

**ICWSTCSC: 2018**

TEQIP-III Sponsored  
2<sup>nd</sup> International Conference

**Women in Science & Technology:  
Creating Sustainable Career**

June 28-30, 2018

# **248: A Novel Approach of Tesseract-OCR Usage for Newspaper Article Images**

**Presentation by:**

Chaitanya Tejaswi

Dr. Bhargav Goradiya

Prof. Ripal Patel

ICWSTCSC - 2018



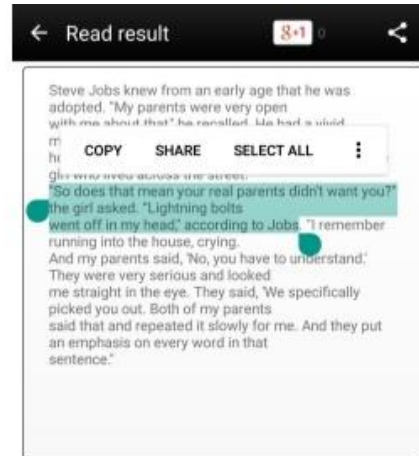
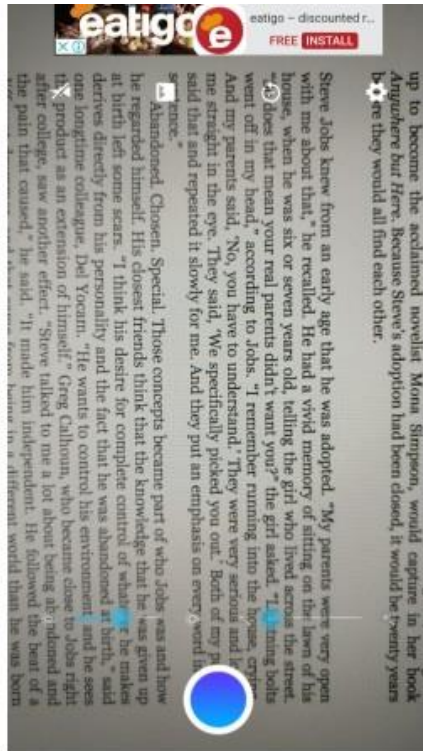
# Outline of the Presentation

- Introduction
  - Scope
  - OCR: Steps
- Objective
- Tests
- Proposed Scheme
- Data Sets
- Results
- Example
- Limitations
- Conclusion & Future Scope



# Introduction

- Our original goal was to create an Android app that gives a text to speech output for text extracted from printed newspaper articles, using open-source software tools.
- The intent was to achieve this for English language & extend it to Indian languages – (Hindi, Gujarati, ...).





# Scope

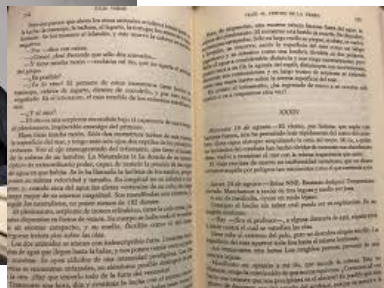
- *People with visual impairment*



- *Journalists, who intend to deal with specific stories*



- *Daily commuters, who wish to consume daily news in an audio form*

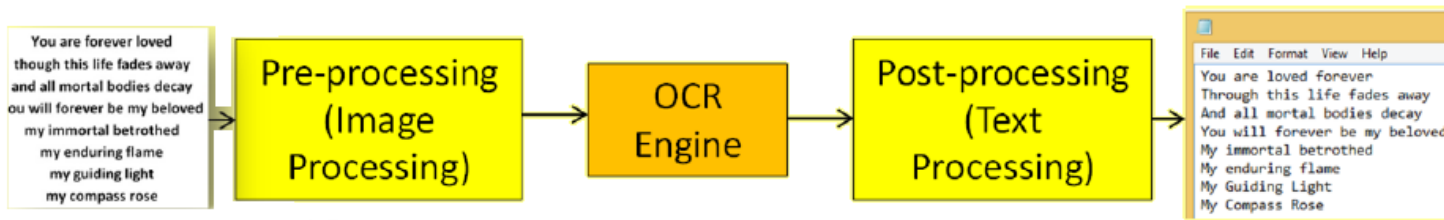


- *Students, who wish to compile notes from a variety of sources*



# OCR: Steps

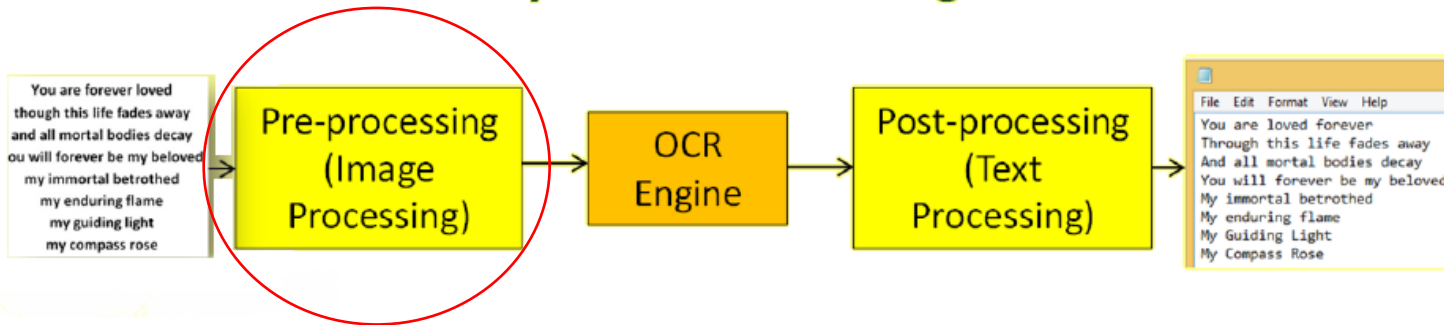
## OCR System: Process Diagram



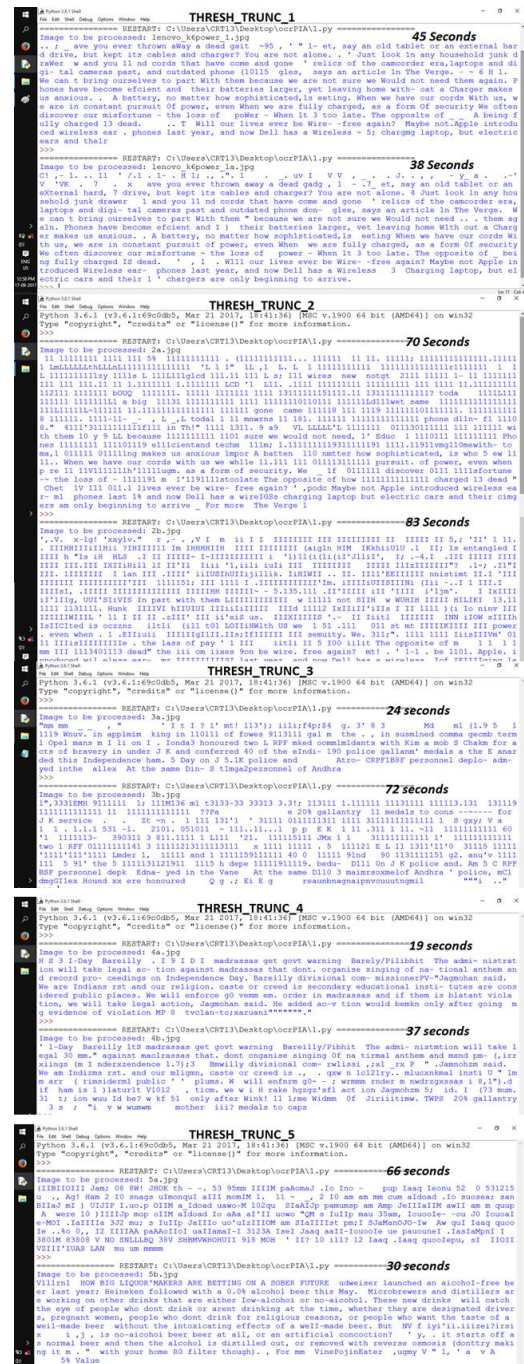
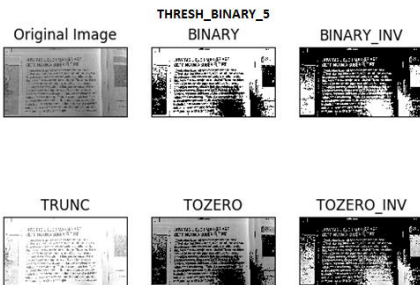
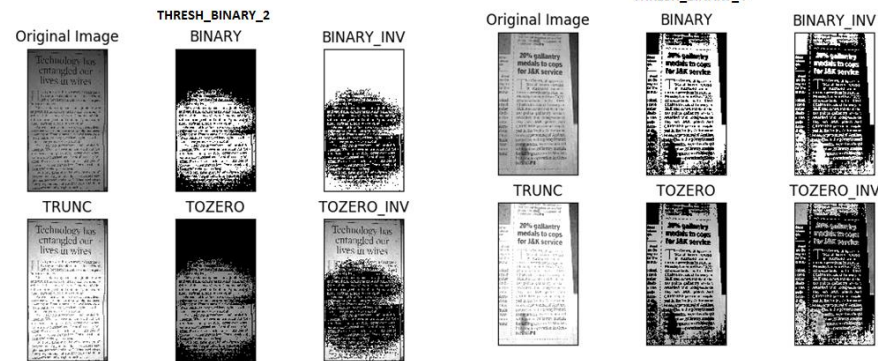
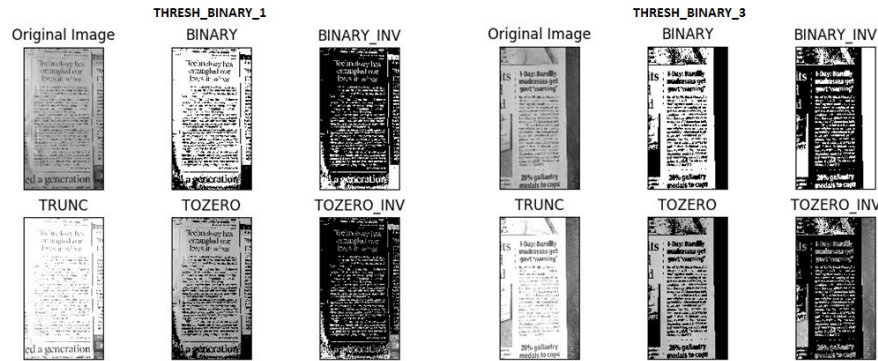
# Objective

- Over time, we shifted our attention to the *Pre-Processing* stage.
- Optimising OCR output from OCR Engine (Tesseract-OCR), by using suitable Pre-processing techniques.

## OCR System: Process Diagram







# Tests

- Determining the optimal OpenCV binarization scheme

- **THRESH\_BINARY**

$$\text{dst}(x, y) = \begin{cases} \text{maxval} & \text{if } \text{src}(x, y) > \text{thresh} \\ 0 & \text{otherwise} \end{cases}$$

- **THRESH\_BINARY\_INV**

$$\text{dst}(x, y) = \begin{cases} 0 & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{maxval} & \text{otherwise} \end{cases}$$

- **THRESH\_TRUNC**

$$\text{dst}(x, y) = \begin{cases} \text{threshold} & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{src}(x, y) & \text{otherwise} \end{cases}$$

- **THRESH\_TOZERO**

$$\text{dst}(x, y) = \begin{cases} \text{src}(x, y) & \text{if } \text{src}(x, y) > \text{thresh} \\ 0 & \text{otherwise} \end{cases}$$

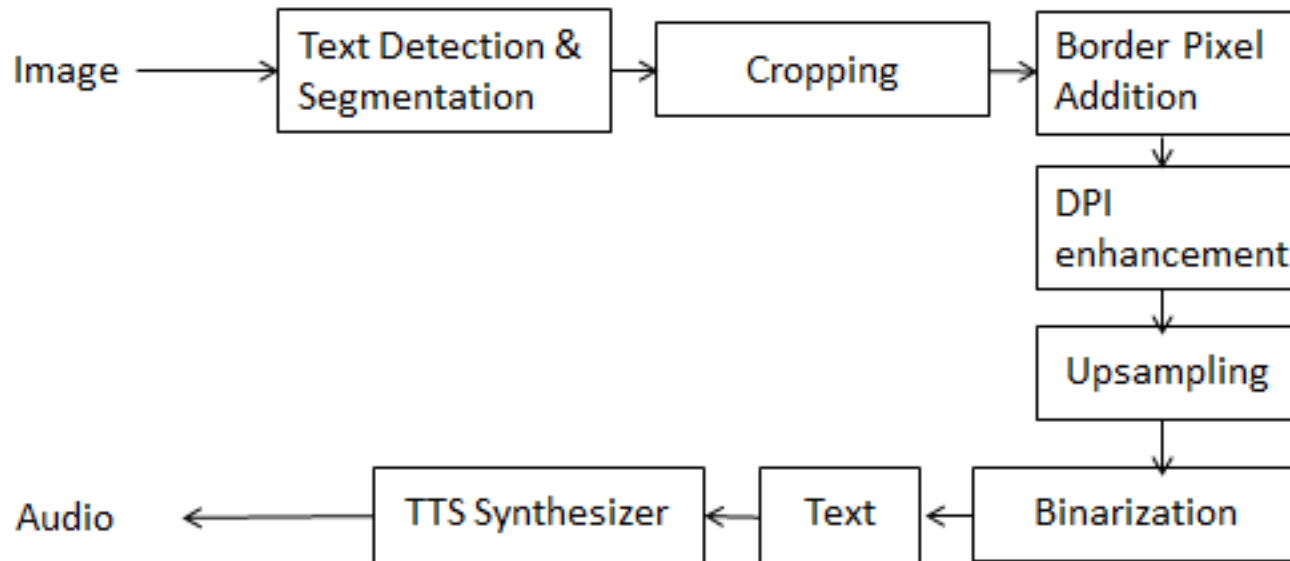
- **THRESH\_TOZERO\_INV**

$$\text{dst}(x, y) = \begin{cases} 0 & \text{if } \text{src}(x, y) > \text{thresh} \\ \text{src}(x, y) & \text{otherwise} \end{cases}$$





# Proposed Scheme



- The method proposed by Lukás Neumann & Jiri Matas [\[16\]](#) is used for Scene Text Detection. Implemented as in-built function (OpenCV-contrib).
- Cropping is done by implementing Maximal Rectangle Area algorithm in Python.
- Upsampling is done using Gaussian Pyramid Upsampling.
- Binarization is done using Otsu's Thresholding.
- Text Synthesis is done using gTTS command line tool.



1<sup>st</sup>

20% gallery  
models to cons

2<sup>nd</sup>

In this context, Union road transport minister Nitin Gadkari's statement that India will not allow driverless cars fearing massive loss of jobs is a giveaway of our confidence levels in the face of disruptive change. The spectre of rampant unemployment is a distinct possibility and it makes eminent sense to tame the divisive forces and violent proclivities that are in danger of cutting loose. Pluralism and education are the keys to securing India's future, and Mukherjee did well to lay down the options for both government and citizens in his expansive fare well speech.

reover, several citizens of Gujarat are suffering due to poor road conditions and damage to bridges and flyovers owing to the extreme weather conditions. Therefore, on humanitarian grounds, we're postponing our strike to August 5, in case the GST council decision does not meet their demands," said Arvind

Soon afterwards, the buggy came out from the northern gate with President Kovind and his predecessor Pranab Mukherjee.

Modi with nearly 30 million followers on his personal handle is the second most followed world leader on Twitter.

The federal agency also plans to resume premium processing of other H-1B petitions.

Crucially all of them will get a one-time exemption with AFC making it clear that they have to adhere to the criteria from next season.

one Zahid as one of the four men who had sodomised the boys on June 6. Meanwhile, Powai police have ex-

A K Gupta, general manager, Western Railway, inspected flood-affected tracks on Viranganam per day.



# Results

	No. of Characters	Without Preprocessing		With Preprocessing		
		Precision	Recall	Precision	Recall	
→ Set 1	3809	14.33	22.47	56.27	48.24	Worst Case
Set 2A	3620	36.2	43.33	94.4	82.13	Best Case
Set 2B	3513	44.28	48.27	69.86	78.37	
→ Set 2C	3508	13.33	37.81	69.86	74.25	
Set 2D	3318	15.75	32.47	64.24	68.92	

Three repetitions were made for each image, and OCR outputs were evaluated against Ground Truth.

## Observations:

- In extreme cases, Precision was higher than Recall value. Vice-versa in general.
- Precision is usually lower than Recall values. This is due to garbage character output.

\*

$$\text{precision} = \frac{\text{number of correct items}}{\text{number of items in OCR output}}$$

$$\text{recall} = \frac{\text{number of correct items}}{\text{number of items in ground truth}}$$





# Example 1



FOLDERS

02

CT01\_results\_original.txt

```
1  ===== ePaper =====
2  ===== Photos =====
3  ===== 1 =====
4  [A]
5  P01 ' t ' Of Empathy
6  > In his farewell I speech, President M ukherjee made
7  an eloquent case for pluralism and education
8  11 his considered farewell speech President Pranab Mukherjee,
9  li'ightly praised by finance minister Arun J aitley for the utterly '
10 non-partisan manner in which he has conducted himself
11 in Rashtrapati Bhavan, brought up the important issue of the
12 relationship between violence and contemporary public discourse,
13 where he invoked the spirit of Mahatma Gandhi to argue that
14 when such discourse becomes tinged with intolerance, it makes
15 the atmosphere conducive to greater violence as well as marginaliv
16 sation of less privileged groups. In today's atmosphere of polarised
17 debate the departing president's fervent plea - "we must free our
18 public discourse from all forms of violence, physical as well as
19 verbal" -needs to be heard,
20 Pluralism has been a hallmark of Indian civilisation, and it's
21 something to which we must hold fast.
22 ,, y f, Thebattleof ideasneednotgetpersonal
23 " 7 ' ' -- there are other, superior ways of
24 y winning it than to suggest that one's
25 , , opponent is anti-nationali Debate has
26 iii, , to be civil and rational, not hysterical
27 "1'1" . '1 _ andparanoid;thosewho suggesttheir
28 " " , "3'," political opponents should be exiled
29 " 'T " to Pakistan are in fact making India a
30 ' ' ' bit like Pakistan, which is unable to
31 tolerate differences. Political difference does not mean all
32 compassion and empathy has to be suspended - this is a truism that
33 1 applies equally to right and left of the political spectrum.
34 1 For the same reason, condemnations of violence should never be
35 ' selective. BJPleaders' condemnations of beef lynchings, for example,
36 ' tend to be qualified and accompanied by caveats - which negate their
37 1 value and can be an incitement to further violence. Another issue that
38 5 is key to the future of the republic that Mukherjee touched upon, is
39 education. Unfortunately this gets stepmotherly treatment in the
40 political discourse, despite its critical and transformative role in
```

Activate Windows  
Go to PC settings to activate Windows

FOLDERS

02

◀▶ CT01\_results\_original.txt

01.txt

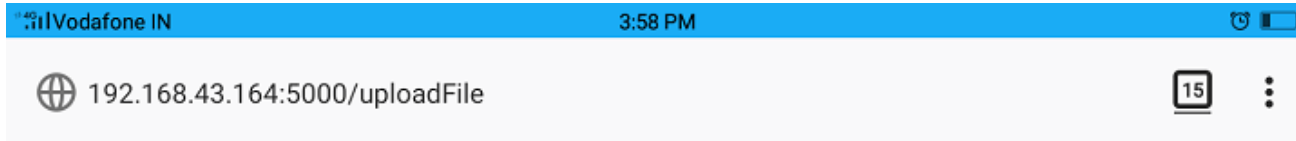
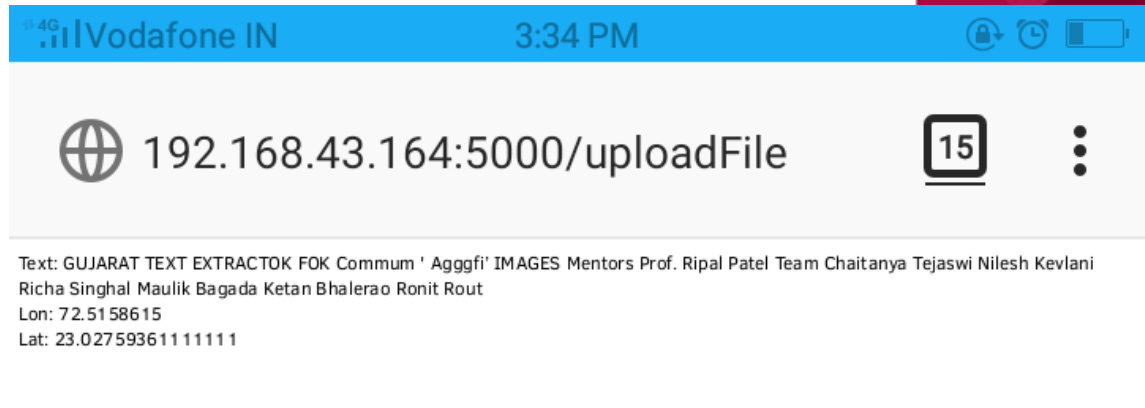
1 [A]  
2 [HEAD]  
3 Politics Of Empathy  
4 [HEAD-sub]  
5 In his farewell speech, President Mukherjee made an eloquent case for pluralism and education  
6 [TEXT]  
7 In his considered farewell speech President Pranab Mukherjee, rightly praised by finance minister Arun Jaitley for the utterly non-partisan manner in which he has conducted himself in Rashtrapati Bhavan, brought up the important issue of the relationship between violence and contemporary public discourse, where he invoked the spirit of Mahatma Gandhi to argue that when such a discourse becomes tinged with intolerance, it makes the atmosphere conducive to greater violence as well as marginalisation of less privileged groups. In today's atmosphere of polarised debate the departing president's fervent plea - "we must free our public discourse from all forms of violence, physical as well as verbal" - needs to be heard.  
8 Pluralism has been a hallmark of human civilisation, and it's something to which we must hold fast. The battle of ideas need not get personal - there are other, superior ways of winning it than to suggest that one's opponent is anti-national. Debate has to be civil and rational, not hysterical and paranoid; those who suggest their political opponents should be exiled to Pakistan are in fact making India a bit like Pakistan, which is unable to tolerate differences. Political differences does not mean all compassion and empathy has to be suspended - this is a truism that applies equally to right and left of the political spectrum.  
9 For the same reason, condemnations of violence should never be selective. BJP leaders' condemnations of beef lynchings, for example, tend to be qualified and accompanied by caveats - which negate their value and can be an incitement to further violence. Another issue that is a key to the future of the republic that Mukherjee touched upon, is education. Unfortunately this gets stepmotherly treatment in the political discourse, despite its critical and transformative role in India being able to meet challenges of the future.  
10 In this context, Union road transport minister Nitin Gadkari's statement that India will not allow driverless cars fearing massive loss of jobs is a giveaway of our confidence levels in the face of disruptive change. The spectre of rampant unemployment is a distinct possibility and it makes eminent sense to tame the divisive forces and violent proclivities that are in danger of cutting loose. Pluralism and education are the keys to securing India's future, and Mukherjee did well to lay down the options for both government and citizens in his expansive farewell speech.  
11  
12 [B]  
13 [HEAD]  
14 Indian Cue Masters League in Ahmedabad from Aug 19  
15 [HEAD-sub]  
16 Solomon.Kumar @timesgroup.com  
17 [TEXT]  
18 Hyderabad: "For a sport that originated in the country, it is a really sad state of affairs that we have to tell people what we do. When I tell people that I am a snooker player, they seem clueless about it," India's top snooker player Vidya Pillai's lament captures the angst of the people in cue sports.  
19 There are about 2.5 million amateur players in India but the sport still has not captured the imagination of the masses. It is considered an elitist game even though Indian players have won 42 world titles over the years. To break the ice and make the sport appeal to a larger section of the society, BSFI and Sportzlive Entertainment launched the Cue Slam Indian Cue Masters League here on



Activate Windows  
Go to PC settings to activate Windows.



# Example 2



Text: ५ इसर्ना डॅडागठ अंतरिक्ष उपयोग केंद्र (इसरो) SPACE APPLICATIONS CENTRE (ISRO) भारत सरकार, अंतरिक्ष विभाग Govt. of India, Dept. of Space  
अहमदाबाद — 380015 AHMEDABAD — 380015  
Lon: 72.5158615  
Lat: 23.02759361111111



# Limitations

- **Limited to News Articles from a Single Newspaper with Consistent Fonts**



# THE TIMES OF INDIA

## Quoting Nehru, Indira & Rajiv

# END OF A LEGEND

- **Inadequate Number of Test Images**  
10 Random, 10 Ordered

- **Angular Defects & Incident Light Intensity Variations**



# Conclusion & Future Scope

- **Conclusion**

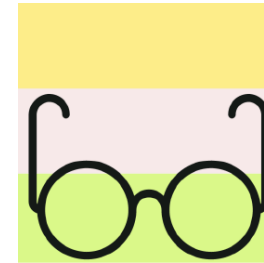
By making use of the proposed scheme on datasets of about 3500 characters each, we have observed a two-fold improvement in the character recognition efficiency of Tesseract OCR, in the case of digital newspaper article images secured using a smartphone camera.

- **Future Scope**

- *Extension (of model) to multiple Indian languages*



- *Extension to other newspaper formats*



- *Comparison & Objective Evaluation with other Binarization Techniques*



Niblack

Sauvola et al.

Wolf et al.





# References

## Patents:

- [\[1\]](#) US6577762B1, Background Surface Thresholding
- [\[2\]](#) US7400768B1, Enhanced optical recognition of digitized images through selective bit-insertion.
- [\[3\]](#) US9298980B1, Image preprocessing for character recognition.
- [\[4\]](#) US20120063690A1, Object-Based Optical Character Recognition Pre-Processing Algorithm.
- [\[5\]](#) US7106905B2, Systems and methods for processing text-based electronic documents.
- [\[6\]](#) US20130329023A1, Text recognition driven functionality.

## Literature:

- [\[7\]](#) Eugene Borovikov, A survey of modern optical character recognition techniques.
- [\[8\]](#) M Seeger, C Dance, Binarising camera images for OCR (ICDAR 2001, Proceedings of the 6th International Conference on Document Analysis and Recognition).
- [\[9\]](#) Ranjith Unnikrishnan, Ray Smith, Combined Script and Page Orientation Estimation using the Tesseract OCR engine (ICDAR '07 Proceedings of the Ninth International Conference on Document Analysis and Recognition).
- [\[10\]](#) Ray Smith, An Overview of the Tesseract OCR Engine (MOCR '09 Proceedings of the International Workshop on Multilingual OCR).
- [\[11\]](#) Ray Smith, Daria Antoniva, Dar-Shyang Lee, Adapting the Tesseract Open Source OCR Engine for Multilingual OCR (MOCR '09 Proceedings of the International Workshop on Multilingual OCR).
- [\[12\]](#) JJ Sauvola, Tapio Seppänen, Sami Haapakoski, Matti Pietikäinen, Adaptive Document Binarization (ICDAR '97 Proceedings of the 4th International Conference on Document Analysis and Recognition).
- [\[13\]](#) Niblack, W (1986), An introduction to Digital Image Processing, Prentice-Hall, pp. 115–116.
- [\[14\]](#) Nobuyuki Otsu (1979). "A threshold selection method from gray-level histograms". IEEE Trans. Sys., Man., Cyber. 9
- [\[15\]](#) OpenCV Reference Manual (OpenCV3.0 Documentation), pg. 265, 294-295.
- [\[16\]](#) Lukáš Neumann and Jiří Matas. Text localization in real-world images using efficiently pruned exhaustive search. In *in Document Analysis and Recognition, 2011 International Conference on. IEEE, 2011*, pages 687–691.



# Thank you

ALL THE LOVE  
ERAN & ENCHERY

