# Text Extraction from Business Cards and Classification of Extracted Text Into Predefined Classes

## C Madan Kumar[a] and M Brindha[b]

[a]Assistant Professor of CSE Department, KITS,Warangal,Telangana, India.
[b]Assistant Professor, CSE Department, NIT Trichy.

Abstract: Optical Character Recognition (OCR) is the technology for identification of characters with utmost accuracy possible by employing suitable pre-processing, processing and post processing refinements. Practical application of OCR is very wide ranging from day-to-day need to scientific research purposes. One very crucial application is to Automate Digitalisation of Business cards. Since Business cards comes in different fonts and sizes, and most importantly with different lighting conditions, applying OCR can be done after careful processing. Avoiding noise from source image is one of the most crucial step in any image-processing process and it has major weightage in the accuracy of further step and thus indirectly has a huge contribution for our final outcome. Further proper noise cancellation in our source image can reduce number of future steps required to attain good accuracy and also avoid problem of iterating our sample over cycles to avoid better contrast or to distinguish text in our source lying in a noisy matrix. Digitalising business cards aims at classification of the text extracted from our source image in the hard copy of the business card directly into respected following classified fields so that it becomes a lot easier to proceed any desired function ones aims to do with that business card in this online era.

Keywords: Optical Character Recognition (OCR), Text Extraction, Biggest Contour Method, Houglines Method, Canny Edge Detection, Image Segmentation

## 1. Introduction

Digital image processing refers to use of computerized algorithms to perform image processing on digital images. Optical character recognition or OCR(Tesseract algorithm) is a form of information entry for business cards, e-mails, pan cards, Id cards, which scans a document in written form or printed form and retrieving the text out of it. The idea of Business card Digitalisation has been evolved from Automatic License Plate Recognition.

Digitizing the Business Cards is a real challenge. They come in different formats and fonts. Text extraction in different lighting conditions is very difficult. A formal structure for the business card reader which was innovated is reported. A Boundary Detection method is proposed called Biggest Contour method and Houghlines Transformation method, image extraction and segmentation technique based on a statistical method called Connected Component Method.

The task is to detect the boundary of the card by eliminating the background. The image was subdivided into an array of smaller blocks, over which gray thresholding is used to compute local thresholds. These thresholds are then stored in another array. This method of dividing the image into smaller images and then applying threshold, instead of applying threshold on the image as a whole is known as adaptive threshold.

## 2. Review of Literature

An OCR system for Business Card was introduced by (Hisashi Saiga et al., 2003) Yasuhisa Nakamura et al., where the arrangement of the bounding rectangles is twice arranged into subsets in various ways in view of the separation between them: once to choose the course of lines in the image and once to extract the lines. After that, they are incorporated into character fragments. One direction is selected which each character on the card is expected to look by matching some character fragments with the standard patterns. The system works on a "line by line" basis through the remaining acknowledgment stage. Character fragments in a line and consolidation made by incorporating neighbouring ones are coordinated against the standard patterns. Their coordinated scores are distinguished in order to choose which fragments consists of which character. The matching scores and hopeful classifications compared to these combinations are stored in a 2-dimensional array. For that associated part extraction strategy was utilized where putting away the directions of the two closures of a corner to corner of the rectangle shape encompassing it. The framework disregards little rectangles, seeing them as noise. The rectangles are the premise of character extraction in that all characters are created in later stages by incorporating them.

After that instead of applying character fragmentation better result was obtained using ROI (Tong li et al., 2017) followed character recognition which was developed by Tao et al. utilize the probability function and transform the identification into a grouping issue based on corner detection. But, the corner detector fundamentally centers around the zone with rich texture data, and hence dismisses the global portrayal of the whole image. To overcome the disadvantage that the salient points are not constantly accumulated in textured regions. The wavelet transform based method is suggested with respect to the ROI lies in the region changing essentially after change. Segmentation based algorithm was presented to address this problem. Then for linguistic processing it was presented by Hisashi Saiga et al.

## 3. Proposed Methodology

Language processing phase is done line by-line. This implies it is calling

upon each time it gets acknowledgment results for a line from the acknowledgment stage. The business cards contain the following fields. First name of the Person, Last name of the Person, Organisation, Designation, Email, Address, Office Number, Direct Dial, Fax Number, Mobile Number, Website, City, State, Country and Zip code.

The output received from Tesseract-OCR has to be matched to the respective fields for which suitable logics are performed. First create an empty dictionary. There are 15 fields. So there must be 15 keyword where each Keyword represent each field and assign the values for all keywords to None. Also create a list which contains data extracted in Tesseract-OCR separated by newline"\n".

### 3.1 Extraction of e-mail

Step 1: For extracting emails Regex is used. For that use re which is a package for Regex. So for extracting emails find if there is '@' present in the word.

Step 2: Add the email extracted to value present in dictionary where keyword is E-mail.

### 3.2 Extraction of website

Step 1: For extracting website Regex is used. For that use re which is a package for regex. So for extracting website find if the word starts with 'www.' or 'Web.'

Step 2: If Step 1 is satisfied then from email scrape the words present after '@' and add it to website.

Step 3: Add the website extracted to value present in dictionary where keyword is Website.

### 3.3 Extraction of firstname and lastname

Step 1: For extracting first name and last name Natural Language Processing is used.

Step 2: For that use difflib which is a package for matching similar strings especially known as Sequence Matcher.

Step 3: Most of the cards have email starting with the name of the person. First make the Tesseract output into lowercase for easier comparison.

Step 4: Convert the word extracted into lowercase and pass the word extracted as second argument to Sequence matcher function.

Step 5: Run the loop and compare every item in list with word extracted from email. Then extract the item which has the highest similarity ratio.

Step 6: Split the word extracted separated by space (" ").

Step 7: If Step 6 fails then add the word extracted from email present before '@' to dictionary where keyword is First name and set Last name to None.

### 3.4 Extracting organization

Step 1: When the Business Cards was analysed it was found that second part of email that is the word present after '@' and before domain names

like '.com' ,'.in' are mostly organization names.

Step 2: The name of the organization is scraped directly from email. Convert the organization to uppercase.

Step 3: Add the organization extracted to dictionary where keyword is Organisation.

### 3.5 Extraction of designation

Step 1: A new text file was created which list of designations has separated by a newline.

Step 2: The text file is read line by line and compare it with the items in Tesseract-OCR list.

Step 3: If the pattern matching fails it means the word is not present in text file. So add the designation which was previously not present to text file and run the code again.

Step 4: If Step 3 is satisfied add it to dictionary where keyword is Designation.

### 3.6 Extraction of numbers

Step 1: For extracting phone numbers use phone numbers which is a package for extracting phone numbers. It finds all phone numbers present in the Tesseract-OCR and creates a list of phone numbers.

Step 2: Next matching the phone numbers to respective fields has to be done.

Step 3: Create a list of prefixes like 'Office', 'Off', 'Main', 'Telephone'.

Step 4: Create an Empty list named matchnumbers.

### 3.7 Matching office-number

Step 1: For matching office numbers first the index of the numbers is found in the phone number list.

Step 2: From that index move ten index forward and check if prefixes like 'Office', 'Off', '(O)', 'Telephone' 'Tel', 'tel', 'T' are present.

Step 3: Remove the prefixes like 'Office', 'Off', '(O)', 'Telephone' 'Tel', 'tel', 'T' from prefix list to avoid redundancy.

### 3.8 Matching direct-dial number

Step 1: For matching direct dial numbers first find the index of the numbers in the phone number(Yamaguchi et al., 2015)list.

Step 2: From that index move ten index forward and check if prefixes like 'Direct', 'Dir', 'Direct Dial', 'Dial', 'Main' ,'D', 'toll free', 'direct', 'Phone'.

Step 3: Remove the prefixes like 'Direct', 'Dir', 'Direct Dial', 'Dial'.

### 3.9 Matching fax numbers

Step 1: For matching fax numbers first find the index of the numbers in the phone number list.

Step 2: From that index move ten index forward and check if prefixes like

'Fax', 'fax', 'FAX', 'F:', '(F)', 'Fax', 'F', 'f' are present.

**Step 3:** Remove the prefixes like 'Fax', 'fax', 'FAX', 'F:', '(F)', 'Fax', 'F', 'f' to avoid redundancy.

### 3.10 Matching mobile numbers

**Step 1:** For matching mobile numbers first find the index of the numbers in the phone number list.

**Step 2**: From that index move ten index forward and check if prefixes like 'Mobile'.

**Step 3:** Remove the prefixes like 'Mobile', 'mobile', 'Mob', 'Cell', '(M)','M', 'm', 'cell' to avoid redundancy.

**Step 4:** Now the matchnumber list contains matched numbers but with prefixes. For that iterate the matchnumber list.

**Step 5:** If anyone of the items in 'Office', 'Off', '(O)', 'Telephone' 'Tel', 'tel', 'T' is present as a substring in one of the item of matchnumber, add it to dictionary where key word is Officenumber.

**Step 6:** If anyone of the items in 'Mobile', 'mobile', 'Mob', 'Cell', '(M)','M', 'm', 'cell' is present as a substring in one of the item of matchnumber.

**Step 7:** If anyone of the items in 'Fax', 'fax', 'FAX', 'F:', '(F)', 'Fax', 'F', 'f' is present as a substring in one of the item of matchnumber.

**Step 8:** If anyone of the items in 'Direct', 'Dir', 'Direct Dial', 'Dial'. 'P' is present as a substring in one of the item of matchnumber.

**Step 9:** If it is found that matchnumber list is empty and its analysed that 'Office'.

**Step 10**:So the same Matching procedure where instead of moving ten indexes backward as well as move ten indexes forward and followed same algorithm.

### 3.11 Extraction of zipcode

**Step 1:** First install package libpostal,   a package for parsing a string. In that string first convert the Tesseract-OCR output.

**Step 2:** Remove the matchnumber list from Tesseract-OCR output.

**Step 3:** Now the string has only address. Then parse the string which have list of tuples where first item in tuple is word second item in tuple is type of the word.

**Step 4:** If found that second item in tuple as 'postcode' then extract the first item in tuple and added to dictionary where keyword is Zip code

**Step 5:** If the Step 4 fails use regex. For that use a package for regex. So an expression is written to extract those numbers and add it to dictionary where keyword is Zip code.

### 3.12 Extraction of city

**Step 1:** First install package libpostal, a package for parsing a string. In that string first convert the Tesseract-OCR output.

**Step 2:** Remove the matchnumber list from Tesseract-OCR output.

**Step 3:** Now the string has only address. Then parse the string which have list of tuples where first item in tuple is word second item in tuple is a type of the word.

**Step 4:** If found that second item in tuple as 'city' then extract the first item in tuple and added to dictionary where keyword City.

**Step 5:** If the Step-4 fails use another package zip code which is used for finding state and city of the corresponding zip code unless zip code is not None

**Step 6:** In some cards it is found that there are more than one city detected.

**Step 7:** Extract the city and added it to dictionary where keyword is City.

### 3.13 Extraction of state

**Step 1:** First install package libpostal, a package for parsing a string. In that string first convert the Tesseract-OCR output.

**Step 2:** Remove the matchnumber list from Tesseract-OCR output.

**Step 3:** Now the string has only address.

**Step 4:** If  found that second item in tuple as 'state'  extract the first item in tuple and added to dictionary where keyword State.

**Step 5:** If the Step 4 use another package zip code which is used for finding state and city.

**Step 6**: Its analysed and found the index of zip code in Tesseract-OCR string and found the city which has closest distance to zip code index.

**Step 7:** Extract the city and added it to dictionary where keyword is State.

### 3.14 Extraction of country

**Step 1:** First install package libpostal, a package for parsing a string. In that string first convert the Tesseract-OCR output.

**Step 2:** Remove the matchnumber list from Tesseract-OCR output.

**Step 3:** Now the string has only address. Then parse the string which have list of tuples where first item in tuple is word second item in tuple is type of the word.

**Step 4:** If found that second item in tuple as 'country' extract the first item in tuple and added to dictionary where keyword Country.

**Step 5:** If the Step 4 fails do pattern matching. Then create a list of tuples where first item in tuple is abbreviated form of the country and second item in tuple is full form of the country. So if the pattern matching occurs at first tuple add the second tuple of the corresponding item to dictionary where keyword is Country or else if the matching does not occur then   pattern matching is done on second item in tuple to Tesseract-OCR output.

**Step 6:** Country extracted is added to dictionary where keyword is Country.

### 3.15. Extraction of address

**Step 1:** First install package libpostal, a package for parsing a string. In that string first convert the Tesseract-OCR output.

**Step 2:** Remove the matchnumber list from Tesseract-OCR output.

**Step 3:** Now the string has only address. Then parse the string which have list of tuples where first item in tuple is word second item in tuple is type of the word.

**Step 4:** Validate that first item to be added to that address string is

'house_number'.

**Step 5:** Add the address string to dictionary where keyword is Address.

**Step 6:** If Step 5 fails first convert the Tesseract-OCR output, firstname, lastname, organization, designation all to lowercase and then remove the firstname, lastname, organization, designation from Tesseract-OCR output.

**Step 7:** Remove the matchnumber list from Tesseract-OCR output. Now the string has only address. Without parsing directly add the address string to dictionary where keyword is Address.

## 4. Results and Discussions

### 4.1 Experimental setup

The proposed system was implemented using OpenCV-Python which is a package used for Image Processing. The important packages used are cv2 for Image Processing, NLTK for Natural Language Processing, Pytesseract for converting Image to Text. When the text is read from a Normal Image, the extracted text quality is low.



**Fig. 1. Input (Drag and Drop a Card in Web Interface)**

So first, the normal image is converted (Fig. 1) to Grayscale Image (Fig. 2) which is to be used as a input for Erosion (Fig. 3) which in turn improves the quality for reading as it makes the Foreground Black and Background White so that Boundary Detection is made easier for Contours (Fig. 4) which connects regions which are of same intensity like all black together and all white together. Then after detecting the Boundary of the card Perspective Transformation is applied (Fig. 5) which makes the Image aligned to 90 degree vertically and perfectly aligns with screen to recognise character easily.



**Fig. 2. Input Card**



**Fig. 3. Grayscale Image**

If the first method fails to detect card boundary in situations like when Foreground and Background are of same color, then Houghlines is applied which uses Canny Edge Detection to detect edges.
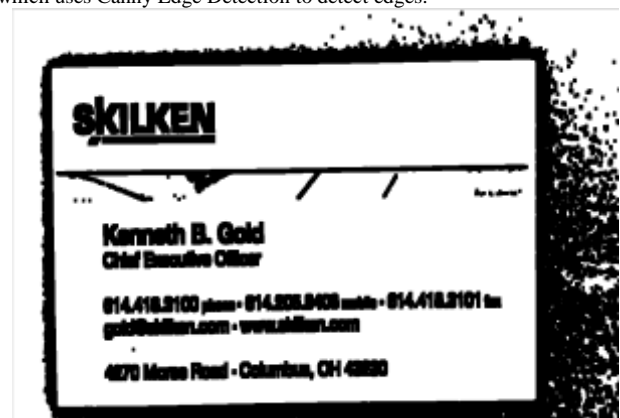


**Fig. 4. Erosion**

**Fig. 5. Contours**



**Fig. 6. Perspective Transformation**

Then the houghlines which are of threshold greater than 100 are found and followed by vertical and horizontal lines and from that the perpendicular lines which detect 4 points are found (Fig. 6).
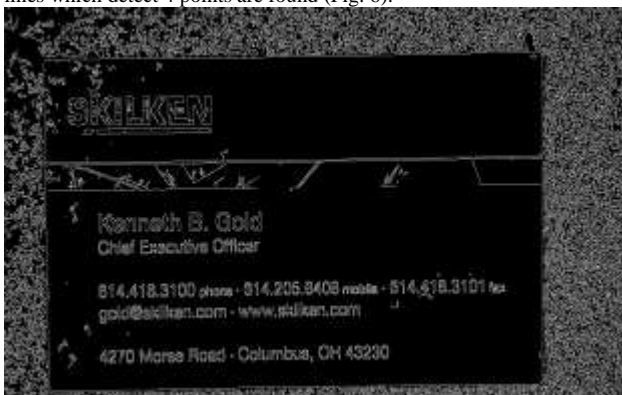


**Fig.7. Canny Edge Detection**

Then, Image processing for the boundary detected is performed (Fig. 7) by dilating it so that white pixels gets closer and regions are detected easily. Find all regions of interest (Z. L. Chen et al., 2013) and apply

morphological gradient followed by sharpening (Fig. 8) to connect boundaries of character.



**Fig. 8. Opening**



**Fig. 9. Connected Component Extraction**



**Fig.10. ROI**



**Fig. 11. Before enhancement ROI-1**



**Fig. 12. After enhancement ROI-1**

Then OCR is done followed by Text Classification. The classified text is displayed in Fig. 19. The read Business card(Wumo Pan et al.,2015) is then

stored in the database which can be queried as Output Report for future uses(Fig. 20-Fig. 21).



**Fig. 13. Before enhancement ROI-2**



**Fig. 14. After enhancement ROI-2**



**Fig. 15. Before enhancement ROI-3**



**Fig. 16. After enhancement ROI-3**



**Fig. 17. Before enhancement ROI-4**



**Fig. 18. After enhancement ROI-4**



**Fig. 19. Classified Text**





**Fig. 20.   Search previous cards in DB based on any field**

**Fig.21. Queried Output Report**

The accuracy was measured by Batch processing all the images in the given directory and it is marked as 1 if output is not NULL and 0 if NULL. Then we calculate the average for all the fields and find the accuracy.

**4.2 PERFORMANCE ANALYSIS**

**Table 1 Results**

| Fields | Accuracy (%) |
|---|---|
| Tesseract OCR | 98 |
| Firstname | 97.3 |
| Lastname | 97.3 |
| E-mail | 97 |
| Website | 97 |
| City | 96.1 |
| State | 96.2 |
| Country | 96.4 |
| Address | 96.1 |
| Phone number | 97.3 |
| Direct dial | 97.2 |
| Fax | 97.2 |
| Mobile | 97.1 |
| Designation | 96.5 |
| Organisation | 97.1 |
| Zipcode | 95.8 |

The performance of the program is given in Table 1. It can be seen that our output has better results than Camcard which recognises on line-by-line basis whereas our system works on ROI (W.Li et al., 2016) where we can easily classify the text using Python string manipulation. It is best used for Book fair, Real Estate, Job fair so that we just scan the card and store the data in a database and send them offers if the company employer find it very interesting or else for Real Estate they can sell the lands if they match the user criteria requirements. It has broad applications and Digitalising Business Card is a growing field in this world.

## 5. Summary and conclusion

In this work, a linguistic processing of text is proposed using Natural language processing technique and machine learning. This processing phase uses various new methods to classify each extracted text into its respective field. Boundary cropping and connected component method played an important role in the segmentation and extraction of text from business card. Traditional scanners use scan lines to read text, whereas in this work it is shown that the accuracy and quality of the text extracted is improved by using connected component method. Instead of reading text line by line, regions of interest is found and block processing of those ROIs is done.

In the future, the proposed system can be extended to include deep learning techniques other than simple machine learning. The accuracy of the system can be improved by making novel changes to the architecture of the linguistic processing phase.

REFERENCES

Hisashi Saiga, Yasuhisa Nakamura, Yoshihiro Kitamura, Toshiaki Morita (1993). An OCR System for Business Cards. Information Technology Research Laboratories Corporate Research and Development Group Sharp Corporation. 0-81864960-7193 $3.00 0 1993 IEEE.

Kise, Sugiyama, Yamaoka, Momota, Babaguchi and Tezuka"Mode1 Based Understanding of Document Images," Proceedings of IAPR Wodshop on Machzne Vision Application '90, pp.471-474, 1990.

Pan, W., Jin, J., Shi, G. and Wang, Q. R. (2014). A System for Automatic Chinese Business Card Recognition. Proceedings of the International Conference on Document Analysis and Recognition. USA. 577-581, 2014.

Yamaguchi, T., Nakano, Y., Maruyama, M., Miyao, H., and Hananoi, T. (2015). Digit Classification on Signboards for Telephone Number Recognition. *Proceedings of the International Conference on Document Analysis and Recognition*. UK, pp. 359-363.

Shan Du, Member, Mahmoud Ibrahim, Mohamed Shehata. (2013). Automatic License Plate Recognition (ALPR): A State-of-the-Art Review, *IEEE Transactions on Circuits and System for Video Technology*, Vol. 23, No. 2, pp. 311-324.

A. F. Mollah, S. Basu, M. Nasipuri. (2015). Segmentation of Camera Captured Business Card Images for Mobile Devices, *Int'l J. of Computer Science and Applications*, pp. 33-37.

Wumo Pan, Jianming Jin, Guangshun Shi, Q. R. Wang. (2015). A System for Automatic Chinese Business Card Recognition, *ICDAR*, pp. 577-581.

Tong Li, Junping Zhang, Xiaochen Lu and Ye Zhang, Member. (2017). SDBD: A Hierarchical Region-of-Interest Detection Approach in Large-Scale Remote Sensing Image, *IEEE GeoScience and Remote Sensing Letters*, Vol. 14, No. 5, pp. 699-703.

W. Li, P. Dong, B. Xiao, and L. Zhou. (2016). Object recognition based on the Region of Interest and optimal Bag of Words model. *Neurocomputing*, Vol. 172, pp. 271–280.

Z. L. Chen, B. J. Zou, J. F. Li, H. L. Shen, and Y. Mao. (2013). Hybrid ROI
extraction with bottom-up and top-down visual attention. J. Inf. Comput. Sci.,
Vol. 8, No. 15, pp. 3481–3488.