# Optical Character Recognition Engine to extract Food-items and Prices from Grocery Receipt Images via Templating and Dictionary-Traversal Technique

Ali Sohani[1]     Rafi Ullah[2]     Faraz Ali[3]     Athaul Rai[4]     Richard Messier[5]

## Abstract

This paper proposes a mix of some old and few novel techniques to nail down the fundamental problem of Food-Items and Prices recognition and eventual extraction of them from the Grocery Receipts. Considering in our research we didn't find any existing OCR engine that is up to that standard let alone specialized for this specific purpose. Since the target was to create a specialized OCR system, we began with an idea of creating the wrappers around basic OCR system to empower it with context of Grocery Receipt. For this, we've built pre-function and post-function wrappers over existing system called Tesseract-OCR. Our system follows specific work-flow to enhance basic OCR output. First it runs the provided image to image filters to make it most suitable for Section-level extraction. Our system then bifurcates the image into sections (like Price, Item-Names, Quantity are dealt separately from one another) according to given template layouts. Specific portion of images (sections) are then forwarded to Tesseract engine for basic OCR. Then text-extracted is forwarded to a contextual pattern matcher, to make sense of the text-extracted in a contextual manner. After testing system on particular grocery stores receipts, we successfully conclude that our techniques significantly improve on both the accuracy of overall context based text recognition and close-match detection when compared to an unassisted/ vanilla Tesseract OCR. Proposed system will empower Food-Kitchen Assistance Mobile Apps in the market.

**Keywords:** Accurate image to text converter, Receipt parsing using template matching, OCR using receipts template, Text retrieval from receipts images

## 1       Introduction

The objective of our work was grocery receipt's parsing i.e image to text conversion using open source Tesserect OCR [23][4][11]. As Tesserect OCR just retrieve text from images. We have proposed techniques for parsing receipts that is template matching. We stored templates or structures of the templates of different stores. We can easily retrieve items, quantity of each item and price of each item. We used the relative positions of items, quantities and prices in the receipts to find the item etc in new image of the same store. Thus OCR read only that required portion. That reduced the time and improved accuracy.

Before applying simple OCR technique, we first did some image pre-processing techniques. First of all image processing techniques; include image background removal. Background is actually non-textual area of image. Then we applied text deskewing [25] followed by image binarization [1][8]. We also applied resizing technique if image was smaller in size.

[1 2 3 4 5] *Data Science Department, Cubix, Pakistan*

[1]*ali.sohani@cubix.co | [2]rafiullah.khan@cubixlabs.com | [3]faraz.ali@cubixlabs.com | [4]athaul.rai@cubixlabs.com* [5]*richard@cubix.com*

Items, quantity and price portions of image were calculated from the template then Tesserect OCR was applied to each portion in order to retrieve the text. Then context sensitive spelling was applied on result.

Rest of the paper includes Related work, Tesserect OCR open source API, Image Pre-processing techniques, Methodology, Image template, Context sensitive spelling correction, Results, conclusion and future work.

## 2 Related Work

In references [3] and [4] OCRing is done using pattern matching and advance heuristics. These methods are proven to be very successful for generic type of receipts parsing. They have used Regular Expressions to extract different texts. Heuristics in [4] are very helpful to discard garbage data, that are not necessary. Here OCRing is based on assumptions noticed in large number of receipts. These techniques are generic and work on every receipt.

[19] describes the ABBY cloud SDK for receipt recognition. As receipts are not always clear. These images may be noisy due to taken by movable mobile devices. So simple scanning may not give you an accurate results. This API is paid API.

[20] is an R&D about similar purpose. This R&D is basically for receipt parsing. This R&D also inclue similar steps like image binarization, text finding etc.

OCRDroid framework has been proposed in [8]. This uses image processing techniques like deskewing, binarization etc for better results. There is limitation on multiple images OCRing and Text detection from complex backgrounds.

Beside this a lot of work has been done and going on using OCR techniques, OCR improvement using image processing techniques, using advance state of the art techniques like Neural Networks etc.

## 3 Tesseract-Ocr

Tesseract is an open source Optical Character Recognition (OCR) Engine or API, available under the Apache 2.0 license. It can be used directly use or using an API to extract typed text, handwritten text or printed text from images of different formats. It supports a wide variety of languages (we have used python) and almost for all operating systems (have used Ubuntu 16.01) [23] [21].

For configuring pytesserect in Ubuntu, use the following commands:
sudo pip install pytesserect
sudo get-apt install tesserect-ocr
After configuring it, you can select language, configuration according to your need. We have used 'eng' English as a language, "- psm 6" as a config parameter and Image object as a parameter.

## 4      Image Processing

Tesserect OCR is open source library sponsored by Google, It has accuracy issue. It is generic image to text converter. To clear the image in order to be read by OCR accurately same image processing steps have been applied as given in [3] and [4].

### A      *Image Background Removal*

As OCR process is little bit slower and there is accuracy issue in case of noisy background. Accuracy and speed issues have been improved. This will work only if you have image like as given below. When we applied OCR on below JPEG image having dimensions 3936 x 5248, it took 3.231 seconds. When we applied background removal, it took 1.725 seconds. And the result was also improved as mentioned in the table below:



**Figure 1:  Walmart receipts before and after image background removal**

## Table 1: Result Before And After Image Background Removal

| | |
|---|---|
| , K" | \" |
| a mar m | a mar 9.4.. |
| Save money. lee better. | Save money. me better. |
| ( 504 ) 522 ~ 4142 | ( 504 ) 522 - 4142 |
| MANAGER TODD JABBIA | MANAGER TODD JABBIA |
| 1901 TCHOUPITOULAS ST | 1901 TCHOUPITOULAS ST |
| NEW ORLEANS LA 70150 | NEW ORLEANS LA 70130 |
| 5T2 5022 0P# 00005251 TB# 65 TRfi 09552 | ST#  5022  OP#  00005251  TE#  83  TR# |
| 15X12 PAS WC 084705715066 7.97 X | 09552 |
| HAND TOWEL 066572107026 2.97 X | 15x12 PAS wc 084703715088 7.97 x |
| GATORADE 005200055582 F 2.00 X | HAND TOWEL 066572107028 2.97 x |
| GATORRDE 005200055652 F 2.00 X | GATORADE 005200033582 F 2.00 x |
| OXICLEAN VSR 075705751525 7.52 X | GATORADE 005200033832 F 2.00 x |
| MTG CAR FL5G 009474657442 9.97 X | OXICLEAN VSR 075703751523 7.52 x |
| T~SHIRT 088529968450 16.88 X | MTC CAR FLAG 009474657442 9.97 x |
| PUSH PINS 002775501514 1.24 X | T-SHIRT 088329968450 16.88 x |
| ULTRATECH 076526451591 5.97 X | PUSH PINS 002775501314 1.24 x |
| REESE MINI 005400044660 F 2.66 R | ULTRATECH 076326431391 5.97 x |
| 16 02 CUP 064541654511 0.67 X | REESE MINI 003400044860 F 2.88 R |
| COPY PAPER 005650010265 4.22 X | 16 oz CUP 064541654511 0.87 x |
| SRAGRAHS LQ 006700070070 25.47 T | COPY PAPER 003650010285 4.22 x |
| SUBTOTAL 87.96 | SEAGRAMS LQ 008700070070 23.47 T |
| TAX 1 9.000 5 7.66 | SUBTOTAL 87.96 |
| TAX 2 4.500 % 0.15 | TAX 1 9.000 % 7.66 |
| TOTAL 95.75 | TAX 2 4.500 % 0.13 |
| AHEX TEND 95.75 | TOTAL 95.75 |
| ACCOUNT # **** **** ***2 289 8 | AMEX TEND 95.75 |
| APPROVAL # 925741 | ACCOUNT # **** **** ***2 289 3 |
| REF # 420200465440 | APPROVAL # 923741 |
| Beg Bai Tran Amt End 831 | REF # 420200485440 |
| CREDIT 254.54 95.75 156.79 | Beg Bal Tran Amt End Bal |
| TERMINAL # 26005954 | CREDIT 234.54 95.75 138.79 |
| 07/21/14 10:45:26 | TERMINAL # 26003934 |
| CHANGE DUE 0.00 | 07/21/14 10:45:28 |
| . # ITEMS SOLD 13 | CHANGE DUE 0.00 |
| TC# 7757 9454 7517 6470 6445 | # ITEMS SOLD 13 |
| our gnu1rnntcmnd Lew: Pricngs | TC# 7737 9454 7317 6470 8445 |
| Are Unbeatable with Ad Hatch! | Our Guaranteed Low Prices |
| 07/21/14 10:45:24 | Are Unbeatable with Ad Match! |
| '*'UUSTOHHR COPY'*' | 07/21/14 10:45:28 |
| | ***CUSTOMER COPY*** |

## B      Image Binarization

Image binarization is process of converting colored image to black and white image [1] [8]. This is used to clean dirty images i.e images having noisy backgrounds [13]. Tesseract OCR by default use Otsu's Binarization [23]. But we have used this is an extra layer to make results more accurate. And the fact is that we have used images taken by mobile camera which has great chance to be noisy.

**Figure 2: Trader's joe receipt before and after image binarization**

**Table 2: Result Before And After Image Binarization**

| | |
|---|---|
| gunman 401-: '3 | 1111311053 JOE'S |
| 'Mmmmmm | I 44 East Ontario Street |
| Chicago IL 60611 | Chicago IL 50511 |
| Store #696 — (312) 951-6369 | Store #596 ' (312) 951-5359 |
| OPEN 8:00AM TO 10:00PM DAILY | OPEN 8:00AM 10 10:00PM DAILY |
| OLIVE OIL POTATO CHIPS.. 1.99 | OLIVE OIL POTATO CHIPS.. 1.99 |
| HUMMUS GARLIC ROASTED EC 1.99 | HUMMUS GARLIC ROASTED EC 1.99 |
| OHEDDAR NEH ZEALAND SHARP 3.71 | CHEDDAR NEW ZEALAND SHARP 3.71 |
| PITA NHOLE NHEAT 5" 1.69 | PITA WHOLE WHEAT 5" 1.89 |
| OLIVES MANZANILLA 2.29 | OLIVES MANZANILLA 2.29 |
| CREAMY SALTED PEANUT BUTTER 2.49 | CREAMY SALTED PEANUT BUTTER 2.49 |
| SUBTOTAL $14.16 | SUBTOTAL $14.15 |
| STATE TAX 1 $0.32 | STATE TAX 1 $0.32 |
| ITAL m1m | TOTAL $14.48 |
| ITEMS 6 v, Karl | ITEMS 6 v, Karl |
| 05-31-2015 03:11PM 0696 06 1173 0559 | 05'31-2015 03:11PM 0695 06 1173 0559 |
| THANK YOU FOR SHOPPING AT | THANK YOU FOR SHOPPING A1 |
| TRADER JOE'S | TRADER JOE'S |
| www.t,raIJgr'Oe§Im,,, | www.trader'oes.com |

## C    Image and Text Deskewing

Sometimes text in receipts are skewed may be in any direction. In that case, OCR doesn't provide correct results or sometimes doesn't. To avoid such situation we used an additional filter for text deskewing. We have used technique mentioned in [25] for text deskewing. Results have been improved. Following figure (on left) shows the skewed text found in receipt, we deskewed it first and then applied OCR. Accuracy has been improved. Comparison has been shown in table given below.
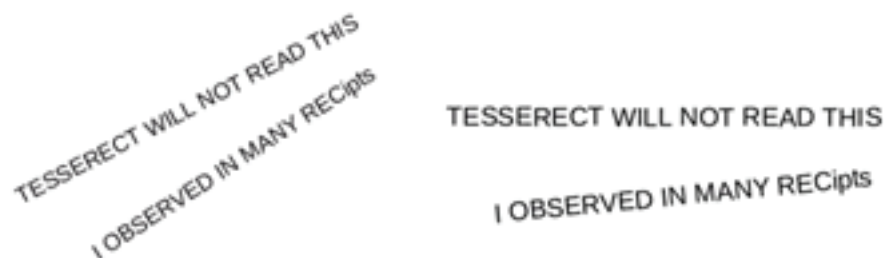
**Figure 3: Skewed text in image**

**Table 3: Results Before And After Text Deskewing**

| | |
|---|---|
| <ee~e~eeec< «£va View "As<br>6<br>\0%9<N$0 WM" «we | TESSERECT WILL NOT READ THIS<br>ECipIS<br>lOBSERVED IN MANY R |

## D      Image Resizing

This technique has been used to speed up the processing. As OCR is a game of playing on pixels so, higher the resolution of image, more will be the processing time and vice versa. So we reduced the size but not too much to effect the OCR accuracy. For example High Definition image of 4000 x 6000 will have the same result as 2000 x 4000 dimension image and processing of image discussed later will be faster than the image discussed earlier. By reducing size, the OCR performed very poor because of information loss in image. Image having DPI (Dots per Inch) greater than 300 has been observed to have good results.

## E      Image Stitching

For long receipts you have two options either you will take photo from far distance or you will take multiple snapshots. For earlier case the quality of image can be disturbed and in the later case, many images should be stitched together first and then OCRed. For this we used image stitching algorithms discussed in [9] [10] [11] and the result was fine. You have also an option to skip image stitching. OCRed multiple snapshots and then clean the result. But this can be very difficult to clean data.



**Figure 4:  First part of the image (Upper portion)**

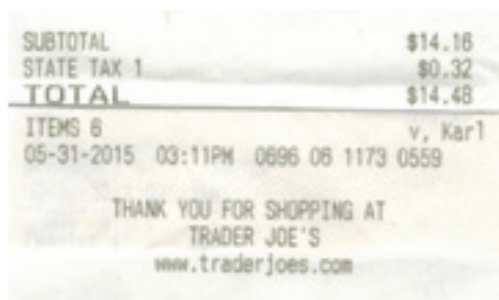**Figure 5: Second part of the receipt (middle portion)**



**Figure 6: last part of the receipt (Lower portion)**

There is no limitation on images. But you must tell us about the base image. So that other images are simply stitched to the base image.

After stitching the result is given below



**Figure 7: Image after stitching all three parts**

Now this image is simply used as an input to Tesseract OCR and text is retrieved from image. In case of high resolution, this process is very slow. To solve this problem, we first resized all the images to low resolution then stitched them together. But resizing to low resolution will affect the OCR accuracy so, that should be done carefully. We resized then to low resolution first, then stitched and then again resized to high resolution.

## 5      Image Template

We have used the stored templates of stores in database.
And while testing the image, we retrieved that specific store template. Store template have the (x, y) coordinate points, width and height information of

Footer: Xf, Yf, Wf, Hf
Item: Xi, Yi, Wi, Hi
Quantity: Xq, Yq, Wq, Hq
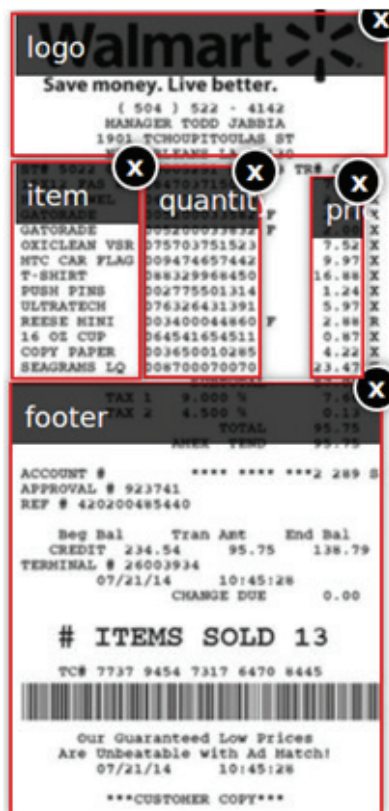Price: Xp, Yp, Wp, Hp
Logo: Xl, Yl, Wl, Hl



**Figure 8: Different regions of image (Image Template)**

We have templates of all the receipts in our database as shown in above Figure 8.

It is clear from figure the item portion is between logo and footer. So we used this is heuristic for other images to be tested. We retrieved items portions mathematically.

item x point = item x point
item y point = item y point
item width = item width
item height =  footer y point (because item is up to the start of footer)

## 6      Methodology

Complete methodology is given in the below flow diagram. In case of multiple images, image stitching is applied first then background is removed. Background contains all the non text area of image. Template may be smaller than or larger than the image that is processed.  We had the location of items, prices and quantity in the template. We used this knowledge to retrieve image's specific portion. For example we had the image portion having (x, y, w, h) in template.

Where "x", "y" is position in image and "w" is the width of image and "h" is the height of image. We found the percentage of items portion in template image.

Following things are known from template,
- Template Picture Width  = $W_t$
- Template Picture Height = $H_t$
- Template Picture "x" = $X_t$
- Template Picture "y" = $Y_t$
- Template Picture items portion width = $W_i$
- Template Picture items portion height = $H_i$

Percentage of items portion width = $PW_i$
Percentage of items portion height = $PH_i$
Percentage of items portion X = $PX_i$
Percentage of items portion Y = $PY_i$

$$PW_i = (W_i / W_t) * 100 \qquad (1)$$
$$PH_i = (H_i / H_t) * 100 \qquad (2)$$
$$PX_i = (X_t / W_i) * 100 \qquad (3)$$
$$PY_i = (Y_t / H_i) * 100 \qquad (4)$$

We have found that what percent of the image is items in the template. Used relative calculation to find items portion in new image.

Set the current picture width and height with respect to template image percentage of width and height
Current  image is the image being processed
Current image width  = $W_c$
Current image Height  = $H_c$
Item portion width in new image = IPW
Item portion height in new image = IPH

Item portion X in new image = IPX
Item portion Y in new image = IPY
These values can be calculated by

$$IPW = ( W_c / 100 ) * PW_i \qquad (5)$$
$$IPH = ( H_c / 100 ) * PH_i \qquad (6)$$
$$IPX = ( IP_W / 100 ) * Px_i \qquad (7)$$
$$IPY = ( IPH / 100 ) * Py_i \qquad (8)$$

Using above technique we retrieved the image portions/sections using template information independent of requested image size, whether greater or lesser than template image. This gave us better results.
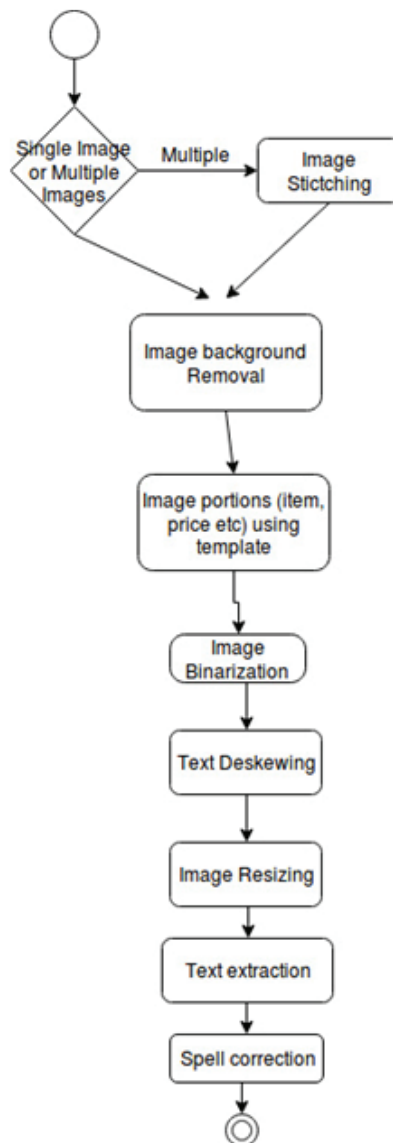


**Figure 9:  Flow chart of the proposed technique**

The next challenge was to read the item and quantity portion if the items in tested image is lesser or greater than items in template image. For example in template image we had 10 items and in tested image we had 2 items. Now portion covered by 10 items will not be the same as portion in image covered by 2 items. If OCR read same portion, result will not be acceptable. To tackle this issue, we used some heuristics, It is clear from the receipts (Figure 7 and Figure 8) that portion between logo and footer is always items portion in receipts.

## 7    Context Sensitive Spell Correction

To make Tesserect OCR results more accurate, we used context sensitive spelling correction. Context sensitive spelling correction is a technique of correcting OCR results by matching them with dictionaries of stores. For example we scanned wallmart Receipt, OCR returned result, we matched the result with wallmart dictionary i.e words/terms used in that store. Each stores name terms/words as well as matching score is stored. For the next read, result was automatically corrected. After scanning a lot of receipts, at last this method perfectly worked. Beside this we have used two levels of dictionaries for spell correction, store specific and grocery related dictionary.

Beside this we have used a corpus of text that are part of receipt but not our required text. Words such as tax, total, subtotal, discount etc are included in that corpus. These words are excluded at the very first stage from OCR result.

## 8    Results

We have tested our algorithm on various receipts template of various sizes and resolution. We observed it is performing best in many cases. We have stored template of wallmart having 397 width and 804 height. Tested using wallmart image having width 2043, height 4128 and noisy background.



**Figure 10:  Test Image (Walmart Receipt)**

**Table 4: Result of Walmart Receipt In Figure 10**

| Items Name | Prices |
|---|---|
| 15x12 fas wc | 7.97 |
| hand towel | 2.97 |
| gatorade | 2.00 |
| gatorade | 2.00 |
| oxiclean vsr | 7.52 |
| mtc car flag | 9.97 |
| t-shirt | 16.88 |
| push pins | 1.24 |
| ultratech | 5.97 |
| reese mini | 2.88 |
| 16 oz cup | 0.87 |
| copy paper | 4.22 |
| seagrams lq | 23.47 |

Given below is the receipt of Trader's joe



**Figure 11:  Test Image (Trader's Joes Receipt)**

**Table 5: Result Of Trader's Joe Receipt In Fig. 11**

| Items Name | Price |
|---|---|
| olive oil potato chips | 1.9 |
| hummus garlic roasted ec | 1.9 |
| cheddar new zealand sharp | 3.71 |
| pita whole wheat 5\ | 1.8 |
| olives manzanilla | 2.2 |
| creamy salted peanut butter | 2.4 |

For both of the above receipts there is no quantity portion, so quantities have been returned empty. If there were quantities in template images, they will be returned.

We have observed that, OCR accuracy has been improved up to great extent using image processing techniques, that was pre-processing step in our proposed technique and then the OCR result has been adjusted using our post-processing technique that is spell correction of OCRed Text.

## 9    Conclusion

A In this paper we worked on template based matching of receipt and retrieving text from receipt images. For the new image (receipt) first we retrieved the structure of receipt and then calculated the different portions of images i.e items portion, price portion and then parsed the new receipt accordingly. Retrieved text is then cleaned by filtering (context sensitive spelling correction). We have tested system for 10 different stores receipts and our proposed template based matching and parsing gave good results. It also worked best in case of noisy images.

## 10    Future Work

We have shown this idea seems good in case of noisy and complex receipts. Future work is to do generic receipt parsing and making template based matching efficient. Because receipts may vary in content from time to time. On some occasions receipts may have discount portion while in normal situation it may be simple containing only items, prices and quantities. Generic receipt parsing works without template. Our proposed system cannot parse receipts whose templates are not available in our database.  We will also do receipt recognition using Machine Learning techniques and then parse it.

## References

[1] Chaki, Nabendu, Soharab Hossain Shaikh, and Khalid Saeed. "A comprehensive survey on image binarization techniques." In Exploring Image Binarization Techniques, pp. 5-15. Springer India, 2014.

[2] Troller, Milan. "Practical OCR system based on state of art neural networks." (2017).

[3] Rafi, Ali, Faraz, Athaul "OCR Engine to extract Food-items and Prices from Receipt images via pattern matching and heuristics approach" In International Conference of Coputing and Related Technologies, December 2017, SMIU, Karachi, Pakistan

[4] Rafi, Ali, Faraz, Athaul "OCR Engine to Extract Food-items, Prices, Quantity, Units from Receipt Images, Heuristics Rules Based Approach" in IJSER Volume 9, Issue 2, February 2018 (accepted)

[5] Stadermann, Jan, Denis Jager, and Uri Zernik. "Hierarchical Information Extraction Using Document Segmentation and Optical Character Recognition Correction." U.S. Patent Application 15/620,733, filed September 28, 2017.

[6] Modi, Hiral, and M. C. Parikh. "A review on optical character recognition techniques."Int J Comput Appl 160, no. 6 (2017): 20-24.

[7] Oudah, Nabeel, Maher Faik Esmaile, and Estabraq Abdulredaa. "Optical Character Recognition Using Active Contour Segmentation."Journal of Engineering 24, no. 1 (2018): 146-158.

[8] Zhang, Mi, Anand Joshi, Ritesh Kadmawala, Karthik Dantu, Sameera Poduri, and Gaurav S. Sukhatme. "OCRdroid: A Framework to Digitize Text Using Mobile Phones." In MobiCASE, pp. 273-292. 2009.

[9] Kumar, Asit, and Sumit Gupta. "Detection and recognition of text from image using contrast and edge enhanced mser segmentation and ocr."IJOSCIENCE (INTERNATIONAL JOURNAL ONLINE OF SCIENCE) Impact Factor 3, no. 3 (2017): 3.

[10] Farahmand, Atena, Hossein Sarrafzadeh, and Jamshid Shanbehzadeh. "Noise removal and binarization of scanned document images using clustering of features." (2017).

[11] Wang, Fu-Bin, Paul Tu, Chen Wu, Lei Chen, and Ding Feng. "Multi-image mosaic with SIFT and vision measurement for microscale structures processed by femtosecond laser."Optics and Lasers in Engineering 100 (2018): 124-130.

[12] Zhang, Jing, Guangxue Chen, and Zhaoyang Jia. "An image stitching algorithm based on histogram matching and SIFT algorithm."International Journal of Pattern Recognition and Artificial Intelligence 31, no. 04 (2017): 1754006.

[13] ZHAO, Yan, Yue CHEN, and Shi-gang WANG. "Corrected fast SIFT image stitching method by combining projection error."Optics and Precision Engineering 6 (2017): 029.

[14]   Sharma, Manoj, Anupama Ray, Santanu Chaudhury, and Brejesh Lall. "A Noise-Resilient Super-Resolution framework to boost OCR performance." In Document Analysis and Recognition (ICDAR), 2017 14th IAPR International Conference on, vol. 1, pp. 466-471. IEEE, 2017.

[15]   Brisinello, Matteo, Ratko Grbić, Matija Pul, and Tihomir Anđelić. "Improving Optical Character Recognition Performance for Low Quality Images." In 59th International Symposium ELMAR-2017. 2017.

[16]   Patel, Amit, Burra Sukumar, and Chakravarthy Bhagvati. "SVM with Inverse Fringe as Feature for Improving Accuracy of Telugu OCR Systems." In Progress in Intelligent Computing Techniques: Theory, Practice, and Applications, pp. 253-263. Springer, Singapore, 2018.

[17]   GOCR - A Free Optical Character Recognition Program. http://jocr.sourceforge.net/.

[18]   OCR resources (OCRopus). http://sites.google.com/site/ocropus/ocr-resources

[19]   OCRAD - The GNU OCR. http://www.gnu.org/software/ocrad/.

[20]   Simple OCR - Optical Character Recognition. http://www.simpleocr.com/.

[21]   https://ocrsdk.com/documentation/quick-start/receipt-recognition/

[22]   http://rnd.azoft.com/applying-ocr-technology-receipt-recognition/

[23]   Tesseract OCR Engine. http://code.google.com/p/tesseract-ocr/

[24]   http://opencv-python-tutorials  last visited 10-Oct-2017

[25]   https://github.com/tesseract-ocr last visited 6-Oct-2017

[26]   http://pyimagesearch.com last visited 9-Oct-2017

[27]   All images are from http://google.com