# Localization of Articles from a Newspaper Image

November 29, 2019

A MINI-PROJECT REPORT SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS FOR THE

## B.Tech Degree in Computer Science and Technology

Submitted by:
**Sushmita Mondal(Roll No:510518028)**
**Subhadeep Roy(Roll No:510518033)**
**Soumik Sen(Roll No:510518036)**
**Visesh Sharma(Roll No:510518061)**
**Devesh Singh(Roll No:510518073)**

under the supervision of

# Prof. Samit Biswas



1856 १८५६
उत्तिष्ठत जाग्रत प्राप्य वरान् निबोधत
INDIAN INSTITUTE OF ENGINEERING SCIENCE AND TECHNOLOGY, SHIBPUR
भारतीय अभियांत्रिकी विज्ञान एवं प्रौद्योगिकी संस्थान, शिवपुर

November 29, 2019

# Localization of Articles from a Newspaper Image

**Sushmita Biswas(Roll No:510518028)**
**Subhadeep Roy(Roll No:510518033)**
**Soumik Sen(Roll No:510518036)**
**Visesh Sharma(Roll No:510518061)**
**Devesh Singh(Roll No:510518073)**

under the supervision of

## Prof. Samit Biswas

Department of Computer Science and Technology

IIEST-Shibpur

India-711103

# ACKNOWLEDGEMENT

(Sushmita Mondal-510518028)

(Subhadeep Roy-510518033)

(Soumik Sen-510518036)

(Visesh Sharma-510518061)

(Devesh Singh-510518073)

# CERTIFICATE

This is to certify that the mini-project titled "*Segregation of Articles from a Newspaper Image*" has been carried out by **Sushmita Mondal(Roll No - 510518028), Subhadeep Roy(Roll No - 510518033), Soumik Sen(Roll No - 510518036), Visesh Sharma(Roll No - 510518061) and Devesh Singh(Roll No - 510518073)** under my supervision and guidance in their third and fourth semesters in partial fulfillment of the requirements for the B.Tech Programme in the Department of Computer Science and Technology, IIEST Shibpur.This work has not been reported anywhere for any other purpose.

(Samit Biswas)

# Contents

# List of Figures

# A   INTRODUCTION

The newspaper is the most powerful of all the means of expression of the news and views about men and matters.Newspapers are regarded by economists as a necessity of modern life. With the growth of literacy and the development of the means of communication they are playing a very important role in nowadays' society.

The functions performed by newspapers cannot be over-emphasized in any society in general and in academic institution in particular. So several organizations like libraries and colleges try to preserve the important articles of old newspapers. But they encounter several problems in managing newspapers.

Newspaper preservation is a challenge because newsprint is an inherently unstable paper. Formulated to be inexpensive and expendable, newsprint is manufactured with large percentages of unpurified wood pulp which contains impurities that remain in the paper after processing.These impurities when exposed to light, high humidity and atmospheric pollutants, promote discoloration and acidic reactions in the paper. Acidity causes the paper fibres to weaken and break, and is the major culprit in causing the paper to become brittle.All these shortcomings make newspaper preservation a tough challenge.

When newspapers or clippings are valued most for the information they contain, and not as artifacts, copying the information onto a more permanent quality paper should by undertaken.Still they can decay. But an image of a newspaper will remain as long as it is kept. So taking images is a far better alternative of preserving old newspapers.

The project mainly focuses on separating the different articles present in a newspaper image so that they can be analysed and segregated under different genres.

## A.1   Motivation

The budget that goes for the management of newspapers in an average library is a lot yet it is not enough to properly preserve the newspapers. Libraries encounter several problems in managing newspapers. Librarians are confronted with the problems of handling the papers because of their fragility; preserving them because of their information contents, and housing them because of lack of space and the rapid rate at which they are churned out. The system of storage and retrieval increases the incidence of destroying and or tearing the newspapers most especially when tying and untying the papers.

Also these newspapers are needed by readers who in times feel the nessecity to take them home. But this causes the newspapers to be damaged making it more and more hard to preserve. Even searching old articles become tough and time consuming for the fragility of old newspapers.

Therefore, the motivation behind the project is to simplify the process of newspaper preservation and make it more economical.

## A.2   The Proposed Project

The programme that is being proposed under this project takes the image of a newspaper page and separates all the articles present in it. It then makes images of all these articles and saves them.

## A.3  Outline Of Project

Section 2: provides the details of all the programmes and applications used to make the project.

Section 3: covers the programme written and all the methods written.

Section 4: discusses the limitations of the project and how it can be improved.

Section 5: discusses how the project can be used.

Section 6: References.

# B  PREREQUISITES

All the programmes and the applications used to make the project are listed below:

## B.1  Eclipse IDE

Eclipse is an Integrated Development Environment (IDE) used in computer programming. It contains a base workspace and an extensible plug-in system for customizing the environment. Since the final project is written in JAVA Programming Language, Eclipse is used as the IDE.
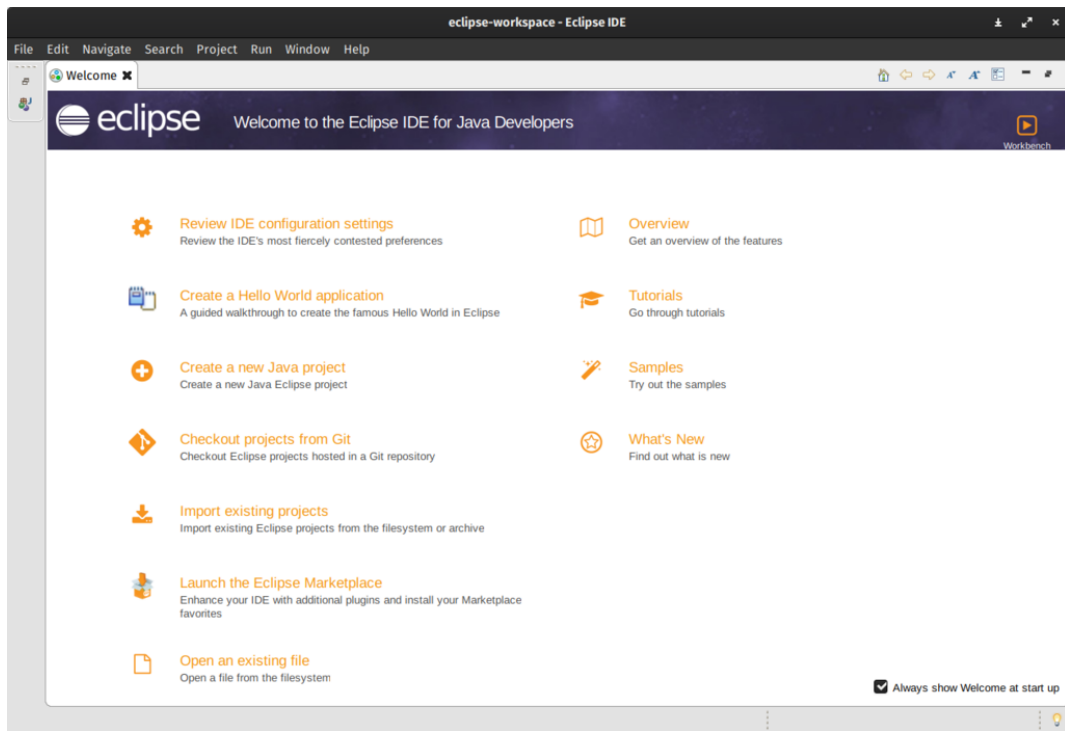
Figure 1: Eclipse IDE

## B.2 Knowledge about images and Image Handling

To get a little amount of knowledge about images and image handling programmetically a side project was done. This project consisted of writing a programme in C using VI Text Editor. The programme takes an .netpbm format image as input does the following operations:

1. **Duplicate an Image**
2. **Making a Histogram of all the colours used in the Image**
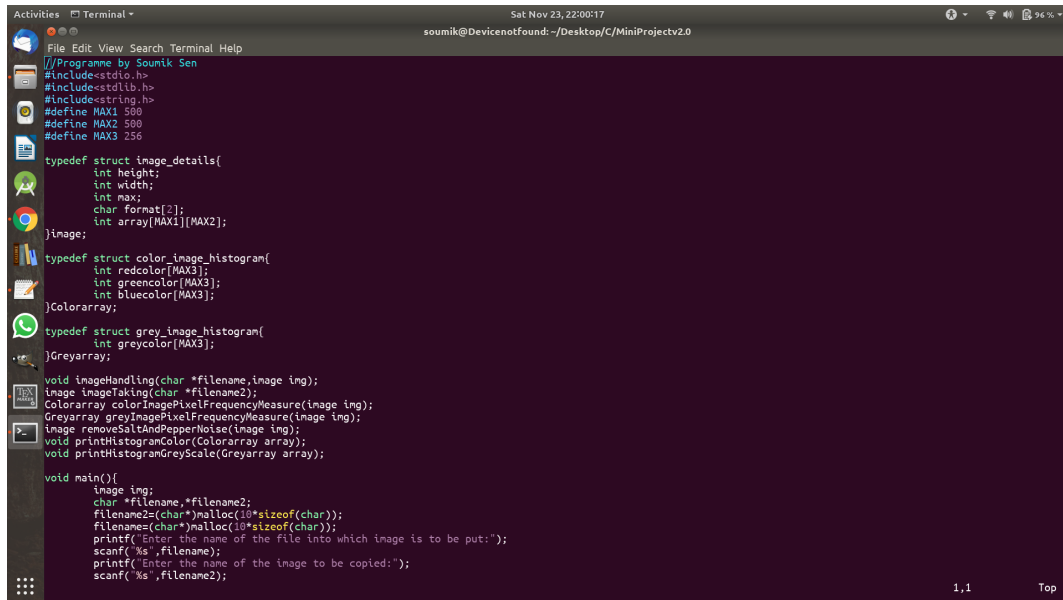3. **Removing Salt and Pepper Noise from the Image**

Figure 2: Code picture

## B.3    LaTeX

LaTeX is a document preparation system. When writing, the writer uses plain text as opposed to the formatted text found in WYSIWYG word processors like Microsoft Word, LibreOffice Writer and Apple Pages. The writer uses markup tagging conventions to define the general structure of a document. LaTeX with TEXMAKER as IDE has been used to write the project document.



Figure 3: Latex Logo

11

# C   PROGRAMME

The programme has been written in JAVA 11 Programming Language on Eclipse IDE. The following steps were followed:

## C.1   Creating Project

**Step 1:** In the *Welcome to the Eclipse IDE for Java Developers* window, Click *Create a new Java Project.*

**Step 2:** In the New Java Project screen, enter the following values:

*Project name:* Article Segregation

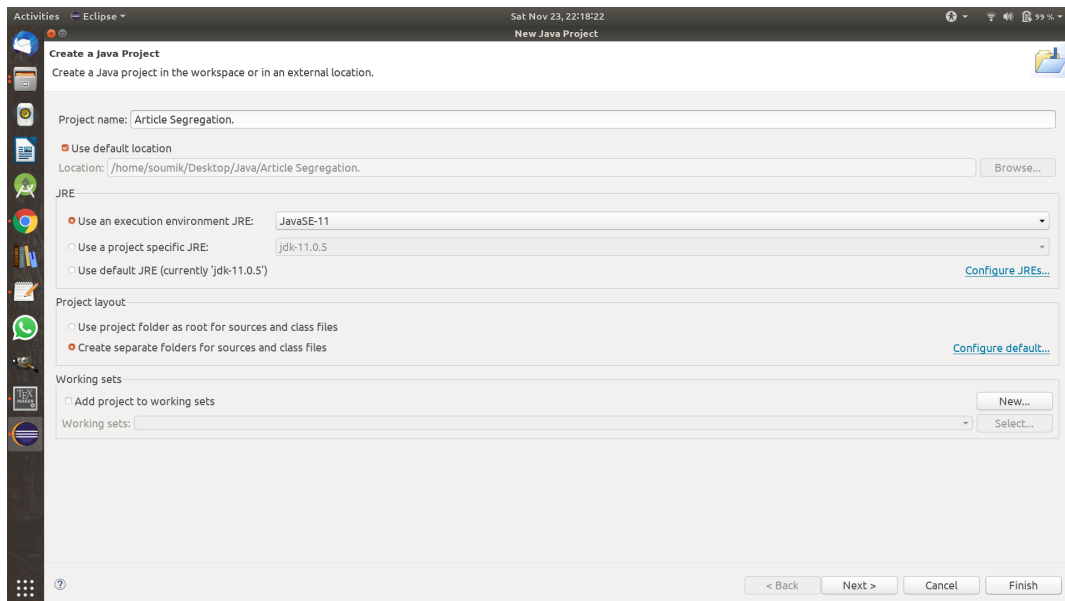*Use an execution environment JRE:* JavaSE-11



Figure 4: Creating Project

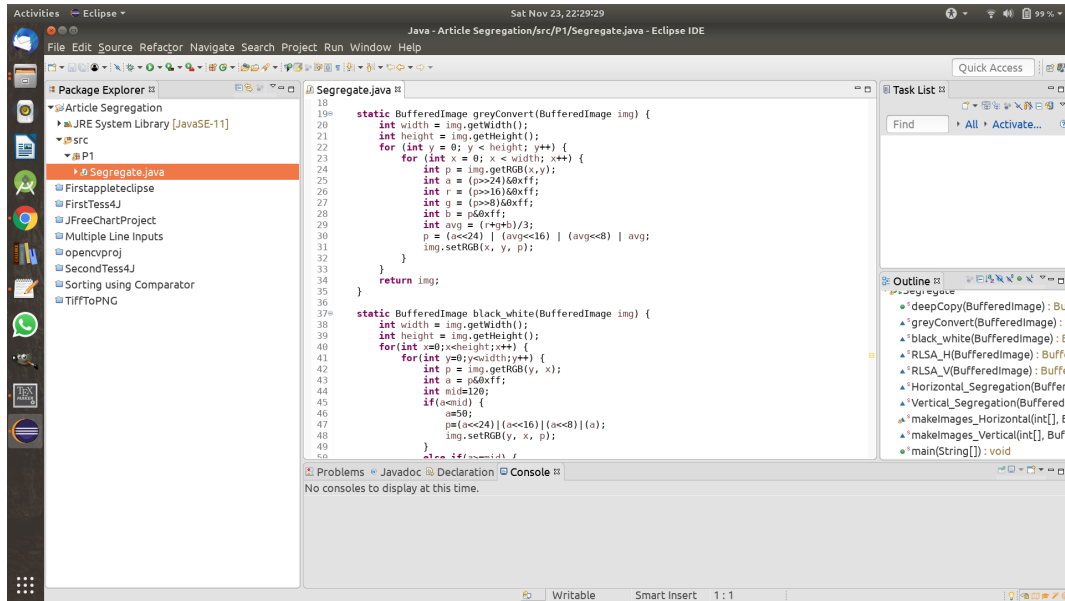Finally click finish project.

## C.2   Article Segmentation



Figure 5: Project Code

The left tab shows the Project Structure. The right tab named *Outline* shows the different methods used in the project.

## C.3   Explaining the Methods

***deepCopy:*** This method copies the contents of a BufferedImage object into another BufferedImage object.

***black-white:*** This method converts a colour image into a black and white image.

***Horizontal-Segregation:*** This method cuts the image into horizontal segments according to white spaces.

***Vertical-Segregation:*** This method cuts the image into vertical segments according to white spaces.

***makeImages-Vertical:*** This creates small subimages of the vertically segmented images obtained from *Vertical-Segregation* method.

***makeImages-Horizontal:*** This creates small subimages of the horizontally segmented images obtained from *Horizontal-Segregation* method.

## C.4   Process of Separation

The main idea of the project is to separate the different articles of the input images. To do so we try to identify the articles by the blank spaces between them. If there is a horizontal space such that it reaches from left end to right end then the *Horizontal-Segregation* method cuts the image there. Similarly if there is a vertical space that extends from the top to the bottom of the image then the *Vertical-Segregation* method cuts the image. Then all the cut pieces are converted into images by the *makeImage* methods.

**Example:**



Figure 6: Input Image

We first convert the colour image to binary image by taking 120 pixel as a midvalue and making any value above 120 pixel to 255 pixel(white) and all values below 120 pixel to 0 pixel(black).

We take 120 as midvalue as it is the median of all the pixels present in the image.



Figure 7: Binary Image Conversion

After the Conversion to Black and White Image we first cut the image horizontally. To show that we use the *horizontal summation graph* of the image.

Figure 8: Horizontal Graph

This graph has 900*255 = 229500 as the x-axis length. This value is the value that comes when we add all the white pixels in a row of length 900. So the image has 3 such rows at 0 height, at 889 height and at 1089 height from the bottom. The segmented images are :



Figure 9: First Segment

Figure 10: Second Segment

Now we localize the two segments vertically. To show that we use the *vertical summation graph* of the image.



Figure 11: Vertical graph

The y axis has maximum value of 1089*255 = 277695 pixels which is the summation value when an entire column of the image is white. So whenever the value reaches this maximum value there is a white column and the image is segmented from there, similar to horizontal segmentation. The resulted segments are given below:

18

(a) Segment

(b) Segment

(c) Segment

(d) Segment

Finally we collect all the segments and recreate them from the original colour image to get all the localized articles. Thus, the algorithm works in this manner.

**RESULTS:**



Figure 13: The Section



(a) Segment



(b) Segment



(c) Segment



(d) Segment

# D  LIMITATIONS

## D.1  Limitations

There are several limitations in the project. They can be listed as below:

**Image Resolution:** The input image must be of high resolution and large sized. If otherwise the segmentation will not take place.

**Image Orientation:** The Image must be in vertical orientation and must be perfectly straight.

**Noise:** The Image must be free of any noise or else erroneous outputs will come.

**Nested Articles:** In case of non-aligned nested articles the subimages must be again run to get all the articles separately.

## D.2  Solution to the limitations

The limitations can be overcome by taking some measures like :
**1.** Removing noises by noise removal techniques.
**2.** Linking all the words of a same article by blackening the spaces between them. This will help identify even nested articles easily.

# E  FUTURE OF THE PROJECT

The project can be easily used as a starting platform for other projects like converting the articles into *.txt* files by using some *Optical Character Recognition* (OCR) tools. Then these *.txt* files can be segregated under different topics and genres.

Also this project will be a good platform to start other projects like Segregating Chapters of a Book.

# F REFERENCE

[**1**] Wikipedia
*https://www.wikipedia.org/.*
[**2**] Stack Overflow
*https://stackoverflow.com*
[**3**] Overleaf
*https://www.overleaf.com*
[**4**] Image Processing in Java
*Douglas A. Lyon*