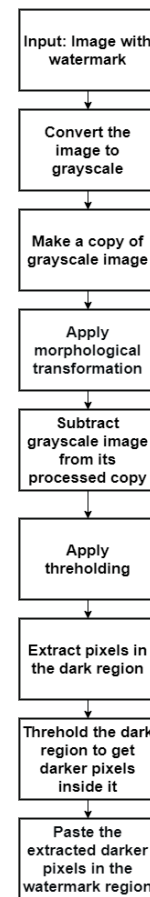


Fig. 2: Process flow of watermark removal

Shadows can also cause the efficiency of the Tesseract OCR Engine to reduce significantly, but these shadows can be removed by filtering out the noise. Shadows can be removed in

some way similar to watermark removal, i.e. by increasing the contrast and brightness of the text image as the text containing the information about the transactions are generally in darker colors.

For images with watermarks, series of the process including conversion of image to grayscale, morphological transformations, thresholding(binary inversion and otsu transformation), extracting pixels in the dark region to get the darker pixels inside it, and then pasting the darker pixels obtained in the watermark region. For invoice images with shadow, the grayscale input image is dilated so as to get rid of the texts present in the invoice. This step, however, preserves the bar code. Additionally, on applying median blur with a decent sized kernel further texts can be suppressed. With this, a good background image remains that contains all the shadows or any discolorations. On computing difference between the original and the background image obtained, the identical bits will be black (minute difference), and the text region will be white (large difference). This result is inverted as black on white is preferable. Finally, after applying thresholding to remove gray regions and normalizing the output, image without any shadows present is obtained. The entire process is summarized in Fig. 2 and Fig 3.

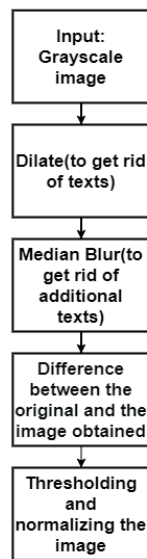


Fig. 3: Process flow of shadow removal

B. Bifurcation of long Invoices

For long invoices (longer than 26cm in length), the algorithm divides the invoice into two parts as shown in Fig. 4, so that extracting data becomes convenient. The algorithm scans the invoice and calculates the dimension of the entire invoice. For the bifurcation process, the invoice is divided into two parts-the first 55 percent of its length and then the last 50 percent of the length. The 5 percent in the middle, is preserved for information retention. In this way, the data of the

invoice can be extracted without any loss of information and the technique used is unique and efficient for longer invoices.



Fig. 4: Image bifurcation and then processing the resized images for better extraction of data.

C. Text Recognition and Extraction

A CNN model can be created and trained on printed text found in bills and this model can be used for detecting text from other bills with similar font [10]. The Tesseract OCR engine has been used to retrieve text from the processed bill images. For OCR, text localization, character segmentation, and character recognition needs to be done. All this is done by the Tesseract OCR engine. The processes done by Tesseract are shown in Fig. 5. It can extract lines of text with word and line segmentation in the same format i.e. the words appear in the same order as in the images. Tesseract OCR engine provides very high accuracy when working with printed text rather than handwritten text.

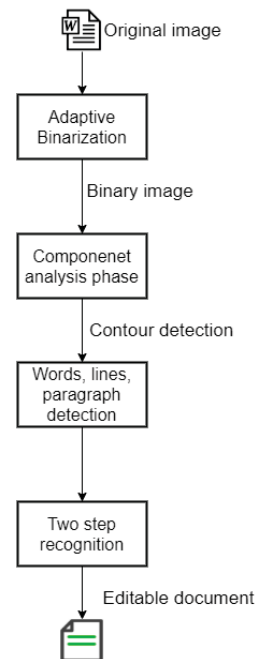


Fig. 5: Process flow of extraction of information by Tesseract

D. Getting Relevant Information

From the whole text, information like the date of the bill can be easily extracted using various regular expressions, for all the possible date formats. The total amount value from the bill can also be extracted using regular expressions owing it to the fact that it appears at the end of the bill. This information would be stored according to the date so that it is easily accessible.

V. SOFTWARE DESIGN

The proposed system can be hosted on a server and then could be connected to an Android app. This method would do all the processing on a remote server [11]. But, nowadays a standard smartphone has enough processing power required for our proposed system to process directly in it. That is why an offline approach have been taken which means all the processing would place locally on the phone without the need for an internet connection. Tesseract-android-tools SDK for text recognition and OpenCv4Android SDK(version 2.4.11) have been used for processing the image.



Fig. 6: An example bill which has been imported in the android application.

The user can use already taken photos or take a new photo of a bill then import the bill in the application. The bill shown in Fig. 6 has been input into the app and the extracted text which will be displayed and stored on the phone. The order of the text is not affected in the stored text as can be seen in Fig. 7.

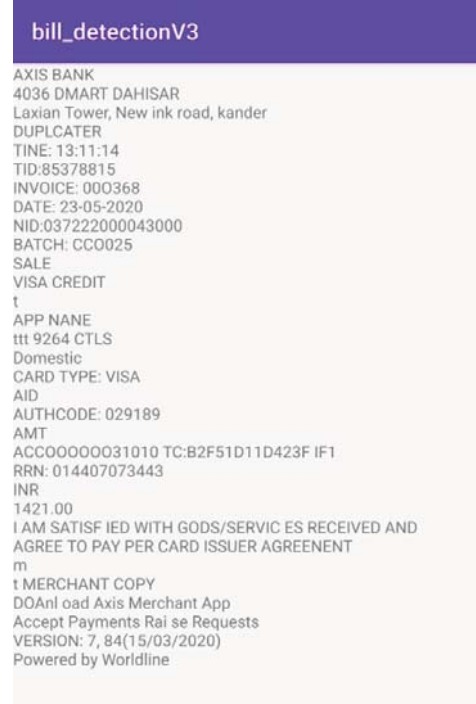


Fig. 7: The extracted text is shown retaining the format from the inputted image of a bill shown in Fig. 6

VI. EVALUATION METRICS

A. Accuracy

Accuracy is defined as the number of words that are extracted by the Tesseract OCR and that are present in the invoice to the total number of words and characters present in the invoice. Higher the accuracy, higher is the efficiency of pre-processing techniques to extract the data efficiently. $x\%$ accuracy means out of 100 characters present in the original invoice, OCR was able to extract x characters efficiently. $(100-x)\%$ would give the error made by the OCR in either misclassifying the characters or inability to extract the characters.

B. Peak Signal-to-Noise ratio (PSNR)

PSNR expression is the ratio between the maximum power of a signal to the power of distorting noise that weakens its quality. The value of PSNR gives information about the quality of image. Higher the value of PSNR indicates better information retention before and after the watermark has been removed from the invoice. Value of PSNR is inversely related to the mean squared error or MSE. MSE computes the comparison of the actual pixel values of the original image with the image retrieved after removing the watermark. The error so obtained represents the amount by which the values of the original image differ from the retrieved image.

$$PSNR = 20 \log_{10} \left(\frac{MAX_f}{\sqrt{MSE}} \right)$$

Fig. 8: Representation of PSNR mathematically.

MAX_f in Fig. 8 represents the maximum value of the pixel in the original image (grayscale version of the original image). In the representation of MSE in Fig. 9, $f(i,j)$ denotes pixel values of original image and $g(i,j)$ denotes pixel values of a retrieved image. The values m and n represent the number of rows and columns respectively.

$$MSE = \frac{1}{mn} \sum_0^{m-1} \sum_0^{n-1} \|f(i,j) - g(i,j)\|^2$$

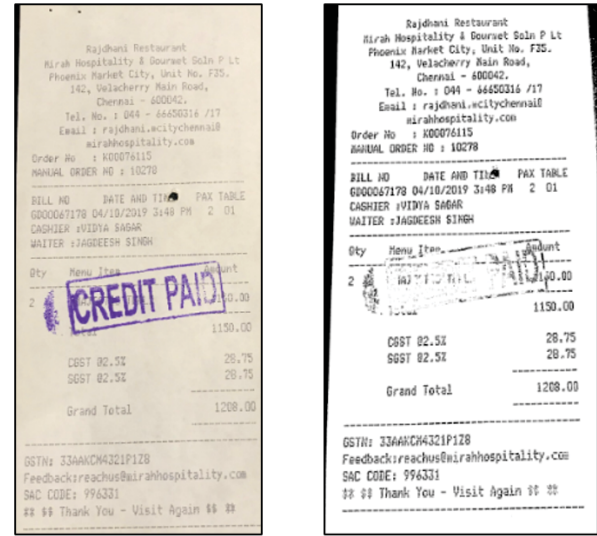
Fig. 9: Representation of MSE mathematically.

VII. RESULT ANALYSIS

The effect of increasing the contrast can be seen in Fig. 10 i.e. the removal of watermarks and shadow. This shows that the effect of watermarks and shadows can be significantly reduced by processing the image first. Sometimes there are stamps on bills which are very similar to watermarks. Fig. 11 shows that these can also be reduced using the same technique used for removal of watermarks and shadows.



Fig. 10: Before and after processing example of a bill with watermarks and shadows.



Raw Image

Processed Image

Fig. 11: Before and after processing example of a bill with a stamp on it.

TABLE I: Results of PSNR

Test Bill Image	PSNR
Fig. 10	38.13 dB
Fig. 11	42.57 dB

For the PSNR range shown in TABLE I, which is expected to be between 30 dB to 50 dB, a value of 38.13 dB is achieved for Fig. 10 and a value of 42.57 dB is achieved for Fig. 11 which indicates good reconstruction of an overall image after removing the watermark.

TABLE II: Results and Comparison

Methodology	Accuracy for Small bills	Accuracy for Large bills
Proposed System with image pre-processing for shadow and watermark removal	97%	83%
Android-Based Text Recognition on Receipt Bill for Tax Sampling System [11] (Small similar to Nearer to camera, Large similar to Farther away from it)	95%	85%

This process does not take much time because it is offline and does not depend on having a stable internet connection. Many bills and receipts can be stored this way in chronological order by the date that can be found in the bills. Having tested the entire process on 25 bills, an average matching ratio of 90% has been gained.

The results and its comparison to other systems is shown in TABLE II. The proposed system is has a accuracy of 97% when working with small bills and 83% when on larger bills. The reduced accuracy in larger bills is partly due to the fact that the font appears smaller when trying to fit the whole bill in the image. We have also found that the system never misses or incorrectly reads the total amount or the date on the bill which is being extracted. This confirms that the accuracy of the overall process is high where it is needed.

VIII. CONCLUSION

This system shows us that technologies like optical character recognition can be used information from bills accurately. This is especially helpful in keeping track of financial transactions done by cash. This ensures the ease of storage and accessibility of such information as they are now easier to search through. With all this achieved, this application holds value in retail shops and household where tracking bills is a problem with an additional feature of tracking one's budget which is an important component that is being solved.

The proposed system could be further developed to be for verifying financial documents i.e. to find any discrepancies in the documents. This would require some complex regular expressions to process the extracted text for all the numeric data present and to find a relationship between them. Tables can also be detected in their correct format to store tables that are frequently found in financial documents [12]. Extracting information from old bills is very difficult as the printed text starts to smudge or become lighter. Some techniques use Augmented Reality and Virtual reality [13] for the restoration of this type of partial text. The conversion of text images to computer-readable text makes the task of text-to-speech conversion for the visually impaired people very simple [14] and can even be used to recite whole books [15].

REFERENCES

- [1] M. A. Ozdil and F. T. Y. Vural, "Optical character recognition without segmentation," *Proceedings of the Fourth International Conference on Document Analysis and Recognition*, Ulm, Germany, 1997, pp. 483-486 vol.2, doi: 10.1109/ICDAR.1997.620545.
- [2] J. Xie, "Optical Character Recognition Based on Least Square Support Vector Machine," *2009 Third International Symposium on Intelligent Information Technology Application*, Shanghai, 2009, pp. 626-629, doi: 10.1109/IITA.2009.327.
- [3] T. K. Hazra, D. P. Singh and N. Daga, "Optical character recognition using KNN on custom image dataset," *2017 8th Annual Industrial Automation and Electromechanical Engineering Conference (IEMECON)*, Bangkok, 2017, pp. 110-114, doi: 10.1109/IEMECON.2017.8079572.
- [4] Weiliang Liu, Xueguang Yuan, Yangan Zhang, Mengya Liu, Zhenyu Xiao, Jianlan Wu, "An End to End Method for Taxi Receipt Automatic Recognition Based on Neural Network", *Automation Control Conference (ITNEC) 2020 IEEE 4th Information Technology Networking Electronic and*, vol. 1, pp. 314-318, 2020.
- [5] S. Afroge, B. Ahmed and F. Mahmud, "Optical character recognition using back propagation neural network," *2016 2nd International Conference on Electrical, Computer and Telecommunication Engineering (ICECTE)*, Rajshahi, 2016, pp. 1-4, doi: 10.1109/ICECTE.2016.7879615.
- [6] R. Smith, "An Overview of the Tesseract OCR Engine," *Ninth International Conference on Document Analysis and Recognition (ICDAR 2007)*, Parana, 2007, pp. 629-633, doi: 10.1109/ICDAR.2007.4376991.
- [7] H. Lu, B. Guo, J. Liu and X. Yan, "A shadow removal method for tesseract text recognition," *2017 10th International Congress on Image and Signal Processing, BioMedical Engineering and Informatics (CISP-BMEI)*, Shanghai, 2017, pp. 1-5, doi: 10.1109/CISP-BMEI.2017.8368720.
- [8] T. C. Wei, U. U. Sheikh and A. A. A. Rahman, "Improved optical character recognition with deep neural network," *2018 IEEE 14th International Colloquium on Signal Processing and Its Applications (CSPA)*, Batu Feringghi, 2018, pp. 245-249, doi: 10.1109/CSPA.2018.8368720.
- [9] Opencv.org, "Information related to Open CV", 2020. [Online]. Available: <https://opencv.org/> [Accessed: March 2020].
- [10] H. Sidhwa, S. Kulshrestha, S. Malhotra and S. Virmani, "Text Extraction from Bills and Invoices," *2018 International Conference on Advances in Computing, Communication Control and Networking (ICACCCN)*, Greater Noida (UP), India, 2018, pp. 564-568, doi: 10.1109/ICACCCN.2018.8748309.
- [11] R. F. Rahmat, D. Gunawan, S. Faza, N. Haloho and E. B. Nababan, "Android-Based Text Recognition on Receipt Bill for Tax Sampling System," *2018 Third International Conference on Informatics and Computing (ICIC)*, Palembang, Indonesia, 2018, pp. 1-5, doi: 10.1109/IAC.2018.8780416.
- [12] Shafait, Faisal and Smith, Ray. (2010). "Table detection in heterogeneous documents." *ACM International Conference Proceeding Series*. 65-72. 10.1145/1815330.1815339.
- [13] Shakya, Subarna. "Virtual Restoration Of Damaged Archeological Artifacts Obtained From Expeditions Using 3D Visualization." *Journal of Innovative Image Processing (JIIP)* 1, no. 02 (2019): 102-110.
- [14] Manoharan, S. (2019)," Smart Image Processing Algorithm For Text Recognition, Information Extraction And Vocalization For The Visually Challenged", *Journal of Innovative Image Processing (JIIP)*, 1(01), 31-38.
- [15] A. Domale, B. Padalkar, R. Parekh and M. A. Joshi, "Printed Book to Audio Book Converter for Visually Impaired," *2013 Texas Instruments India Educators' Conference*, Bangalore, 2013, pp. 114-120, doi: 10.1109/TIIEC.2013.27.