

Shahjalal University of Science and Technology  
Department of Computer Science and Engineering



Document Decomposition System for Bangla OCR

Biswajit Debnath, Reg. No: 2011331009, Dept. of CSE  
Md. Sajid Shahriar, Reg. No: 2011331048, Dept. of CSE

Adviser:

Dr. Muhammed Zafar Iqbal, Professor, Dept. of CSE

Co-Adviser:

Mr. Sabir Ismail, Assistant Professor, Dept. of CSE  
Mr. Md. Saiful Islam, Assistant Professor, Dept. of CSE

18<sup>th</sup> October, 2016

# Document Decomposition System for Bangla OCR



A Thesis submitted to the Department of Computer Science and Engineering,  
Shahjalal University of Science and Technology, in partial fulfillment of the requirements  
for the degree of Bachelor of Science in Computer Science and Engineering.

Name: Biswajit Debnath, Reg. No: 2011331009, Dept. of CSE  
Name: Md. Sajid Shahriar, Reg. No: 2011331048, Dept. of CSE

Adviser:

Dr. Muhammed Zafar Iqbal, Professor, Dept. of CSE

Co-Adviser:

Mr. Sabir Ismail, Assistant Professor, Dept. of CSE  
Mr. Md. Saiful Islam, Assistant Professor, Dept. of CSE

18<sup>th</sup> October, 2016

## **Recommendation Letter from Thesis Adviser**

The Thesis

entitled "**Document Decomposition System for Bangla OCR**"  
submitted by the students

**1. Biswajit Debnath**

**2. Md. Sajid Shahriar**

is a record of research work carried out under my supervision and I, hereby, approve  
that the report be submitted in partial fulfillment of the requirements for the award  
of their Bachelor Degrees.

---

Adviser : Dr. Muhammed Zafar Iqbal, Professor, Dept. of CSE  
Date : / 10 / 2016

---

Co-Adviser : Mr. Sabir Ismail, Assistant Professor, Dept. of CSE  
Date : / 10 / 2016

---

Co-Adviser : Mr. Md. Saiful Islam, Assistant Professor, Dept. of CSE  
Date : / 10 / 2016

## **Certificate of Acceptance of the Thesis**

The Thesis

entitled "**Document Decomposition System for Bangla OCR**"

submitted by the students

**1. Biswajit Debnath**

**2. Md. Sajid Shahriar**

on \_\_\_\_\_

is, hereby, accepted as the partial fulfillment of the requirements for the award of their Bachelor Degrees.

---

Head of the Department

Dr. Md. Reza Selim  
Professor, Dept. of CSE

---

Chairman, Thesis Committee

Dr. Md. Reza Selim  
Professor, Dept. of CSE

---

Adviser

Dr. Muhammed Zafar Iqbal  
Professor, Dept. of CSE

## ABSTRACT

Resources like newspaper or printed documents contain important contents. Content Digitizing is a possible solution of protecting and retrieving the endangered documents. Optical Character Recognition or in short OCR is the best option or process for digitizing these printed documents. Some pre-processing procedures like Segmentation or Decomposition help in Optical Character Recognition methodology. By following these pre-processing methodology, textual part of the document can be identified. These textual parts are the main component to document digitize. Our work is based on segmenting the Bangla Documents. Our goal was to decomposing the contents from a document along with de-skewing it, retrieving the tabular data and rebuilding it. The methodology we implemented work 96.88% perfectly.

**Keyword:** Bangla Document, Digitizing, OCR, Segmentation, Decomposition, DE skew, Tabular Data.

## Acknowledgements

This thesis work and paper is made possible through the help and support from all sentient beings. Especially, please allow us to dedicate my acknowledgment of gratitude toward the following significant advisers, co-advisers and contributors who was with us throughout the whole time:

First and foremost, we would like to thank to our Adviser Dr. Muhammed Zafar Iqbal for his most support and encouragement. His comments, remarks, engagement, guidance and constant supervision as well as for providing necessary information regarding the thesis was really helpful to us to do the analysis and finding out what to do.

Second, we would like to thank our Co-Adviser Mr. Sabir Ismail and Mr. Md. Saiful Saif to advise us about the thesis and to provide valuable information which was really help to accelerate our work.

Finally, we sincerely thank to our seniors Fahad Hasan and Tasmin Afroz Tannee. Who previously contributed with this work and after that help us to catch up with the work.

# Contents

Chapter 1 .....	1
Introduction.....	1
1.1    Meaning of Optical Character Recognition .....	1
1.2    Process of an Optical Character Recognition System .....	2
1.2.1    Preprocessing .....	3
1.2.2    Character Recognition.....	3
1.2.3    Post Processing.....	3
1.3    Importance of Optical Character Recognition.....	4
Chapter 2 .....	5
Preprocessing.....	5
2.1    What is preprocessing in OCR.....	5
2.1.1    De-Skew .....	6
2.1.2    Noise Removal.....	7
2.1.3    Binarization .....	8
2.1.4    Layout Analysis & Segmentation.....	9
Chapter 3 .....	10
Document Decomposition .....	10
3.1    What is Document Decomposition .....	10
Chapter 4 .....	11

Background Study.....	11
4.1    Reason of Studying Previous Research Works .....	11
4.2    Research on Document Analysis and Recognition.....	12
4.2.1    Decomposing Newspaper using Image Processing & Document Analysis	12
4.2.2    Decomposing using a Split and Merge Approach .....	15
4.2.3    Automated Text-Line Height Calculation for Font Size Detection	
	18
Chapter 5 .....	20
Research on Bangla Document Decomposition .....	20
5.1    Research on Bangla Document Decomposition in SUST .....	20
5.2    Decomposition Research from Batch 2006.....	21
5.3    Decomposition Research from Batch 2007.....	22
5.4    Decomposition Research from Batch 2010.....	25
Chapter 6 .....	29
Findings on Previous Methodology.....	29
6.1    Findings on 2011 batch Methodology .....	29
Chapter 7 .....	35
Current Methodology: Part 1.....	35
7.1    Noise Reduction.....	35
7.2    Shaded pixel removed.....	38
7.3    Edge Detection .....	39
7.3.1    Canny Edge Detector.....	39
7.4    Conversion of input image .....	42
7.5    Discerning the Element Separators.....	43

7.6	Segmenting the Black Boxes.....	43
7.7	Distinguishing the Elements .....	43
7.8	Line Height Calculation .....	48
7.9	Comparing Threshold Values with Previous Methodology .....	51
Chapter 8 .....		53
	Current Methodology: Part 2.....	53
8.1	Skew Detection and DE Skew.....	53
8.2	Using Hough Transform for DE Skew .....	53
8.3	Detecting Skew Contour and rotating by its slop .....	55
8.4	Reconstructing the Document.....	57
Chapter 9 .....		58
	Result Analysis.....	58
9.1	Accuracy for Current Decomposition Methodology.....	58
9.2	Accuracy Comparison with Previous (Batch 2010) Decomposition Methodology.....	60
Chapter 10.....		61
	Conclusion .....	61
10.1	Findings in our Developed Methodology.....	61
References.....		62

# List of Figures

Figure 1.2.1 The different areas of character recognition [1]. .....	2
Figure 1.2.2 Steps of Optical Character Recognition.....	4
Figure 2.1.1: A Sample Scanned Document. ....	5
Figure 2.1.2 An Example of Image De-skew [3] .....	6
Figure 2.1.3: Image Noise Removal Example.....	7
Figure 2.1.4: Binarization Example .....	8
Figure 4.2.1: Newspaper Image Decomposition: (a) Original, (b) Segmented Image .....	14
Figure 4.2.2: In a content locale, there are no overlapping boxes of associated components .....	16
Figure 4.2.3: Splitting Newspaper image into small zones.....	17
Figure 4.2.4: A part of correctly segmented. Each part is each labeled with different color.....	18
Figure 5.2.1: Input and Output for using the work of 2006 Batch.....	22
Figure 5.3.1: Input and Output for work from 2007 Batch.....	23
Figure 5.3.2: Sub Headline between Columns .....	23
Figure 5.3.3: A Complex Document. .....	24
Figure 5.4.1: An Input Image we used in 2010 Batch Methodology.....	26
Figure 5.4.2: Image detected from Document using 2010 Batch methodology..	27
Figure 5.4.3: Headline detected from Document using 2010 Batch methodology .....	28
Figure 5.4.4: Column detected from Document using 2010 Batch methodology. .....	28
Figure 6.1.1: A Scanned Image is zoomed to preview the noise.....	30
Figure 6.1.2: 96dpi image Binarized .....	31
Figure 6.1.3: Overlapping case are highlighted in red area.....	32
Figure 6.1.4: A Sample of a Table Component .....	33
Figure 7.1.1: Input Image for Document Decomposition .....	36
Figure 7.2.1: Image Matrix of Non Shaded Pixel.....	38
Figure 7.3.1: Edge Detected of the Input Image.....	41
Figure 7.4.1: Conversion Image .....	42
Figure 7.7.1: Detected Images.....	45

Figure 7.7.2: A Table with Vertical and Horizontal Histogram.....	47
Figure 7.7.3: Detected Headline.....	48
Figure 7.8.1: Line Overlapping Case.....	49

# **Chapter 1**

## **Introduction**

Humans like to dream big and their dreams are persistently something critical like machine replication of human limits, including reading, composing and so forth. From these dreams, machine reading has created into reality all through the latest five decades. Optical Character Recognition has wound up a champion among the best applications of advancement in the field of Computer Vision<sup>1</sup> [1].

### **1.1 Meaning of Optical Character Recognition**

Optical character recognition (optical character reader) (OCR) is the mechanical or electronic transformation of images/pictures which may be as typed, handwritten or printed text into machine-encoded text [2]. For data digitization OCR is now the most common method for printed texts so that it can be electronically edited, searched, stored more compactly, displayed on-line, and used in machine processes such as Machine Translation<sup>2</sup>, Speech Synthesis, Key Data and Text Mining [2]. It is a field of research in Pattern Recognition<sup>3</sup>, Artificial Intelligence<sup>4</sup> and Computer Vision. [2].

---

<sup>1</sup> Computer Vision - A field that includes methods for acquiring, processing, analyzing, and understanding images.

<sup>2</sup> Machine Translation - Performs simple substitution of words in one language for words in another.

<sup>3</sup> Pattern Recognition - Focuses on the recognition of patterns and regularities in data.

<sup>4</sup> Artificial Intelligence - It is the intelligence exhibited by machines or software.

It is widely used as a form of data entry from printed paper data records, whether passport documents, invoices, bank statements, computerized receipts, business cards, mail, printouts of static-data, or any suitable documentation.

## 1.2 Process of an Optical Character Recognition System

Optical Character Recognition is mainly an “offline” process, in which it performs some analysis on a static document. The offline work is performed after the writing or printing has been completed, as opposed to on-line recognition where the computer recognizes the characters as they are drawn. Both hand printed and printed characters may be recognized, but the performance is directly dependent upon the quality of the input documents.

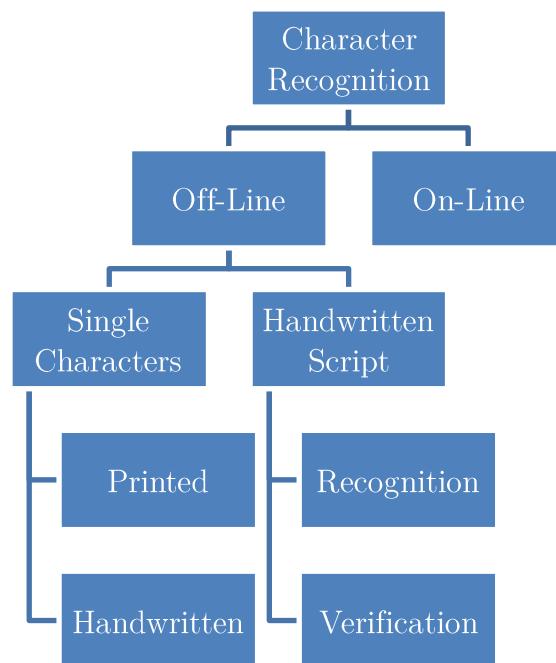


Figure 1 The different areas of character recognition [1].

The software made for OCR often follows three techniques or methods before giving the final output. It deals with the problem of recognizing optically processed characters.

### 1.2.1 Preprocessing

- **De-skew** - Aligning the document properly.
- **Binarization** - Convert the image from color or grayscale to black and white.
- **Layout Analysis** - Identifies columns, paragraph, captions, images etc.
- **Segmentation** - Multiple characters that are connected in the image artifacts must be separated.

### 1.2.2 Character Recognition

- An image is stored on pixel by pixel basis.
- Input glyph being isolated from the rest of the image then used pattern matching techniques for recognition.

### 1.2.3 Post Processing

- **Grouping** - Merging the symbols output to string.
- **Error Analysis** - Error detection and correction

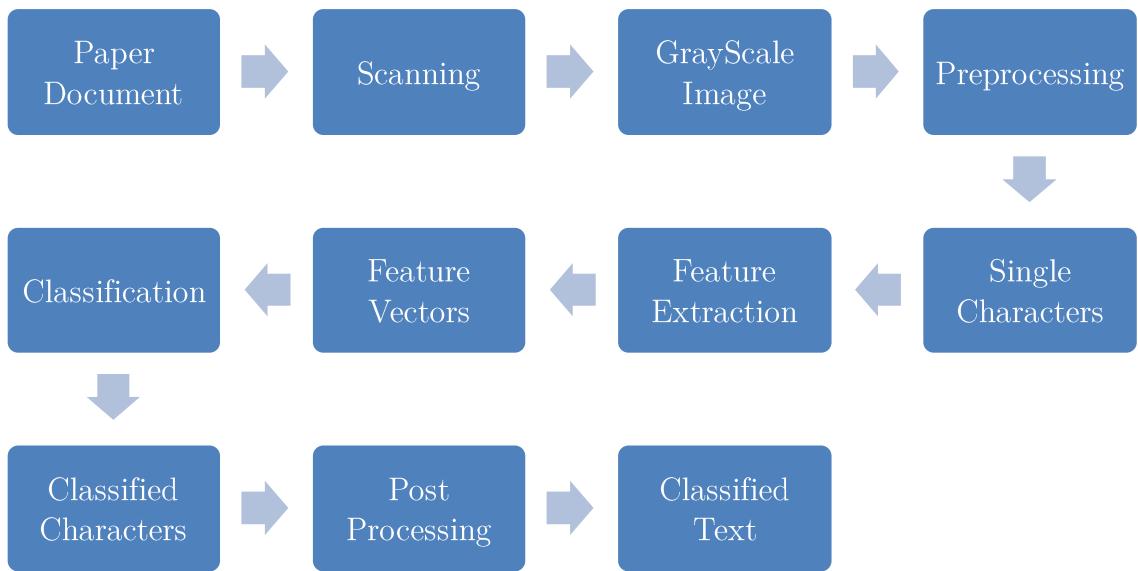


Figure 2 Steps of Optical Character Recognition.

### 1.3 Importance of Optical Character Recognition

If we follow the very traditional way for digitizing the data, it will need a huge manpower, because we have to create each digital document by giving the inputs manually. For cases like these, automatic identification may be an alternative and OCR is the best alternative of that process. OCR converts normal scanned documents to text-searchable, so it allows us to search content on the same. A scanned file is usually stored as an image if not converted into a text-searchable file with the help of OCR technology. Once that is done, the document allows the search engine to perform content search inside the file and retrieve appropriate results. OCR also can identify any kind of text inside the document if it is wrong.

## Chapter 2

# Preprocessing

Preprocessing the most crucial part to start the digitizing the scanned or hand written document data. Process like De-skew, Binarization, Layout-Analysis are included in pre-processing.

## 2.1 What is preprocessing in OCR

Technically documents like newspaper not only contains just some textual data but also it is populated with figures, tabular datasheets, and images may contain many other contents. Along with this kind of component variation scanned document could be skewed. Sometimes document may contain unnecessary points which could create a messy result.



Figure 3: A Sample Scanned Document.

To improve the chances of successful digitizing process of the documents images need to be de-skewed, noise removed and segmented. These steps are called the preprocessing in OCR. A brief description is given below about these and the steps are done for it:

### 2.1.1 De-Skew

While scanning a document, there is always a chance it to be skewed. So a problem could be arise. In later steps of preprocessing each components need to be separated, which is much needed for corrected analysis the document. Part of image or part of a textual area could be cut off because of a skewed document. As a result the output we will get may not the match with the document. Checking and de-skewing the skewed image will lower the error rate at a great rate.

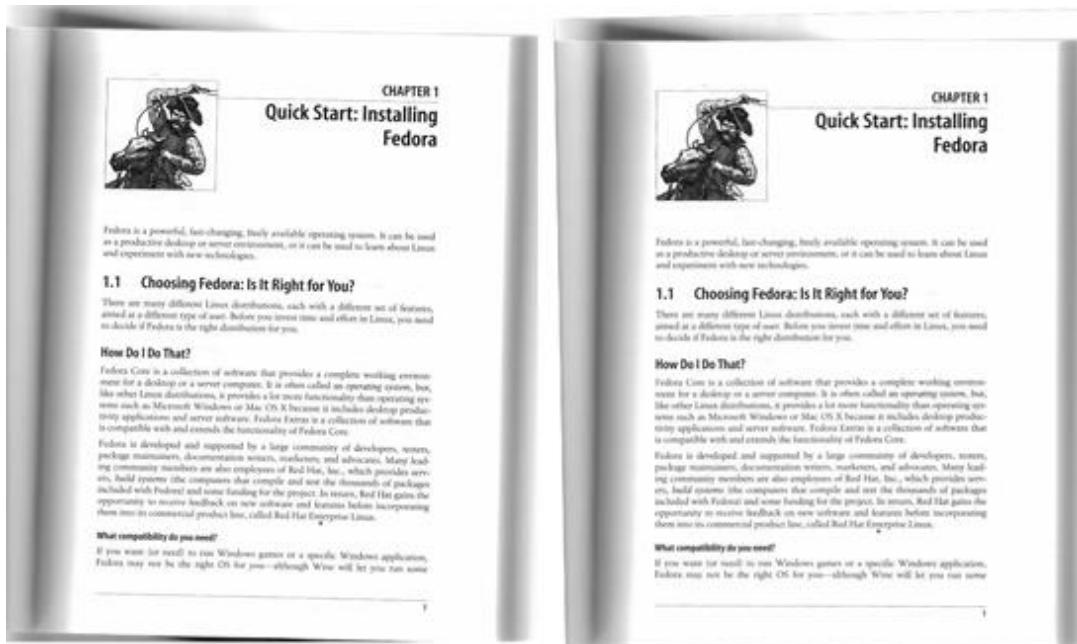


Figure 4 An Example of Image De-skew [3]

De-skewing a textual image is done by casting rays from left to right and find how many blackish pixels the rays intersect with. Then, by rotating the angle of the rays by a small amount and cast the rays again, and doing the same process again and again we can finally de-skew our document [4].

### 2.1.2 Noise Removal

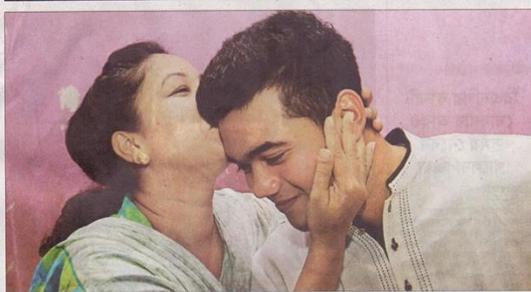
Newspaper like documents are printed in Newsprint paper. It contains a lot of noise like positive or negative spots. Old documents also have a high rate of having noise in it. In the steps of recognition these positive or negative spots increase the chances of bad impact of analysis. Documents should be and needed to be preprocessed by removing those noise and have smoothed edges.



Figure 5: Image Noise Removal Example

### 2.1.3 Binarization

Binarization means creating a “binary image” from the colored or grey scale image document. The reason it is called binary image is mainly it contains only black and white color. Like 0 and 1. The task of Binarization is performed as a simple way of separating the text (or any other desired image component) from the background. Binarization influences the OCR stage to a significant extent the quality [2, 5].



টি-টোরেটি বিশ্বকাপ শেষ হয়ে গেল হচ্ছেই, আইনিসির বোঝি নিষেধাজ্ঞা মাধ্যমে ওপর। হতাশা আছে। কিন্তু মা সবিনা ইয়াসমিনের হেফ্পথ অর প্রেরণা তাসকিন অভিযন্তে নিষ্পত্তি প্রত্যায়ী করে ভল্লুকে মাটে ফিরতে। কাল তাসকিনের মোহাম্মদপুরের বাসায় ৩ শামসুল হক

## দুজনের দেশে ফেরা

ରାନ୍ଧା ଆକହାସ

‘ଭାରତେ ଡିଭାର ଖେଳେ ଖେଲିଲେ, ଏକିବାର  
ମା କୋଣେ ଦୂରାକ୍ଷାନ ନିଯେ ଖିଲେବେ ଓ  
ଆମେ ଆମିଶରେ କଟି ଆକପ୍ତ । ଯାହା  
ଉଠେଇଛନ୍ତି । ଅପେକ୍ଷା କରିଲେ ହେଲେ  
ଆହମ୍ବଦୀର ଜଳା । କଲକାତା ଥେବେ ସବୁ  
ଶାହମନ୍ଦିର ବିମାନଘରେ ଏବଂ ପୌରୀ  
ତାର । ମୋହାରିଙ୍ଗର ଜାରି ହେଲେ ନାହାଇଲା  
ହେଲେର ଜଳା ଅପେକ୍ଷା ଆବଶ୍ୟକ ରଖିଲା ।

A portrait of a man with dark hair and a beard, wearing a light-colored shirt. The background is a textured, reddish-brown surface.

বিশ্বকল খেতেও প্রয়োজন। অশালোরা বিপ্লবে প্রথম যাচাই উল্ল সন্দেহের মধ্যে আঁকড়ে থাকেন অভিযন্তা। প্রথমে নির্মাণ করাকেও বিশ্বাস করতে পারেনি তিনি। সামগ্রে হেসেইলেন প্রথম পূর্বের সংজ্ঞায় ১৫ মার্চ হোয়াইয়ে নিতে যান পর্যবেক্ষণ ধর্ম না কাটিও এবং দিন ইন্ডেন পরিকল্পনার বিপ্লবে মাটি ও পরীক্ষা করে হয়েছে, যদি তাঁর পক্ষেই আসবে—এই ই-

হতাম। এভাবে একটি কৃতিত্ব হবে, তা কে  
বেরিয়ে দেয়। কৃতিত্বের ধারণা কৃতিত্ব  
দেয়। কিন্তু এমন কৃতিত্ব সহজ কি এসেছে? উত্তর  
হাতের তাপমাত্রা বলেন, “নাই।” কোর্টে  
পতেকে, ধারণা সহজ দেয়। কিন্তু এমন কৃতিত্ব  
পরিষ্কারভাবে পরিচিনি। এমন অভিজ্ঞতা এই প্রথম।

দুই দিন আগে সেখে বিভিন্নের অভিযোগে  
আক্ষেপে অভিযুক্ত হয়ে থাকলেও সেখনে আক্ষেপ  
বোলাবে আরাফাত সামিন। এসেই দেন

Figure 6: Binarization



ଟି-ଟୋରେଟି ବିଷ୍ଟକାପ ଶେ ହେଁ ଗେଲ ହାତ୍ତେ, ଆଇସିମିର ବେଳିଙ୍ ନିଷେଧାଳୀ ମାଧ୍ୟମ ଓପର । ହତାଶ ଆଛେ । କିନ୍ତୁ ଯା ସାବିନୀ ଇହାସମିନେର ମେହମର୍ମ ଅର ପ୍ରେସର ତାତ୍କାଳିକ ଅହମଦଳ ନିର୍ଦ୍ଦେଶ ପ୍ରତ୍ୟାକ୍ଷର କରେ ତଥା ମାଟେ ଫିରିବେ । କଳ ତାତ୍କାଳିନେର ମୋହାର୍ଦଦପରେ ବାସାର ॥ ଶାମମଳ ହକ

## দুজনের দেশে ফেরা

ଶାନ୍ତି ଆକାଶ ୧

A black and white portrait of a man with short, dark hair, wearing a light-colored shirt. He is looking slightly to his left. The background is dark and textured.

### Image Before Binarization

### Image After Binarization

Figure 6: Binarization Example

#### **2.1.4 Layout Analysis & Segmentation**

After de-skewing, removal of the noise and the Binarization, identifying the columns, paragraphs, captions, images etc. need to be detected as distinct blocks. This identification process is done by analyzing the document.

As OCR is the system through which any document is transformed into digital data. Now it is a fact that OCR can process texts only. If it is given an input, which contains non-textual components, garbage texts will be produced as output. So, if we want OCR to work perfectly is must be ensured that it is not fed any kind of non-textual input.

## Chapter 3

# Document Decomposition

In one of the previous discussion we have said that, Documents aren't always contain a simple layout. It may be as complex as having table, columns, images etc. So avoiding error on character detection stage, decomposing is the most needed pre-processing step. Automation of decomposition for a document plays a vital role when the system meets with a huge archive of documents. At that time automated decomposition process with higher accuracy ensures the most time saving decomposing.

### 3.1 What is Document Decomposition

Document decomposition is mainly a part of Layout Analysis. As it said that documents could have various contents, and not all the contents are needed in the step of recognition. So if these components are separated through e process than we can only use the textual part for the recognition stage.

At the document decomposition phase, the document image is decomposed in its basic principal components, such as text areas, titles, articles and the close contact of different segments hinders the application of many standard page segmentation algorithms. [6, 7].

## **Chapter 4**

# **Background Study**

To drive the study or action of the breaking down, research is a verifiable necessity required in order to progress further. For separating the reports and in system of making digitize data center, various examination is so far going on to breakdown the records. Many of these research work was done for English documents. As there are always structural difference between documents to documents, working or contemplating on decomposing complex Newspaper archives may come convenient as far as future work. In this section we discussed about the necessity of studying about other research works what type of research is done in the past and what we observed from various research.

### **4.1 Reason of Studying Previous Research Works**

If we can find out that there has been some of the work was done before or if the work fits with some of the basic points with any research work, studying those research works will help a lot to get on the track of own research. If there is any kind of similarity, then it is obvious that studying those work will surely give a concept to develop own work methodology. Understanding on the work also depends on studying about the topic.

## **4.2 Research on Document Analysis and Recognition**

To get a brief knowledge about our work, we started to look out others work on decomposing scanned documents. We find out a various number of literature for working in order to learn and work. These papers are focused on decomposing an English based document. On Document analysis and decomposition an international academic conference in the banner of “International Conference on Document Analysis and Recognition” is held in every two years of gap [8]. Our target was to study on segmenting documents like newspaper, as it is compiled with a lot of complex contents. In September, 1999 this conference was held on Bangalore of India [8]. The specialty of this conference was arranging the First International Newspaper Segmentation contest. So this paper came out handy for us to start or study.

### **4.2.1 Decomposing Newspaper using Image Processing & Document Analysis**

To recognize a newspaper articles and segmenting its pages, an approach was introduced by B. Gatos, S.L. Mantzaris, K.V. Chandrinos, A. Tsigris and S.J. Perantonis whom are from Lambrakis Press, Athens Greece and Institute of Informatics and Telecommunications, National Research Center “Demokritos”, Athens Greece. They introduced their work on the Fifth International Conference of ICDAR’99 held on Bangalore [9].

They proposed a rule-based approach to article tracking and identification. This proposal was mainly focused on newspaper images. The most significant findings of them is the separating the complex layout of newspaper pages, in particular the oldest ones, where text columns are located very close to each other in a haphazard way, as well as the poor scanning results obtained because of low print quality or deterioration through time. The rules they set to decompose the newspaper is described below:

- **Page Segmentation:** This division is centered around decomposing the essential segments of the daily paper, for example, segments, features pictures, lines and so on [9]. In a newspaper as the layout is organized in a haphazard way, the contact of different segments are so close that the application of many standard page segmentation algorithms fails [10, 11]. So they proposed smearing and labeling of regions and on gradual extraction of image components in the following order [9]:
  1. Lines.
  2. Images and Drawings
  3. Background Lines
  4. Text and headline Blocks.

They find out that vertical background need to preserve as a vertical border if the text columns nearly in touch of each other.

- **Line Extraction:** In their proposal they identified that horizontal and vertical foreground lines are usually too close to other kinds of segments, they are to be the first elements to be identified then all other region categories will be more effectively extracted. The approach they took for line identification is based on Hough Transform as well as on Morphological transformations [12]. This technique gives the accurate results combining with fast implementation and can be summarized in the following steps:

They subsampled the input image with respect to foreground pixel and then find out the vertical and horizontal lines of a newspaper. Another sub sampling was done with respect to background pixel. From the result the image need to be extract into two grayscale images which will be Foreground Horizontal ( $F_H$ ) and Foreground Vertical ( $F_V$ ). For each line in vertical and horizontal the pixel length needed to be assigned by creating a property of minimum length and maximum

width line with interesting characteristics will be identified. The segment whose height are greater than the largest expected headline letter was chosen then.

Again with the respect to background pixel another sub-sampling is done. With this sub-sampled image is the horizontal projections are examined with the FFT. If there are no dominant frequencies found then the block was classified as image.

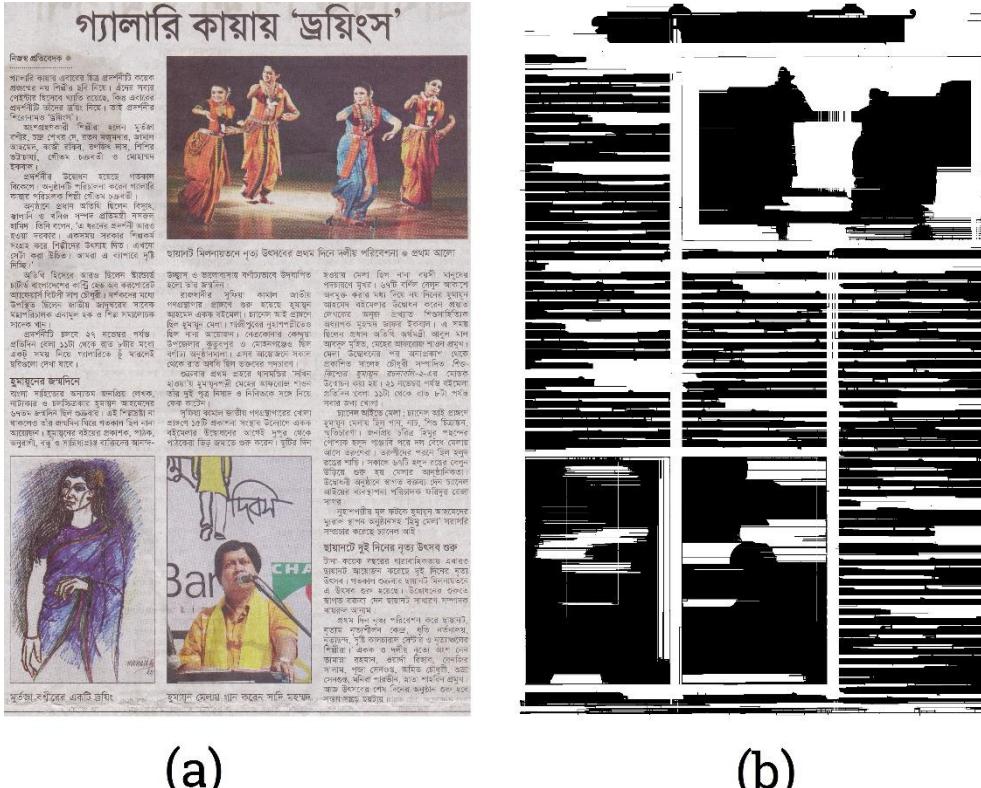


Figure 7: Newspaper Image Decomposition: (a) Original, (b) Segmented Image

Then for the recognition of the text blocks, they have taken 4 threshold values.

$$T_s = \begin{cases} 0, & \text{if } H_s(x,y) < t_1 \\ 1, & \text{if } H_s(x,y) \geq t_1 \wedge H_s(x,y) < t_2 \\ 2, & \text{if } H_s(x,y) \geq t_2 \wedge H_s(x,y) < t_3 \\ 3, & H_s(x,y) \geq t_3 \end{cases} \quad (1)$$

[9]

Here  $t_1$ ,  $t_2$  and  $t_3$  are respectively minimum expected letter height of text segment, a title and a title region. Depending on those values they decided whether a text is from headline, sub headline or columns.

#### 4.2.2 Decomposing using a Split and Merge Approach

On January 2001, another conference was held in Seattle of USA on the banner of ICDAR. In this Conference Karim Hadjar, Oliver Hitz and Rolf Ingold from DIUF, University of Fribourg presented their work on newspaper documentation. They proposed decomposing the documents on using the Split and Merge Approach [13]. Their research mainly follows the work described below:

In their research work, they proposed a system for daily paper page division in view of part and consolidating zones. They slice the info picture as per the even and vertical line. At that point to distinguish pictures they ascertained associated part. Each associated component is depicted by its rectangular jumping box. [13] For each letter, they got isolated rectangular box. They picked a limit estimation and if a rectangular box is greater than this esteem then that was chosen to be picture. Further they concluded that in a picture, the bouncing boxes of the encased associated parts are probably going to cover, whereas in writings areas there are not very much covering it.



২০০৮ সালের ফাইল ছবিতে প্রাচীন পালমিরা  
থিয়েটারে একদল পর্যটক ● রয়টার্স

## প্রাচীন শহর পালমিরা

পালমিরা শহরের প্রত্ন কবে হয়েছিল, তার নিশ্চিত সঠিক সময় কারও জানা নেই। তবে খ্রিস্টপূর্ব প্রথম শতক থেকে প্রবর্তী ৪০০ বছর এ মরুভূমি শহরটি উন্নতির শিখরে পৌঁছেছিল বলে ঐতিহাসিকেরা মনে করেন। মসলা, সুগন্ধি, রেশম ও হাতির দাঁতের শিল্পের ওপর ভিত্তি করে এক সম্পদশালী শহরে পরিণত হয় পালমিরা।

‘মরুর মুক্তি’ নামে পরিচিত এ শহরে আছে দুই হাজার বছরেরও বেশি সময়ের পুরোনো মন্দির, কারুকার্যমণ্ডিত স্তুতিশুভ নানা নির্দশন। জাতিসংঘের শিক্ষা, বিজ্ঞান ও সংস্কৃতি সংস্থা ইউনেস্কো ১৯৮০ সালে পালমিরাকে বিশ্ব ঐতিহ্য ঘোষণা করে। ২০১১ সালে সিরিয়ায় গৃহযুদ্ধ শুরু হওয়ার আগে প্রতিবছর গড়ে দেড় লাখ পর্যটক পালমিরা ভ্রমণ করতেন।

Figure 8: In a content locale, there are no overlapping boxes of associated components.

To distinguish between text and title blocks. Text blocks are blocks, which contain text that is smaller or equal to the dominant letter height of the whole image. Blocks with text greater than the dominant letter height are title blocks.

# রাষ্ট্রদ্বারে মামলায়

## খালেদাকে সমন

আদালত প্রতিবেদক >

রাষ্ট্রদ্বারে অভিযোগে করা একটি মামলা আমলে নিয়ে বিএনপি চেয়ারপারসন খালেদা জিয়াকে আগামী ৩ মার্চ আদালতে হাজির হওয়ার জন্য সমন ভরি করা হচ্ছে। ঢাকার মেট্রোপলিটন ম্যাজিস্ট্রেট রাশেদ তালুকদার গতকাল সোমবার এ আদেশ দেন। মুক্তিযুক্ত শহীদের সংখ্যা নিয়ে সংশয় প্রকাশ এবং বঙ্গবন্ধুর ভূমিকা নিয়ে কটক্ষ করার আদালতে আজ মামলাটি দায়ের করেন। বাদার জবাবদি প্রাপ্ত করে শুনানি শেষে সব অভিযোগ আমলে নিয়ে সমন ভরি করেন বিচারক। একই সঙ্গে হাজিরার দিন ধার্য করেন।

মামলায় দণ্ডবিধির ১২৩ (ক) ধারায় বাংলাদেশ সৃষ্টির নিম্না ও তার সর্বভৌমত বিলাপে সমর্থন করার এবং ১২৪(ক) ধারায় রাষ্ট্রদ্বারের অভিযোগ আনা হচ্ছে। পাশাপাশি জনগণের অনিষ্টসাধন সহায়ক বিবৃতি দেওয়ার অভিযোগ আন হচ্ছে। ৫০৫ ধারায় ১২৪(ক) ধারায় অভিযোগ প্রমাণিত হলে সামোচরিত কার্যান্বয় হচ্ছে পারে।

- মুক্তিযুক্ত শহীদের সংখ্যা নিয়ে সংশয় প্রকাশ এবং বঙ্গবন্ধুর ভূমিকা নিয়ে কটক্ষ করার এই মামলা
- অভিযোগ প্রমাণিত হলে সর্বোচ্চ শাস্তি যাবজ্জীবন



না করে বঙ্গবন্ধুর প্রতি ইঙ্গিত করে খালেদা জিয়া বলেন, ‘তানি স্বাধীনতা চাননি, তিনি পাকিস্তানের প্রধানমন্ত্রী হতে চেয়েছিলেন।’ আরজিতে বলা হয়, ওই বজ্রবা পরের দিন বিভিন্ন পত্রিকায় সর্বিদ্বারে প্রকাশিত হয়। খালেদার হেই বজ্রবা ১৯৭৫ সালের ২৬ মার্চের বাংলাদেশের স্বাধীনতা সংগ্রামের মূল ঘোষণার শুধু পরিপন্থী নয়, অবজ্ঞাপূর্ণ বটে। তার একপ বজ্রবা মহান মুক্তিযুদ্ধের মাধ্যমে প্রতিষ্ঠিত

কেটের আগল বিভাগের আইনজীবি ড. মোমতাজ উদ্দিন মেহেন্দী। আওয়ামী লীগের কার্যনির্বাহী কমিটির সদস্য ড. মেহেন্দী দণ্ডবিধির ১২৩(ক)/১২৪(ক)/৫০৫ ধারার অভিযোগ মামলাটি দায়ের করেন। বাদার জবাবদি প্রাপ্ত করে শুনানি শেষে সব অভিযোগ আমলে নিয়ে সমন ভরি করেন বিচারক। একই সঙ্গে হাজিরার দিন ধার্য করেন। মামলায় বাদীগকে শুনানি করেন নামেক ব্রাষ্টমন্ত্রী অ্যাডভোকেট সাহেব খাতুন এন্পি, ঢাকা জেলা পাবলিক প্রসিকিউটর খেলাদার আবেদুল মামান, ঢাকা মহানগর পাবলিক প্রসিকিউটর আবদুল্লাহ আব, বিশেষ পিপি ফেরকান মিয়া প্রাপ্ত। মামলার আরজিতে একটি তাবোচা সত্য খালেদা জিয়া মুক্তিযুক্ত শহীদের সংখ্যা নিয়ে সংশয় প্রকাশ করে বলেন, ‘স্বাধীনতাযুক্ত শহীদের সংখ্যা নিয়ে বিতর্ক আছে। আজকে বলা হয় এত সক্ষেত্রে শহীদ হয়েছে। কিন্তু প্রক্রিয়াক কর্তজন লোক শহীদ হয়েছেন তা নিয়ে ‘বিতর্ক আছে।’ বিএনপি নেটো আরো বলেন, আওয়ামী লীগে কোনো মুক্তিযোক্তা নেই, যারা তাৎক্ষে তারা সর্বাই ভুয়া। বঙ্গবন্ধু শেখ মুজিবুর রহমানের নাম ডেখে

► পৃষ্ঠা ১৩ ক. ৬

Figure 9: Splitting Newspaper image into small zones.

To achieve the labeling, their algorithm computes the height of all lines inside the blocks to classify as title or text. They have said that their proposed methods will work flawlessly if and only if every input image was noiseless and was not skewed.



Figure 10: A part of correctly segmented. Each part is each labeled with different color.

At the end they proposed that the result they obtained is an automated newspaper decomposition but can create problems in some other cases.

#### 4.2.3 Automated Text-Line Height Calculation for Font Size Detection

In this segment, we have to study for the crucial part of the document reconstruction. While reconstructing the document, we have to work it size of the image used in the document, different font size for labeling headline, sub-headline and paragraph. OCR intelligence needs to detect the font size in text documents in order to understand the document image and extract knowledge from the document texts [14] [15].

Working with the scanned documents is doing the processing with an uncompressed image document. D. S. B. F. R. Chen at el. [15] Proposed a technique which use the Line Segmentation process in order to detect the font size.

On the other hand, compressed document gives us the high performance and storage adaptability, so that machines like fax machines [16], Photocopier machines and many others uses the compressed form of the documents [17]. If the documents are like run length compressed TIFF documents detection of the font size come handy [17] in terms of complexity like time and space [18], [19], [20]. Mohammed Javed at el. [17] proposed an approach to automate the font size detection straight from the run length compressed TIFF documents by working with the binary documents. These methodology worked without doing any decompression of the document [17].

## Chapter 5

# Research on Bangla Document Decomposition

In the year of 1994 the work on OCR was done by U. Pal and B.B. Chaudhuri in title associated with “OCR in Bangla: an Indo-Bangladeshi Language” [21]. In later they again worked with the title “A complete printed Bangla OCR system” [22]. But both of their research work was mainly based on recognition system. And many of the published research work also focused just for the recognition system.

In our study we find out that there is lack of digitize Bangla document, many research like Machine Translation, Text Mining work in Bangla didn’t advanced and for many special cases the decomposition process done in other language. It is a matter of fact that as Bangla letters are not separated from each other. Therefore, in Bangla newspaper rectangular boxes bounding texts are likely to overlap. This is the main reason documents in Bangla needs to be decomposed uniquely.

### 5.1 Research on Bangla Document Decomposition in SUST

At Shahjalal University of Science and Technology, Sylhet Department of Computer Science and Technology, development of Bangla OCR has been started since 1998. But they work of decomposition mainly started from the Batch 2006. Batch from 2006 and 2007 has worked on headline and column detection. Batch from 2008 Syed Rezwanul Haque Rubel, Arifuzzaman Sohel started working with “Tesseract” an OCR engine developed by Hewlett Packard Lab and sponsored by Google. [23]. From the Batch 2010 which was done by Fahad Hasan and Tasmin Afroz Tanee, They borrowed some of the steps research work of the document decomposition and added extra rules which was set by them.

## 5.2 Decomposition Research from Batch 2006

This batch proposed their methodology by setting a Fixed Space Width, which means the minimum space between headlines and news body. Then fixed space height which means the minimum space between columns. From their proposal at first Space width set as 0 and starting point is  $(0, 0)$ . At starting point it check the pixel vertically and find the white line. If this process find white line then it increase Space Width. At the point  $(x, y)$  where Space Width is greater than Fixed Space Width, they cut the image in a rectangular shape from  $(x, y)$  to starting point. Then the process save this image as a sub image and set starting point as  $(x, y)$ . If it found a line, which is not fully white, then they set Space Width to 0 and continued this process for whole image. If number of sub images from previous step is greater than 1, they considered the first sub image as headline image and other sub images as input for last step. They load each sub images for checking the pixel horizontally. To do so, set Space Height as 0 and starting point  $(0, 0)$ . If we find the white line, we increase Space Height. At the point  $(x_1, y_1)$  where Space Height is greater than Fixed Space Height, we cut the image in a rectangular shape from stating point to  $(x_1, y_1)$ . After that, they saved this image as a final child image and set starting point as  $(x_1, y_1)$ . If we find a line which is not fully white, we set Space Height as 0 and continue this process for whole image.



Figure 11: Input and Output for using the work of 2006 Batch.

We can clearly see that their process is just separating the headline and the column portion of a newspaper document. But from the criteria of recognition part images should be removed must.

### 5.3 Decomposition Research from Batch 2007

Batch from 2007 followed the work from the batch 2006 but with setting some different rules to adjust the output. They took the width of first couple of text lines from the top of the image. It compared the widths of this line. If there is a rapid decrease in the width then they decided that headline ends on the top of current line. Afterward to detect columns, they proposed a system to count the number of black pixels and keep the result sequentially. The result would be a sequence of zero and non-zero value. Zero sequence is deduced to be column and nonzero sequence are column space.

It is a matter to notice that to work perfectly with this method, documents supposed to have the headline on the top of the article. So this approach will work for paper images like the below one. Because, here the headline is on the top.



Figure 12: Input and Output for work from 2007 Batch.

Many the time it can happen that the headline or sub headline is not on the top of the article but in betwixt the columns.



Figure 13: Sub Headline between Columns

Moreover, most importantly, newspaper page contains a lot of images along with the texts.

Sometimes images are at the top of the columns, and sometimes they are in between two consecutive columns.

# যশোরের পাখির বাসা যায় ইউরোপে

**মনিরুল ইসলাম, যশোর ০১**

যশোরে তৈরি শৈবিন পাখির বাসা ইউরোপের হয় দেশে রঙনি হচ্ছ। একে ঘোরে কোটি টাকার সম্পর্কিমাণ বেদনের মূল্য হচ্ছে দেশ। এর মাধ্যমে স্থানীয় হচ্ছে এ অর্থনৈতিক জাতোরে নিম্ন আয়ের মানুষ।

ইউরোপের জার্মানি, ফ্রান্স, বেলজিয়াম, নেদারল্যান্ডস, স্পেন ও পর্তুগালে বিভিন্ন শহরে সৌন্ধৰ পাখি উৎপন্নের খামারে যাচ্ছে ৪০ ধরনের পাখির বাসা। তবু পুঁজির সংকটের কারণে এ পণ্যের রঞ্জন বাণিজ্যে যাত্তা প্রসার ঘটার কথা ছিল, ততটা ঘটেনি বলে দারি উদ্বোধনে।

এ পণ্য রঙনির অন্যতম উদ্বোধন চাকার সিদ্ধিষ্ঠৰী একান্তর আয়ের মানুষের মাধ্যমে পাখির বাসা তৈরি করেন। তিনি যশোরের কয়েকটি গ্রামের নিম্ন আয়ের মানুষের মাধ্যমে পাখির বাসা তৈরি করেন। পরে এগুলো রঙনি উপকারণে করে নোপরে দিবেশে পাঠান।

জনতে চালিলে খায়রিল আলম বলেন, 'যশোর সদর উত্তোলন আবাদ করুয়া, সীতারামপুর ও বাহাদুরগ়ুপ্তের গ্রামে তৈরি শৈবিন পাখির বাসা ইউরোপের বাজারের বাজারের বাপাক চাহিদা। অতভে ৪০ ধরনের পাখির বাসা ছয়টি দেশে এখন রঞ্জনি করা হচ্ছ। বর্তমানে বছরের প্রায় ১ কোটি টাকা মূল্যের পাখির বাসা রঙনি হচ্ছে।'

সন্তুষ্টি আবাদ করুয়া ও সীতারামপুর গ্রামে শিয়ে দেখা গেছে, ঘরের ঘরে পাখির বাসা তৈরির কাজ চলছে। ঘরের বারান্দা ও অভিন্নায় বাসা নাহি—পুরুষ মিলে পাখির বাসা বুননের কাজ করছেন। পুরুষেরা বাবের চাটাই দিয়ে বুননের মূল উপকরণ তৈরি করছেন। নারীরা বাসা তৈরির জো (কাজের পৃষ্ঠ) রূপচৰ্ছন। আরেকজন বাসা তৈরির কাজ শেষ করছেন।

আবাদ করুয়া গ্রামের সম্পর্কের কিনারাম দাসের বাড়ির উচ্চানে শিয়ে দেখা গেল, তিনি তলতা বাসের চাটাই চিরে না দিয়ে সুচারু করছেন। পশে বসে তার শ্রী কাঞ্জন দাস ও মা রঞ্জিতা

জার্মানি, ফ্রান্স, বেলজিয়াম, নেদারল্যান্ডস, স্পেন ও পর্তুগালে যাচ্ছে ৪০ ধরনের পাখির বাসা ,



কাজ শিখেছি। এখন এ কাজ করেই 'জীবিত নিবাহ করি'।

চৰীদাস বলেন, 'ওৱল নিকে ঢাকাই সেলিম সাবে একটি নবনি বাসা এনে বলেন, এ ধরনের পাখির বাসা তৈরি করতে পারো কি না দেখো। সেনে সন্দে আমরা কয়েকজন ওই রঞ্জন বাসা তৈরি করে দিলাম। সেই থেকে সেলিম সাহেব আমাদের পাখির বাসা তৈরির কাজ দে। প্রথমে আমরা ৩০ জনকে বাসা তৈরির প্রশিক্ষণ দিই। সেই থেকে কাজ করে যাচ্ছি।'

সুব্রহ্মণ্য বলেন, 'আম থেকে পাখির বাসা সংস্থান করে এ কেবলে এনে শুয়ে রোদে শুলিয়ে রঙনি বাসার প্রক্রিয়া করে প্রাক্কেট করা হয়। প্রথম কাটিন করে প্রাক্কেট মাধ্যমে ঢাকার অভিযন্তে পাঠানো হয়। সেখান থেকে বিদেশে যাব।'

এ শিল্পের ফুল উদ্বোধন খায়রিল আবাদ চাকা বিশ্ববিদ্যালয় থেকে আন্তর্জাতিক সম্পর্ক বিষয়ে স্নাতকোত্তর শেখে করে চাকারিয়ে চোটা করেননি। উদ্বোধন হওয়ার হিসেবে তিনি পাখির বাসা উৎপক্রিয় বিদেশে রঞ্জিত করেন।

খায়রিল আলম বলেন, '৩০ বছর আগে যশোরের বেশৰকাৰি সংস্থা 'বাঁচতে শেখা'ৰ নিৰ্বাচিত পৰিচালক আঙেলা গোমোজের মাধ্যমে যথোচ্চ গিয়ে পাখির বাসা তৈরি কৰেছিলাম। এখনো সে কাজ চলছে।' তিনি বলেন, 'বিদেশ ভিত্তি বাঁচতে বাঁচাবে এ পোষা পাখির চাষ হয়। আমাদের দেশের তৈরি পাখির বাসায় ওই পাখি ডিম পাতে, বাচা ফুটায়। যে করাগে বিদেশে এ বাসা অনেক চাইদান। আগে চীজের দখলে হিল ইউরোপের পাখির বাসা তৈরির কাজ শুরু করেছিলেন।'

সীতারামপুর গ্রামে চৰীদাসের বাড়িতে শিয়ে দেখা গেল, চৰীদাস ও তাঁর মেয়ের জীবাই আৱনত বিশ্বাস নারকেলে ও পাট দিয়ে বাবুৰ পাখির বাসা যাতা বাসা তৈরি করছেন। আৱনত বিশ্বাস বলেন, 'উচ্চমাধ্যমিক পাস করে অন্য চাকারিতে যাইনি। হোটেবেলো থেকে পাখির বাসা তৈরির

Figure 14: A Complex Document.

24

## **5.4 Decomposition Research from Batch 2010**

As there was many demerits on decomposing the documents, Batch of 2010 started to work from scratch and set a methodology to work on a complex document structure. The working proposal they gave is like below:

- Edge Detection.
- Conversion of input image.
- Finding the blank spaced area (element separators) among several elements.
- Finding components adjacent to the separators and segmenting them.
- Using these portions to detect which category they belong to.

So our findings is that, they mainly focused on the work of B. Gatos at el. [9] And Karim Hadjar at el. [13] In this process the result of the input document images started to give a very good result. So this methodology was quite satisfactory. They used some threshold filter to detect the Images, Headlines and columns from the component they separated.

# এক যুবকের চেখে শত মানুষের কান্না



নিজস্ব প্রতিবেদক ॥

শুকে “নাজাবিদার চাই” প্লাকার্ট বালেখা নিয়ে প্রধানমন্ত্রীর বসন্তবাসের কাছে পাঠিয়ে থাকা অসম কলাম আজাদের দরিদ্র শেষ শত মানুষের নাজাবিদার পাঞ্চাংগের আকাঙ্ক্ষার প্রকাশ ঘটিয়েছে। রাজ্যে একা মাড়িয়ে পড়া এই যুবক কর্তৃত চোখে করে নিয়ে এলেক্ষণেন নারায়ণগঞ্জের সেনানীর উপজেলার নামাচী আবুল কালাম আজাদের শত মানুষের কান্না।

শেষ পর্যন্ত আজাদের প্রতিবান কানে ভুলেয়ে গেছে। গতকাল তাঁকে তেকে নিয়ে হাস্তী প্রশংসন। প্রধানমন্ত্রী শেখ হাসিনার মৃত্যু সংবিধি অসম কলাম আজাদ তাঁর কথা দেখেছেন।

এখনে তাঁর কানে “তুম্হ সমস্যা” নিয়ে আজাদ কেন প্রধানমন্ত্রীর মৃত্যুর কানা নাড়ার সাহস দেখিয়েছেন, তা খুঁজতে গিয়ে বেরিয়ে এসেছে অকর্তৃ সব অভিযোগ। উপজেলার সামুদ্র ইউনিয়নের পক্ষবিধাটি বাজারের পাশে নামাচী মৌজায় মেঠো পার্মেটিস আঙ্ক ডাইং নামের একটি কারখানা করেছেন নারায়ণগঞ্জের বাবস্থার অসম পেঁচার। অভিযোগ উঠেছে, তিনি তই কারখানার জন্য এক বছৰ ধরে এলাকাক নিরীহ



মানুষের ফসলি ভাই জোর করে বিদেশ নিয়েছেন। বিহু করতে রাজি না হলে বাল দেখে দখল করেছেন। প্রামাণ্য অভিযোগ করেছেন, ইতিমধ্যে ২১ জন বাসিন্দার ৩৪ বিহু জমি দখল করা হয়েছে। ৩২ জনকে কিন্তু উকা হাতে পরিদেয়ে জমি নিয়ে এলাকা থেকে বিনার করে দেওয়া হয়েছে। এখনে অসমকে জমি, বসন্তবান্ডি হেতু তলে যাওয়া বাস্তিক সেওয়া হচ্ছে।

কৃতজ্ঞকৌশলের অভিযোগ, এ কাজে করখন্দার মালিককে সহায়তা করছেন হাস্তী আবুল কালামের অফসারগুলির কিন্তু সেতা-কর্মী এবং রাবের কর্মকর্তা সময়। তাঁর বিহু করতে রাজি হলেন এমন করখন্দারী পরিবারের সদস্যকে রাব হোজার করে মাসক তাঁকা ও বিহুর অভিযোগ এসে আসলা করেছে। হাস্তী প্রধানমন্ত্রীর কানে গিয়ে কৃতজ্ঞকৌশল কোনো অভিকার প্রস্তু। শেষ পর্যন্ত গত বোবৰুর আজাদ নামের যুবকটি এসে নাড়িয়েছেন প্রধানমন্ত্রীর মৃত্যু। সেটা করতে নিয়েও হাবোনির শিকার হয়েছে আজাদের পরিবার।

এরপর পৃষ্ঠা ২ কলাম ৫

Figure 15: An Input Image we used in 2010 Batch Methodology.

## Image Detection:

To decide whether a block is image or not they have proposed three filters.

It checks whether the length and width of each block is greater than a threshold value or not. Here the threshold value is calculated as 100. If it finds out the pixel ratio of the blocks. If the ratio is greater than 1500 then they anticipate that this block may be an image. Finally they used to determine the ratio of the sum of the total length of each block and total width of each block. If a block is image then the ratio will be greater than 1.25. If any block satisfies all these three filters then they are decided to be an image.



Figure 16: Image detected from Document using 2010 Batch methodology

#### **Headline Detection:**

For filtering out the Headlines they set up three threshold rules:

Here they check the height and width of a block. In addition, if the process have found out that for being a headline a block should have a height should be in between 30 to 99. Here they find out the horizontal histogram of first 30% pixels towards the height of any box. If, the highest value of the horizontal histogram equal or greater than 80% of the width of the block, as well as, the horizontal histogram of last 5 pixels of the box will be 10% less than the width of the box. If any block fulfills these two characteristics then they have decided them to be a headline.

# এক যুবকের চেখে শত মানুষের কান্না

Figure 17: Headline detected from Document using 2010 Batch methodology

## Column Detection:

For A newspaper contains headlines, images, columns mainly. Therefore, they marked all the other blocks as columns.

বুকে 'ন্যায়বিচার চাই' প্রাকার্ত বা  
লেখা নিয়ে প্রধানমন্ত্রীর  
বাসভবনের কাছে দাঁড়িয়ে থাকা  
আবৃল কালাম আজানের নাবি শেষ  
পর্যন্ত শত মানুষের ন্যায়বিচার  
প্রাণ্যার আকাঙ্ক্ষার প্রকাশ  
যাইয়েছে। রাজ্য একা দাঁড়িয়ে  
পড়া এই যুবক কার্যত চেখে করে  
নিয়ে এসেছিলেন নারায়ণগঙ্গের  
সোনারগী উপজেলার মানুষী  
হ্যামের শত মানুষের কান্না।      আবৃল কালা

শেষ পর্যন্ত আজানের প্রতিবাদ কানে তুলেছে  
রাষ্ট্র। গতকাল তাঁকে ভেকে নিয়েছে ছানীয়  
প্রশাসন। প্রধানমন্ত্রী শেখ হাসিনার দুখ্য সচিব  
আবৃল কালাম আজান তাঁর কথা খুঁজতে দেয়েছেন।

গ্রামে জমিজমার 'ত্রুট সমস্যা' নিয়ে আজান  
কেন প্রধানমন্ত্রীর দুয়ারে কড়া নাড়ার সাহস  
দেখিয়েছেন, তা খুঁজতে গিয়ে বেরিয়ে এসেছে  
অর্জনের সব অভিযোগ। উপজেলার সদিপুর  
ইউনিয়নের পক্ষবীঘাট বাজারের পাশে মানুষী  
মৌজায় মেট্রো গার্মেন্টস অ্যান্ড ডাইং নামের  
একটি কারখানা করছেন নারায়ণগঙ্গের বাবসাহী  
অঞ্চল পোকার। অভিযোগ উঠেছে, তিনি ওই  
কারখানার জন্য এক বছর ধরে এলাকার নিরীহ

মানুষের ফসলি জমি জোর করে  
কিনে নিচ্ছেন। বিক্রি করতে রাজি  
না হলে বাল দেলে দখল করেছেন।  
গ্রামবাসী অভিযোগ করেছেন,  
ইতিমধ্যে ২১ জন বাসিন্দার ৩৪  
বিহু জমি দখল করা হয়েছে। ৩৫  
জনকে কিছু টাকা হাতে ধরিয়ে জমি  
নিয়ে এলাকা থেকে বিদায় করে  
দেওয়া হয়েছে। এখনো অনেকেরে  
জমি, বসতবাড়ি হেবে তুলে  
যাওয়ার দুর্ভিক দেওয়া হচ্ছে।

তুর্কভোগীদের অভিযোগ, এ<sup>১</sup>  
কাজে করখানার মালিককে  
সহায়তা করছেন ছানীয় আশ্রয়ী লীগের  
অঙ্গসংগঠনের কিছু নেতৃ-কর্মী এবং রায়বের  
কয়েকজন সদস্য। জমি বিক্রি করতে রাজি  
হননি এমন কয়েকটি পরিবারের সদস্যকে রায়ব  
যোগার করে মানক রাখা ও বিক্রির অভিযোগ  
এসে মামলা করেছে। ছানীয় প্রশাসনের কাছে  
গিয়ে তুর্কভোগীরা কেনো প্রতিকার পাননি।  
শেষ পর্যন্ত গত রোবৰার আজান নামের যুবকটি  
এসে দাঁড়িয়েছেন প্রধানমন্ত্রীর দুয়ারে। সেটা  
করতে গিয়েও হয়েরানির শিকার হয়েছে  
আজানের পরিবার।

Figure 18: Column detected from Document using 2010 Batch methodology.

## **Chapter 6**

# **Findings on Previous Methodology**

The Bangla Document Decomposition work done by the Batch 2010 from the Department of Computer Science and Engineering of Shahjalal University of Science and Technology was the starting point for our research work. They started their work from the scratch and build a new system. So we look into the deep of this system, tried to find the errors how make a decision how to resolve them.

### **6.1 Findings on 2011 batch Methodology**

Our first finding was that they worked on the document “e-Prothom Alo” and “e-Jonokontho”. Which is a big fact on further research work. Our goal is to create a decomposition system where it have work with scanned documents. Now by studying scanned document we find out that, there is will be always some noise in it by default. And for aging, the ratio of the noise increases exponentially. This high rate of noise creates a huge impact on finding out the connected component on segmentation step. It will make the whole document one connected component and the whole segmentation process will fail on the start.



Figure 19: A Scanned Image is zoomed to preview the noise.

As we said previously that, the images are not scanned directly but downloaded from the e-prothom-alo site. These image documents are in 96dpi. Which is a really bad input for the recognition part. Each of the font will be barely recognized. Also from the binarized image we can see that a lot of bits are missing from the lines of the columns due to low resolution. These missing bits will create a very big factor for the recognition part. So for a perfect recognition with a higher rate, higher dpi images needed to be used. So we used higher dpi scanned pictures to check their methodology performances. For higher resolution images, the methodology fails to segment the documents.

# কেন্দ্ৰীয় ব্যাংকেৱ ৭-৮ কৰ্মকৰ্তাৰ দায়িত্বে অবহেলা

## বিশেষ প্রতিনিধি

বিজৰ্ণ চুৱিৰ বটনায় বাংলাদেশ  
ব্যাংকেৱ সাত. আটজন কৰ্মকৰ্তাৰ  
সদস্যেৰ ভালিকাৰ বেগৰহে  
পৰিষেৱ অপৱধি তন্তৰ বিভাগ  
(সিইআইডি) এন্দৰ কৰ্মকৰ্তাৰ  
বাস্তুগত বাংক হিসাব, কেইল  
অসাম, প্ৰদল ও মেৰহিল ফণেৰ  
কাথেপকখন্ত তথ্য খন্তভৰ দেখা  
হিস্তে তৈদেৱ ৬৫৩৩ সকা঳  
জিঞ্জসামাদও কৰেছে সিইআইডি

তন্তৰকাৰী সংষ্ঠিৰ কৰ্মকৰ্তাৰ  
কৰেশা হৰন্তৰ ক'ৰণ নাম আনতে  
হ'লি হ'লি হ'লি হ'লি কেন্দ্ৰীয়  
ব্যাংকেৱ কোন পৰ্যাপ্তৰ কৰ্মকৰ্তা—  
সে হ'লি হ'লি নিতে সহ না সিইআইডি  
সিইআইডিৰ অভিযোগ উপৰিটি



বিজাতীয়  
অর্থ চুৱি

“সন্দেহভাজনেৰ  
ভালিকাৰ আমৰা  
ছেটি কৰে এনেছি

শাহ আলম  
অভিযোগ ডিইআইজি, সিইআইডি

বিভিন্ন সময়া চৰক্ষিক কৰ্মকৰ্তাৰ  
বাস্তুগত ও নাটৰিক কৰ্মকৰ্তাৰ জন্ম  
কৰেছে সিইআইডি অসাম  
কৰ্মকৰ্তাৰ জন্ম কৰা হয় কেলোৱা

Figure 20: 96dpi image Binarized

Because some important threshold values they determined for finding out the connected components only works with the 96 dpi image. Higher DPI like 300, it failed to segment properly.

Another problem is on determining the line height. The method they use to count the line height was to start a track if a black bit is found in a row and then start find is there any row which has no black bit and the height from the tracking point must be greater than 5. But if any kind of overlapping case arise this process fails to find out the line height. A more than 1 line height will be shown as a line height.

## খায়রুল আনাম।

প্রথম দিন নৃত্য পরিবেশন করে ছায়ানট,  
নৃত্যম নৃত্যশীলন কেন্দ্র, ধূতি নর্তনালয়,  
নৃত্যছন্দ, সৃষ্টি কালচারাল সেন্টার ও নৃত্যঞ্চলের  
শিল্পীরা। একক ও দলীয় নৃত্যে অংশ নেন  
তামানা রহমান, ওয়ার্দা রিহাব, বেনজির  
সালাম, পূজা সেনগুপ্ত, অমিত চৌধুরী, শুভা  
সেনগুপ্ত, মনিরা পারভীন, স্বাতা শাহরিন প্রমুখ।

Figure 21: Overlapping case are highlighted in red area.

On next point we find out that it fails to detect a table components. The image filtering case checks if it has a border or not. If anything with border is found as a component it labels it as an image. So tables with borders are detected as an image.

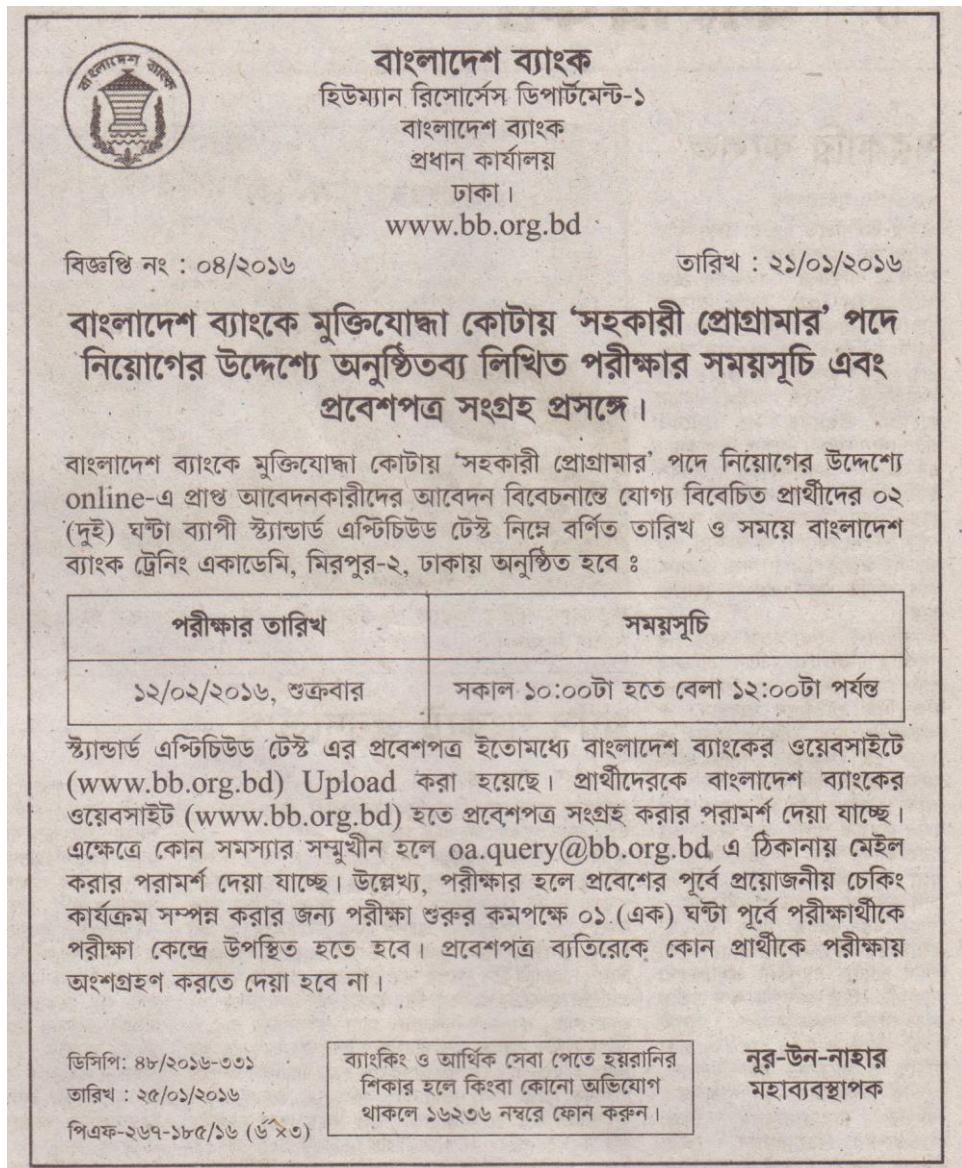


Figure 22: A Sample of a Table Component

The last error we find out that, there is no case for handling the skew case. During scanning an image documents, sometimes there are some skew. These skewness is

a problem both for segmentation and recognition. So we need to DE Skew the image before implementing the segmentation process.

The very last thing we will discuss in this topic is about reconstructing the document. We find out that, there was not any implementation or research was done for reconstructing the document. So we also tried to work on this topic to resolve the issue.

## Chapter 7

# Current Methodology: Part 1

Rather than start over the work from the scratch, we took the previous work and changed some set of rules. Also we took the 300dpi as a standard scanned document input for research work o document decomposition. In this section we discussed about the work we have done for our course CSE 400. On later part of this chapter we discussed what we did to improve the performance of our previous work.

### 7.1 Noise Reduction

At first, the whole image is taken as a form of matrix. In order to create a perfect result, the noises should be removed from the scanned document. In this process Bilateral Filter Method is used. Gaussian Blur also removes noise pixel by smoothing the image, but in Bilateral Filter it preserves the edges and smooth out the noise pixels.



Figure 23: Input Image for Document Decomposition

# ଗ୍ୟାଲାରି କାନ୍ଯା ଡ୍ରାଇସ

ନିଜ୍ୟ ପ୍ରତିବେଦକ

ଗ୍ୟାଲାରି କାନ୍ଯାର ଏବାରେର ଚିତ୍ର ପ୍ରଦଶ୍ନିଟି କହେକ ପ୍ରତିବେଦନ ନ ଶିଳ୍ପୀର ଛବି ନିଯେ । ଏବେର ସବାର ପେଇଟିର ହିମେବେ ଥାତି ରମେହେ, କିନ୍ତୁ ଏବାରେର ପ୍ରଦଶ୍ନିଟି ତାନେର ଡ୍ରାଇସ ନିଯେ । ତାହିଁ ପ୍ରଦଶ୍ନିର ଶିରୋନାମ ଓ ଡ୍ରାଇସ ।

ଅଂଶପ୍ରାତିକରି ଶିଳ୍ପୀର ହିମେନ ମୁଖ୍ୟା ବଶୀର, ଚନ୍ଦ୍ର ପେଇର ଦେ, ରତନ ମଜୁମଦାର, ଜୀମାଲ ଆହମେଦ, କାଞ୍ଜି ରକିବ, ରଙ୍ଗଜିଂ ଦାସ, ଶିଶିର ଭାବାନ୍ଧୀ, ପୌତମ ଚକ୍ରବର୍ତ୍ତୀ ଓ ମୋହାମଦ ଇକବାଲ ।

ପ୍ରଦଶ୍ନିର ଉତ୍ସୋଧନ ହେଁବେ ଗତକାଳ ବିକଳେ । ଅନୁଷ୍ଠାନଟି ପରିଚାଳନା କରେଲା ଗ୍ୟାଲାରି କାନ୍ଯାର ପରିଚାଳକ ଶିଳ୍ପୀ ପୋତମ ଚକ୍ରବର୍ତ୍ତୀ ।

ଅନୁଷ୍ଠାନେ ପ୍ରଧାନ ଅଭିଧି ଛିଲେନ ବିଶ୍ୱାସ, ଆଲାନି ଓ ଖଣ୍ଜ ସମ୍ପଦ ପ୍ରତିମାରୀ ନମ୍ବରଳ ହମିନ । ତିନି ବିଲେନ, ‘ଏ ଧରନେର ପ୍ରଦଶ୍ନି ଆରା ହେଁବା ଦରକାର । ଏକମଧ୍ୟ ସରକାର ଶିଳ୍ପକର୍ମ ସଂଗ୍ରହ କରେ ଶିଳ୍ପୀରେ ଉତ୍ସାହ ଦିତ । ଏଥିନେ ସେଟା କରା ଉଚିତ । ଆମରା ଏ ବ୍ୟାପାରେ ଦୁଇ ଦିଛି ।’

ଅତିଥି ହିମେବେ ଆରା ଛିଲେନ ଟ୍ରୋଡାର୍ଡ ଚାର୍ଟର୍ ବାଲୋଦେଶର କାନ୍ତି ହେତୁ ଅବ କରିପୋରେଟ ଆଫ୍ରିକାର ବିଟପୀ ଦାଶ ଚୌଥୀମୀ । ଦଶକଦେଶ ମଧ୍ୟେ ଉପସ୍ଥିତ ଛିଲେନ ଜାତୀୟ ଜାଦୁଯରେର ସାବେକ ମହାପରିଚାଳକ ଏନାମୁଳ ହକ ଓ ଶିଳ୍ପ ସାମୋଳେଟକ ନାମକେ ଥାଇ ।

ପ୍ରଦଶ୍ନିଟି ଚଲବେ ୨୭ ନଭେମ୍ବର ପର୍ଯ୍ୟନ୍ତ । ପ୍ରତିଦିନ ବେଳା ୧୧୮ ଥେକେ ରାତ ୮୩୮ ମଧ୍ୟେ ଏକଟୁ ସମୟ ନିଯେ ଗ୍ୟାଲାରିତେ ଟୁ ମାରଲେଇ ଛିଲିଗୁଣେ ଦେଖା ଯାବେ ।

## ହୁମାୟନେର ଜୟଦିନ

ବାଲୋ ସାହିତ୍ୟର ଅନ୍ତର୍ମାନ ଜନପିଲ୍ ଲେଖକ, ନାଟ୍କାର ଓ ଚଲଚିତ୍ରକାର ହୁମାୟନ ଆହମେଦର ୬୭ତମ ଜୟଦିନ ଛିଲ ଉତ୍ସାହର । ଏହି ଶିଳ୍ପାନ୍ତରୀ ନାଥକଳେ ଓ ତାର ଜୟଦିନ ଧିରେ ଗତକାଳ ଛିଲ ନାମା ଆୟୋଜନ । ହୁମାୟନେର ବିହେର ପ୍ରକାଶକ, ପାଠ୍ୟ ଏକମଧ୍ୟାନ୍ତରୀ, ବକ୍ତ୍ଵ ଓ ସାମିଦ୍ଧାନ୍ୟାଙ୍କ ବାଜିଦେବ ଆନନ୍ଦ-



ଛାଯାନ୍ଟ ମିଲନାୟତନେ ନୃତ୍ୟ ଉତ୍ସବର ପ୍ରଥମ ଦିନେ ଦଲୀଯ ପରିବେଶନା । ପ୍ରଥମ ଆଲୋ

ଉତ୍ସବ ଓ ଭାଲୋବାସ୍ୟ ବର୍ଗାଦାତାବେ ଉଦ୍ୟାପିତ ହେଁବା ତାର ଜୟଦିନ ।

ରାଜଧାନୀର ସୁଧିଯା କାମାଳ ଜାତୀୟ ଗଣଗ୍ରହାଗାର ପ୍ରାଜଗେ ଶୁରୁ ହେଁବେ ହୁମାୟନ ଆହମେଦ ବୈଇମେଲାର ଉତ୍ସୋଧନ କରେଲା ପ୍ରସାଦେ ଆହମେଦ ଏକବିଷେଳା । ଚାନ୍ଦେଲ ଆଇ ପ୍ରାଜଗେ ଛିଲ ହୁମାୟନ ମେଲା । ଗାଞ୍ଜିପୁରେ ନୁହାଶପରୀତେ ଓ ଛିଲ ନାମା ଆୟୋଜନ । ନେତ୍ରକୋନାର କେନ୍ଦ୍ରର ଉପଭୋଲାର, କୁତୁବପୁର ଓ ମୋହନଗଞ୍ଜେ ଓ ଛିଲ ଉପଭୋଲା ଅନ୍ତର୍ନାମାଳା । ଏହି ଆୟୋଜନେ ସକଳ ଦେଖିବାର ଅବଧି ଛିଲ ଉତ୍ସବର ପଦଚାରଣ ।

ଉତ୍ସବର ପ୍ରଥମ ପ୍ରହରେ ଧାରାବାହିକତା ଦାଖିଲ ହୁମାୟନପରୀତେ ମେଲାର କାମାଳ ଜାତୀୟ ଗଣଗ୍ରହାଗାରର ଖୋଲା ପ୍ରାଜଗେ ୧୫୮ ପ୍ରକାଶନା ସଂହାର ଉତ୍ସବରେ ଏକକ ବୈଇମେଲାର ଉତ୍ସୋଧନେର ଆପେଇ ଦୁର୍ପରେ ଦେଖିବାର ଭିତ୍ତି ଜମାତେ ଶୁରୁ କରେଲ । ଛୁଟିର ଦିନ କେଳକ କାନ୍ଦେନେ ।

ଶୁଭିମା କାମାଳ ଜାତୀୟ ଗଣଗ୍ରହାଗାରର ଖୋଲା ପ୍ରାଜଗେ ୧୫୮ ପ୍ରକାଶନା ସଂହାର ଉତ୍ସବରେ ଏକକ ବୈଇମେଲାର ଉତ୍ସୋଧନେର ପର ଅନ୍ତର୍ମାନ ପେକେ ପ୍ରକାଶିତ ସାଲେହ ଚୌଥୁମୀ ସମ୍ପାଦିତ ଶିଳ୍ପିକାଳର ହୁମାୟନ ରଚନାବଳୀ-୨୦୨୦ ଏର ମୋଡ଼କ ଉତ୍ସୋଧନ କରା ହେଁ । ୨୧ ନଭେମ୍ବର ପର୍ଯ୍ୟନ୍ତ ବୈଇମେଲା ପ୍ରତିଦିନ ବେଳା ୧୧୮ ଥେକେ ରାତ ୮୩୮ ପର୍ଯ୍ୟନ୍ତ ସବାର ଜନ ଖୋଲା ।

ଚାନ୍ଦେଲ ଆଇତେ ମେଲା : ଚାନ୍ଦେଲ ଆଇ ପ୍ରାଜଗେ ହୁମାୟନ ମେଲାର ଛିଲ ଗାନ, ନାଚ, ଶିତ୍ତ ଚିତ୍ରକାଳ, ମୁତ୍ତିଚାରଣା । ଜନପିଲ୍ ଚାରିତ୍ ହିମର ପରମଦେଶର ପୋଶାଳ କହୁନ ପାଞ୍ଜାବି ପରେ ଦଲ ବେଳେ ମେଲାଯ ଆମେ ତରମେରା । ତରମୀନେ ପରମେ ଛିଲ ହଲୁନ ରଙ୍ଗ ଶାଢି । ତରମୀନେ କହେଲେ ୬୭ଟି ହଲୁନ ରଙ୍ଗର ମେଲୁନ ଉଡ଼ିଯେ ଶୁରୁ ହେଁ ଯେ ମେଲାର ଆନୁଷ୍ଠାନିକତା । ଉତ୍ସବର ବସ୍ତାନେ ସାଗତ ବଜର୍ଯ୍ୟ ଦେନ ଚାନ୍ଦେଲ ଆଇରେ ବସ୍ତାପନା ପରିଚାଳକ ଫରିଦୁର ରେଜା ସାଗର ।

ନୁହାଶପରୀତେ ମୂଳ କଟକକେ ହୁମାୟନ ଆହମେଦର ମୂରାଳ ହାପନ ଅନୁଷ୍ଠାନିକ ହିମୁ ମେଲା’ ସରାସରି ସମ୍ପଦର କରିବେ ଚାନ୍ଦେଲ ଆଇ ।

## ଛାଯାନ୍ଟେ ଦୁଇ ଦିନେର ନୃତ୍ୟ ଉତ୍ସବ ଶୁରୁ

ଟାନା କେଳକ ବହୁରେ ଧାରାବାହିକତା ଏବାରେ ଛାଯାନ୍ଟ ଆଇଜନ କରିବେ ଦୁଇ ଦିନେର ନୃତ୍ୟ ଉତ୍ସବରେ । ଗତକାଳ ଉତ୍ସବର ଶୁରୁ ହେଁବେ । ଉତ୍ସୋଧନେର ଶୁରୁତେ ସାଗତ ବଜର୍ଯ୍ୟ ଦେନ ଛାଯାନ୍ଟ ସାଧାରଣ ସମ୍ପଦର ଖାୟାଳ ଆନନ୍ଦ ।

ପର୍ଯ୍ୟନ୍ତ ଦିନ ନୃତ୍ୟ ପରିବେଶନ କରେ ଛାଯାନ୍ଟ, ନୃତ୍ୟ ମେଲାର କେଳ, ଧୂତ ନର୍ତ୍ତାଳୀଙ୍କ ମେଲାର କେଳନ ଅନୁଷ୍ଠାନିକତା ଏବାରେର ପ୍ରଦଶ୍ନିଟି ତାନେର ଡ୍ରାଇସ ନିଯେ । ତାହିଁ ପ୍ରଦଶ୍ନିର ଶିରୋନାମ ଓ ଡ୍ରାଇସ ।



Figure 7.1.2: Noise Filtered Document

## 7.2 Shaded pixel removed

From the noise free image we have created a threshed Matrix of the image. This image doesn't contain any shaded pixel. For Removing shaded pixels, getting chance of correct components increases.

# গ্যালারি কায়ায় 'ড্রয়িংস'

### নিম্ন প্রতিবেদক •

গ্যালারি কায়ার এবারের চিত্র প্রদর্শনীটি কয়েকে  
প্রজন্মের নয় শিল্পীর ছাবি নিয়ে। এদের সবার  
প্রেইসের ছিলেনে খাতি রয়েছে, কিন্তু এবারের  
প্রদর্শনীটি তাদের ড্রয়িং নিয়ে। 'ভাই প্রদর্শনীর  
প্রিয়েরামও 'ড্রয়িংস''

অশ্রুণ্ঠানকারী শিল্পীরা হলেন মুর্তজা  
বশীর, চক্র পেখর দে, বুল মজুমদার, জামাল  
আহমেদ, কাজী রফিক, রশজিদ নাস, শিল্পীর  
ভট্টাচার্য, গোত্তুল চক্রবর্তী ও মোহামেদ  
ইকবাল।

প্রদর্শনীর উরোধন হয়েছে গতকাল  
বিকেলে। অনুষ্ঠানটি পরিচালনা করেন গ্যালারি  
কায়ার প্রধান পিষ্টি প্রেতম চক্রবর্তী।

অনুষ্ঠানে প্রধান অভিন্ন ছিলেন বিদ্যুৎ,  
জালানি ও খনিজ সম্পদ প্রদর্শনী নসকল  
ছায়ান। তিনি বলেন, 'এ ধরনের প্রদর্শনী আরও  
হওয়া দরকার। এককাল সরকার শিল্পকর্ম  
সংগ্রহ করে শিল্পীদের উৎসবের স্থিৎ। এখনো  
সেটা করা উচিত। আমরা এ ব্যাপারে দৃষ্টি  
নিচ্ছি।'

অভিধি হিসেবে আরও উচ্চারণ কাটি বেড় অব কেগোরেটি  
আকেফোর্ম বিপীল নাশ চৌধুরী। দর্শকদের মধ্যে  
উপস্থিত ছিলেন জাতীয় জামাল আহমেদের  
মহাপরিচালক এনামুল হক ও শিল্প সমালোচক  
সদস্যের খান।

প্রদর্শনীটি চলবে ২৭ নভেম্বর পর্যন্ত।  
প্রতিদিন বেলা ১১টা থেকে রাত ৮টার মধ্যে  
একটু সময় নিয়ে গ্যালারিতে টু মারলেই  
ছিলগুলো দেখা যাবে।

### হুমায়ুনের জন্মস্থান

বাংলা সাহিত্যের অন্যতম জনপ্রিয় লেখক,  
নাট্যকার ও চলচ্চিত্রকার হুমায়ুন আহমেদের  
৬৭তম জন্মদিন ছিল উক্তবার। এই শিল্পজগৎটা  
থাকলেও তার জন্মস্থান ঘরে গতকাল ছিল নানা  
আয়োজন। হুমায়ুনের বইয়ের প্রকাশক, পাঠক,  
অনুবালী, বক্তৃ ও সামিধানিক ব্যক্তিদের আনন্দ-



ছায়ানটি মিলনায়তনে নৃত্য উৎসবের প্রথম দিনে দলীয় পরিবেশনা • প্রথম আলো

উচ্চস ও ভালোবাসায় বর্ণিতভাবে উন্ধাপিত  
হলো তার জন্মদিন।

রাজধানীর সুফিয়া কামাল জাতীয়  
গণগ্রহণগার্গ প্রাসাদে তুর হয়েছে হুমায়ুন  
আহমেদ একটি বইমেলা। চালেন আই প্রাসাদে  
ছিল হুমায়ুন মেলা। গাজীপুরের সুব্রহ্মণ্যপুরীতেও তুর  
নানা আয়োজন।

নেতৃত্বের প্রদর্শন কৃত্বপূর্ব ও মোহনগাঁওতে ছিল  
বর্দিচা অনুষ্ঠানমালা। এসব আয়োজনে সকাল

থেকে রাত অবধি ছিল ভজনের পদচারণ।

জুকামের প্রথম প্রেরণে ধানমন্ডির 'দুখিন

হাওয়ার' হুমায়ুনপুরী মেহের আফরোজ শাওন

তুই সুতো নিয়াম ও নিনিতকে সঙে নিয়ে

কেক কাটেন।

সুবিয়া কামাল জাতীয় গণগ্রহণগারের খোলা

প্রস্তরে ১৫টি প্রকাশনা সংস্থার উত্তোলনে একক

ব্যক্তিদের উরোধনের আগেই দুপুর থেকে

পঞ্চকেরা ভড় জামাতে শুরু করেন। ছুটির দিন

হওয়ার মেলা ছিল নানা বয়সী মানবের  
পদচারণে মুখ্য। ৬৭টি বার্ষিক বেলুন আকাশে  
অবস্থৃত করার মধ্য দিয়ে নয় দিনের হুমায়ুন  
আহমেদ বইমেলার উৎসবেন করেন প্রায়ত  
সেখকের অনুজ্ঞা প্রযোজন করেন ইকবাল। এ সময়  
ছিল প্রধান অভিধি অর্ধমাসী আবুল মাল  
আবিনের মুক্তি, মেছে আফোজ সান্দেশ প্রযুক্তি।

এসব আয়োজনের প্রতি অন্যপ্রকাশ থেকে

প্রকাশিত সালেহ চৌধুরী সম্পর্কিত শিঙ্কে

কেরান হুমায়ুন রচনাত্ম-২-এর মোড়ক

উত্তোলন করা হয়। ২১ নভেম্বর পর্যন্ত বইমেলা

প্রতিদিন মেলা ১১টা থেকে রাত ৮টা পর্যন্ত

স্বর্ব জন খেলা।

চালেন আইতে মেলা: চালেন আই প্রকল্পে

হুমায়ুন মেলার ছিল গান, নাচ, শিশি চিতাকন,

স্পিচিচারণা, জাতীয় চরিত্র হিয়ে পছন্দের

পেশাক হৃদয় পাঞ্জাবি পরে দল বৈধে মেলায়

আসে তক্কেশরা। তরলীদের পরামে ছিল হৃদয়

রঙের শাপি। সকালে ৬টার হৃদয় রঙে বেলুন

উড়িয়ে ওক্ত হয় সেলার আনন্দিনিকতা।

উত্তোধনী অনুষ্ঠানে স্বাগত বজ্রব্য দেন চালেন

আইয়ের বাবহাসন। পরিচালক ফরিদুর রেজা

সাগর।

হুমায়ুনের মূল ফটকে হুমায়ুন আহমেদের

মুরাব হাপন অনুষ্ঠানসহ 'হিমু মেলা' সরাসরি

সম্পত্তির করেন চালেন আই।

### ছায়ানটে দুই দিনের নৃত্য উৎসব শুরু

টানা কয়েক বছরের ধারাবাহিকতায় এবারও

অবোজন করেন করেন দুই দিনের নৃত্য

উৎসব। গতকাল শুরু হয়েছে ছায়ানটি মিলনায়তনে

এ উৎসব শুরু হয়েছে। উৎসবের শুরুতে

স্বাগত বজ্রব্য দেন ছায়ানটি সাধারণ সম্পত্তি

খায়কল আমান।

প্রথম দিন নৃত্য পরিবেশন করে ছায়ানটি,

নৃত্য নৃত্যশীলন কেন্দ্র, ধৃতি নৰ্তনালয়,

নৃত্যদ্বন্দ্ব, সৃষ্টি কালচারাল সেকেন্ড ও নৃত্যালয়ের

শিল্পীরা। একক ও দলীয় নৃত্য অধিক অংশ দেন

তামাঙা রহমান, ওয়ার্দি বিহার, বেনজির

সালাম, পুজা সেনগুপ্ত, অমিত চৌধুরী, ওভা

সেনঙ্গ, মনিরা পারভুন, সাতা শাহরিন প্রযুক্তি।



Figure 24: Image Matrix of Non Shaded Pixel

## 7.3 Edge Detection

From the noise free image we find edges in the input image and mark them in the output map using the ‘Canny’ algorithm. The working method of ‘Canny’ edge detector is described below.

### 7.3.1 Canny Edge Detector

The *Canny Edge detector* was developed by John F. Canny in 1986. Also known to many as the *optimal detector*, canny algorithm aims to satisfy three main criteria:

- **Low error rate:** Meaning a good detection of only existent edges.
- **Good localization:** The distance between edge pixels detected and real edge pixels have to be minimized.
- **Minimal response:** Only one detector response per edge.

Steps:

1. Filter out any noise. The Gaussian filter is used for this purpose. An example of a Gaussian kernel of size = 5 that might be used is shown below:

$$k = \frac{1}{159} \begin{bmatrix} 2 & 4 & 5 & 4 & 2 \\ 4 & 9 & 12 & 9 & 4 \\ 5 & 12 & 15 & 12 & 5 \\ 4 & 9 & 12 & 9 & 4 \\ 2 & 4 & 5 & 4 & 2 \end{bmatrix} \quad (2)$$

2. Find the intensity gradient of the image. For this, we follow a procedure analogous to Sobel: a. Apply a pair of convolution masks (in x and y directions):

$$G_x = \begin{bmatrix} -1 & 0 & +1 \\ -2 & 0 & +2 \\ -1 & 0 & +1 \end{bmatrix} \quad (3)$$

$$G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ +1 & +2 & +1 \end{bmatrix} \quad (4)$$

b. Find the gradient strength and direction with:

$$G = \sqrt{G_x^2 + G_y^2} \quad (5)$$

$$\theta = \arctan\left(\frac{G_y}{G_x}\right) \quad (6)$$

3. Non-maximum suppression is applied. This removes pixels that are not considered to be part of an edge. Hence, only thin lines (candidate edges) will remain.
  4. **Hysteresis:** The final step. Canny does use two thresholds (upper and lower):
    - a. if a pixel gradient is higher than the upper threshold, the pixel is accepted as an edge
    - b. If a pixel gradient value is below the lower threshold, then it is rejected.
    - c. If the pixel gradient is between the two thresholds, then it will be accepted only if it is connected to a pixel that is above the upper threshold.
- Canny recommended an upper: lower ratio between 2:1 and 3:1.

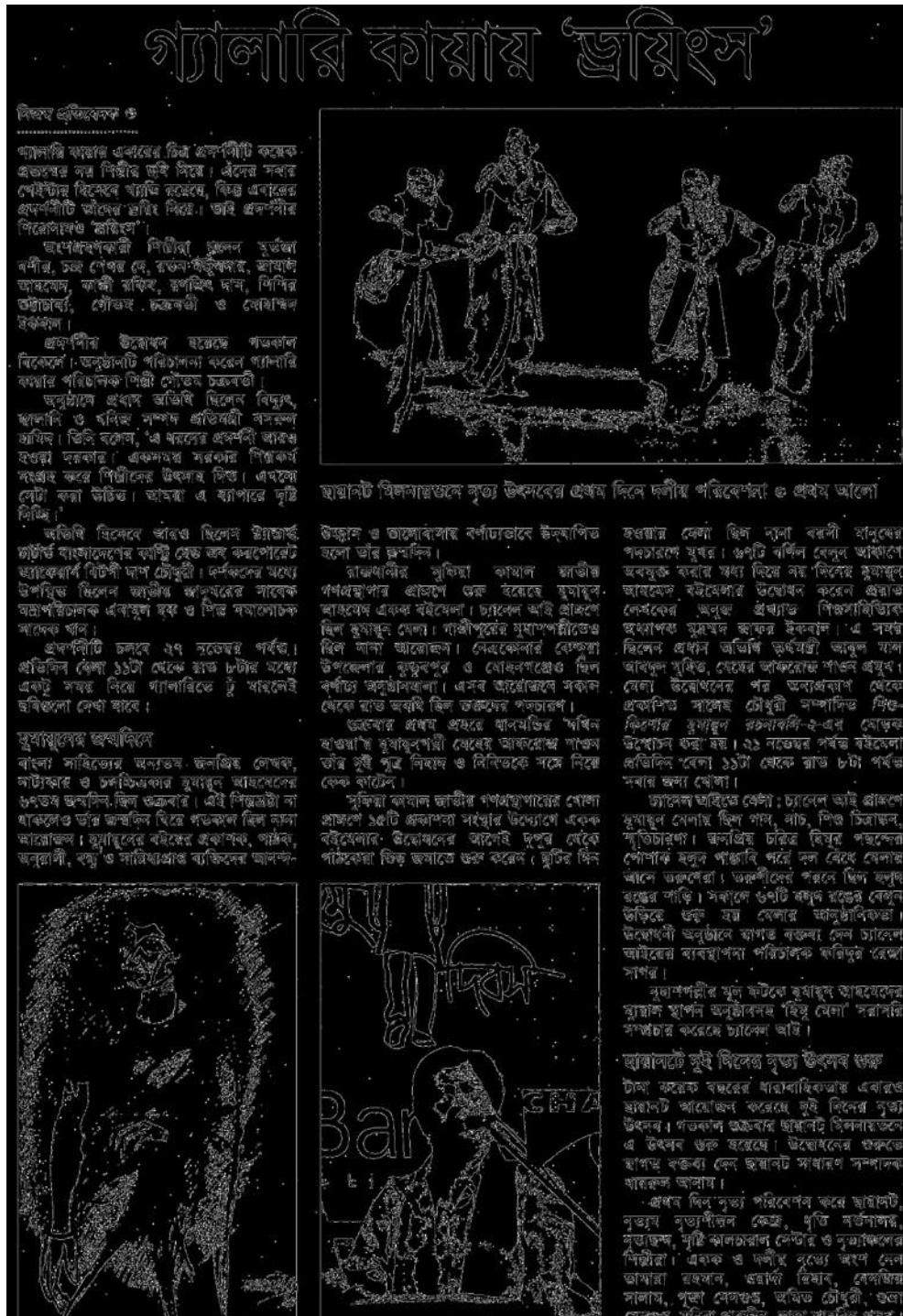


Figure 25: Edge Detected of the Input Image.

## 7.4 Conversion of input image

This process is done for creating an array of connected components. To process this step an array will contain values for every pixel according to the length of the vertical and horizontal line it belongs to. That is every pixel of same line will be given a same value, which is the total length or height of that black line. Here we assigned 0 (black) for pixel to the array if the pixel color is not black ( $>0$ ). Thus we will get an image where, the area which contains texts, images etc. will be black and the area which are not textual area (white spaced, mainly the separating space among articles) will be white.



Figure 26: Conversion Image

## **7.5 Discerning the Element Separators**

Again a binary image of the input image is need to be created for this step. As we know, usually all the textual elements of any page starts after maintaining a certain gap from the page borders [13]. We pick pixels and compare them with some threshold values to determine whether they belong to non-text area (white spaced area of converted image) or not. If any pixel, which is, we compare the height as well as comparing the width, we can anticipate that it lies on the horizontal non-text area. Again when a pixel is found which is, we compare the height as well as comparing the width, we can anticipate that it lies on the vertical non-text area. For comparing this height and we have set a static value. These static values are assigned by observing the document structures.

## **7.6 Segmenting the Black Boxes**

In this stage, we store the top-left and bottom-right coordinates of each black box neighboring the horizontal or vertical element separators from the converted image. Then we cut our image according to coordinates of the black blocks. This cutting process is done by a using the BFS Algorithm.

## **7.7 Distinguishing the Elements**

A newspaper mainly contains three kinds of elements like image, headline and columns. In this stage, we detected the category of each segmented block.

## **Image Detection:**

From the previous methodology, we just change some of the threshold value to work on higher resolution image detection. Because the number of pixels is much more than a low res documents in a high res document.

### **1<sup>st</sup> filter**

It checks whether the length and width of each block is greater than a threshold value or not. The threshold value was set by observing the document structure.

### **2<sup>nd</sup> filter**

It finds out the pixel ratio of the blocks. If the ratio is at much higher rate, then we anticipate that this block may be an image.

### **3<sup>rd</sup> filter**

It is used to determine the ratio of the sum of the total length of each block and total width of each block. If a block is image then the ratio will be greater than 1.25. If any block satisfies all these three filters then they are decided to be an image.



Figure 27: Detected Images

### Table Detection:

In last methodology, in every case table was detected as an image. So it was easy filter for us. After detecting it as an image. Another filter process is done with the sub matrix. In this filter it is checked that is there any line present which takes more than 70% in a row and is the percentage of those line is below 5%. This set of rule is checked for both vertical histogram and horizontal histogram of that sub matrix. The whole process is like below:

1.  $\text{HH}_{\max} = \text{Max( Horizontal Histogram )}$
2.  $\text{VH}_{\max} = \text{Max( Vertical Histogram )}$
3.  $\text{HH}_{\text{Range}} = \text{Threshold Percentage of } \text{HH}_{\max}$
4.  $\text{VH}_{\text{Range}} = \text{Threshold Percentage of } \text{VH}_{\max}$
5.  $\text{HHF} = \text{Frequency (All HH in } \text{HH}_{\text{Range}})$
6.  $\text{VHF} = \text{Frequency (All VH in } \text{VH}_{\text{Range}})$
7.  $\text{SHH} = (\text{HHF}/\text{Frequency (HH)}) * 100.0$
8.  $\text{SVH} = (\text{VHF}/\text{Frequency (VH)}) * 100.0$

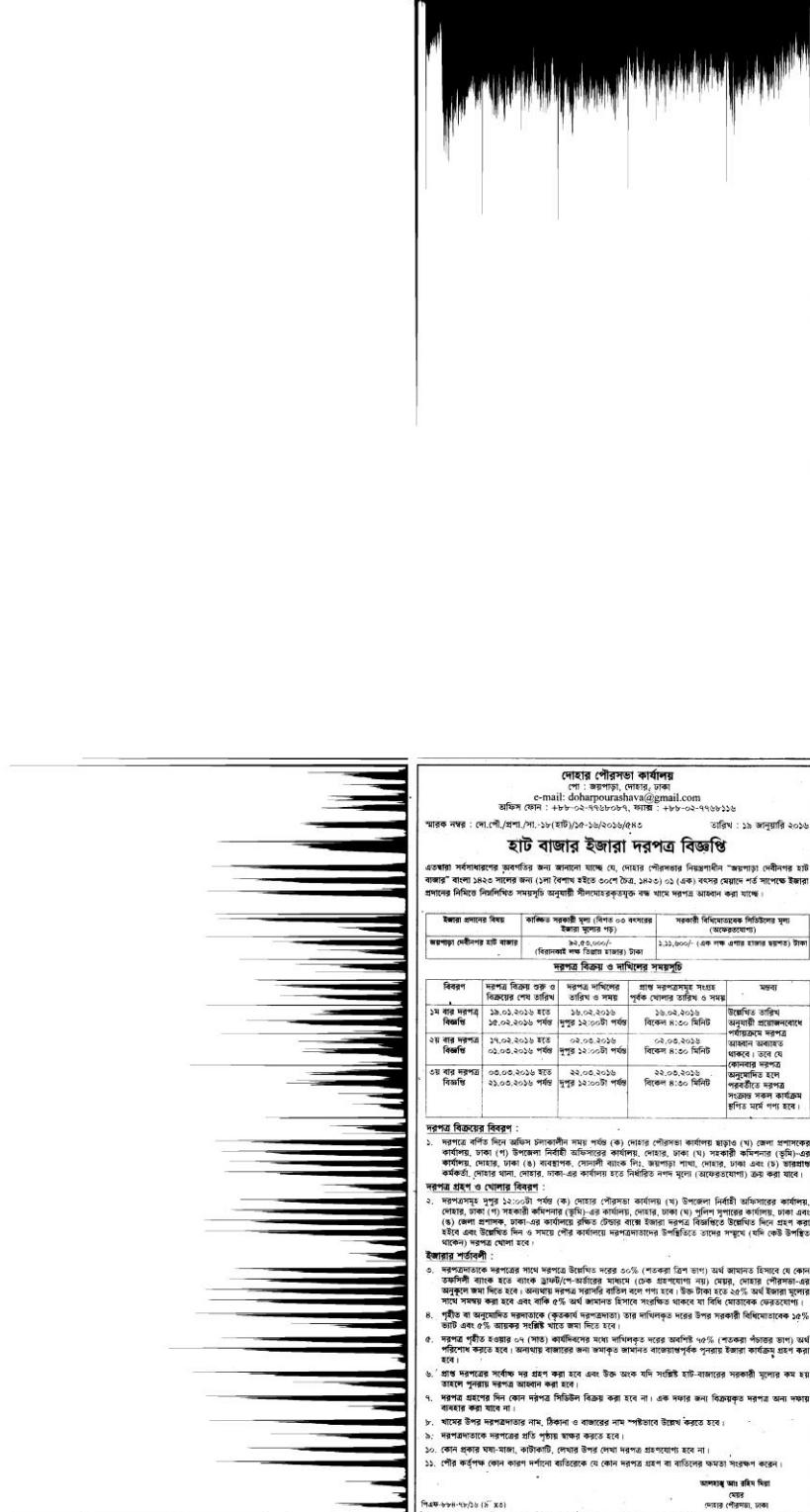


Figure 28: A Table with Vertical and Horizontal Histogram

## **Headlines Detection:**

From the previous methodology we have just change some of the threshold value to work on the case of detecting headlines. To ensure whether a block is headline or not, at first, we checked the height and width of the block and then we checked a specific ratio through the histogram of that block. Only the component left after detecting the images are used in detecting the headlines. Here we check the height and width of a block. In addition, we have found out that for being a headline a block should have a height should be greater than other textual components. The highest range is detected as headline.

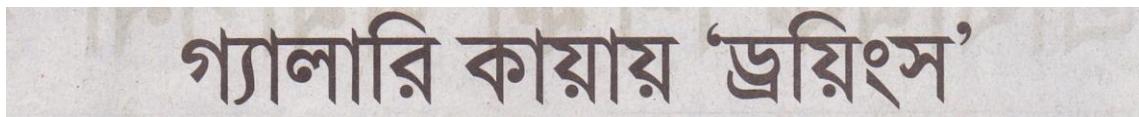


Figure 29: Detected Headline

## **Column Detection:**

A newspaper contains headlines, images, columns mainly. Therefore, other blocks are marked as columns. We can see there are 6 columns in the given image.

## **7.8 Line Height Calculation**

From the methodology of previous work, height of a line was calculated in the process like below:

1. Start from row 1.
2. If Black Bit found start= row position.
3. Go on until row with no black bit found.
4. If case 3 is True then position of that row - start is the height of the line

- Else Height of the total sub matrix is the line height.

This Methodology has a negative side if an overlapping case like below occurs.

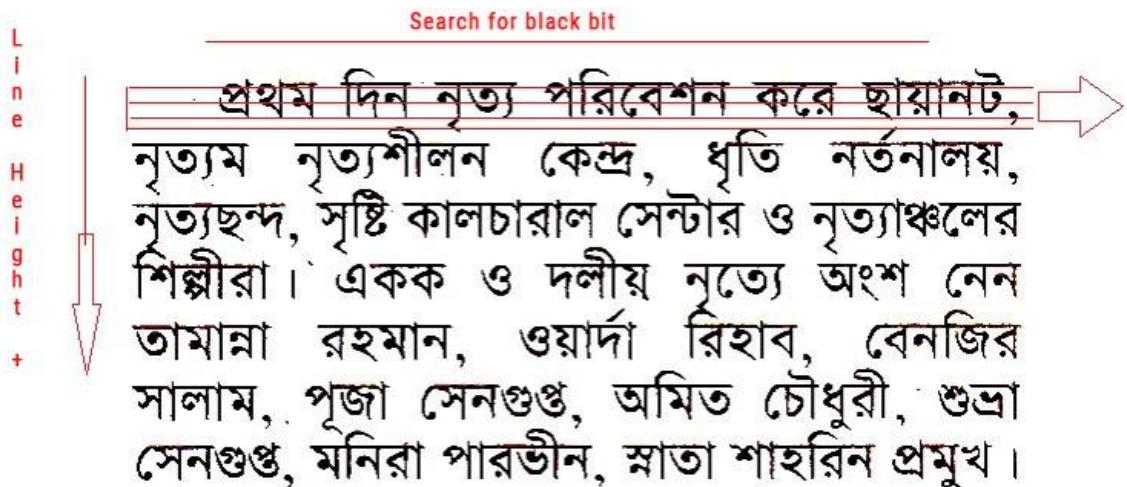
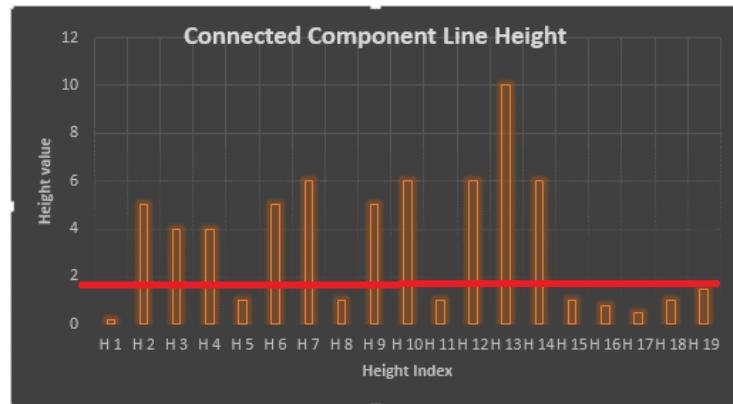


Figure 30: Line Overlapping Case

So we have proposed a new line height calculation technique.

- Take the height of all connected component which is similar to a word.
- $H_{Filter} = \text{Remove Garbage Height } (H, \text{ Threshold})$
- Finding  $H_{Avg}$  from  $H_{Filter}$
- $H_{Standard deviation}$  of  $H_{Filter}$
- Set  $H_{Max} = H_{Avg} + H_{Standard Deviation}$   

$$H_{Min} = H_{Avg} - H_{Standard Deviation}$$
- $H_{Candidate} = H_{Min} < H_{Filter} < H_{Max}$
- $H_{Line Height} = \text{Max of } H_{Candidate}$



প্রথম দিন নতুন পরিবেশন করে ছায়ানটি, নতুন নতুন শীলন কেন্দ্র, ধৃতি নতুনালয়, নতুন চূন্ডি, সৃষ্টি কালচারাল সেন্টার ও নতুনাঞ্জলের শিল্পীরা একক ও দলীয় নতুনে অংশ নেন তামানা রহমান, ওয়ার্দা রিহাব, বেনজির সালাম, পূজা সেনগুপ্ত, অমিত চৌধুরী, শুভা সেনগুপ্ত, মনিরা পারভীন, স্বাতা শাহরিন প্রমুখ।

Figure 7.8.2: Removing garbage line height

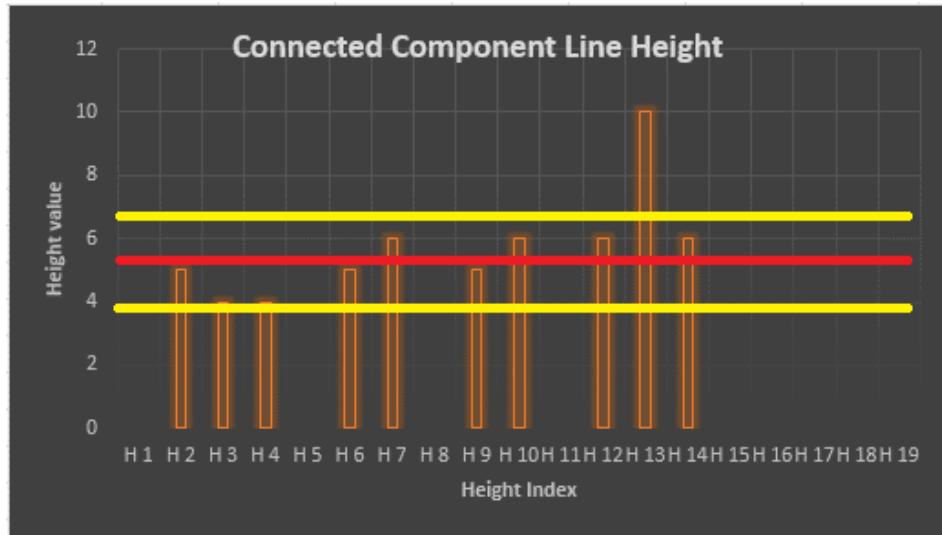


Figure 7.8.3: Line Height Standard Deviation

In this process we get a result which much more correct than previous system.

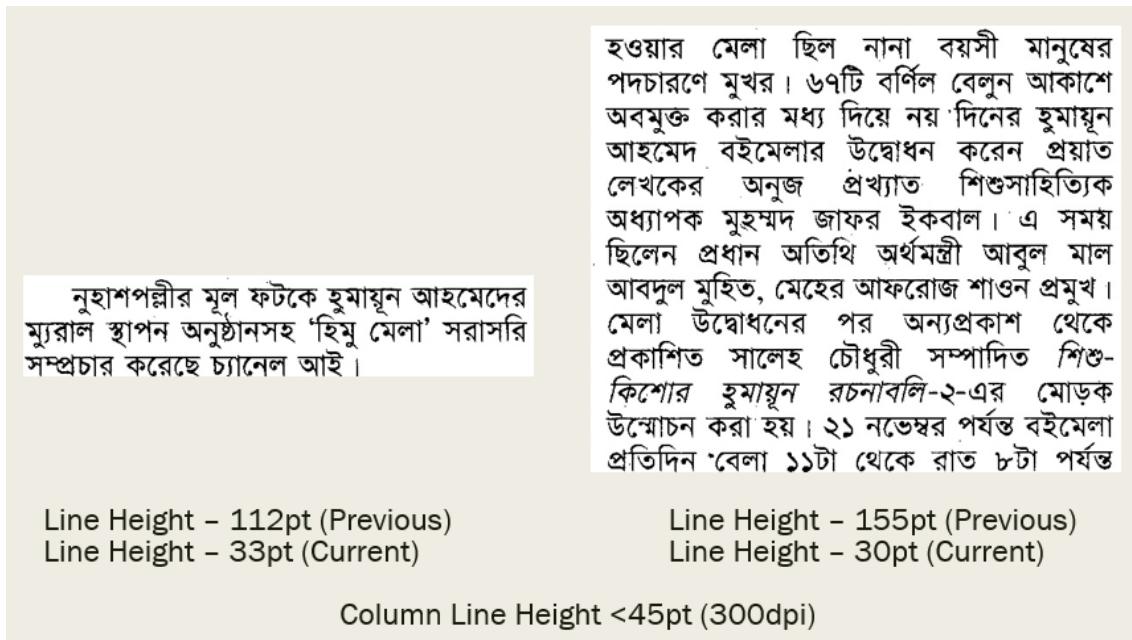


Figure 7.8.4: Line Height in new Methodology

## 7.9 Comparing Threshold Values with Previous Methodology

	Line Height for 96dpi	Line Height for 300dpi
Headline	h>50px	h>100px
Sub-Headline	h>23px	h>45px
Column	h>8px	h>18px

Figure 7.9.1: Line height threshold for categorization.

	Space gap for 96dpi	Space gap for 300dpi
HORIZONTAL NON TEXT AREA VERTICAL VALUE	6	8
HORIZONTAL NON TEXT AREA HORIZONTAL VALUE	150	180
VERTICAL NON TEXT AREA HORIZONTAL VALUE	7	9
VERTICAL NON TEXT AREA VERTICAL VALUE	200	270

Figure 7.9.2: Conversion space threshold.

## **Chapter 8**

# **Current Methodology: Part 2**

In this chapter we discussed how we improved our performances by solving the skew problem. Also we discussed about the document reconstruction system after recognition.

### **8.1 Skew Detection and De-skew**

During scanning an image, if the document wasn't alignment correctly, we have to face the skew problem. So if we start to decompose a skew image, then our methodology will not work perfectly. To going on with our methodology, we must DE Skew it first.

### **8.2 Using Hough Transform for De-skew**

Hough transform is very handful when we need to extract shape features form the image. In most of the cases, it is used for detecting lines but we can also use it to find circles and ellipses by extending the features. We also can extend this methodology to detect the skew of a scan image, then calculate the skew angle. Using this property, we have to rotate the image to remove the skew. The Interline Cross-Correlation concept and the Hough Transform was used by Hong Yan to perform the skew correction option [24]. In the year 1997 Huei-Fen Jiang, Chin chuan Han, Kuo-Chin Fan also proposed a fast approach to detect the skew document using the Hough Transformation [25]. In recent work Chandan Singh, Nitin Bhatia and Amandeep Kaur researched on accurate skew correction also based on the Hough Transform [26]. So we implemented their work to detect any kind of skew. Using this methodology we saw that, we were able to perform correctly if there are no image. But in case of an image, the methodology fails for detecting any kind of line. Rather than it detects others shape.



Figure 8.2.1: Detected line property using Hough Transform.



Figure 8.2.2: Rotated by the angle got from Hough Transform.

### 8.3 Detecting Skew Contour and rotating by its slop

Contours are useful in terms of analyzing shape an object along with detecting and recognize it. Using this system we detect the border of the image. The border of the document a rectangle. So we will say any other detect contours are not useful for. If there is any image, it's contours will also be an rectangle and it will be parallel to it's parent. We then

detect how much the rectangle is rotated perspective to (0, 0) point. The result angle will help us to correct the skew.



Figure 8.3.1: Detected Contour (Green Box) and rotated by the angle.

## 8.4 Reconstructing the Document

After segmentation, The Textual part will be sent to the recognition, it will return an output as text. This text will be used for reconstruction the document again. When the document was segmented, we saved the coordinate points to recognize from where it was segmented and where it was located. We then sorted the position of the component and added to a document one by one.

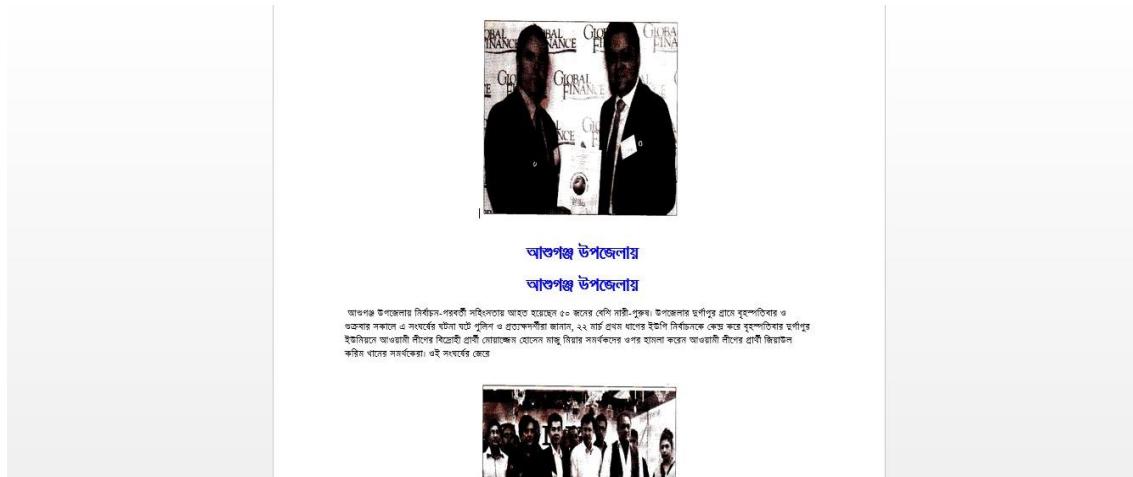


Figure 8.4.1: Screenshot of a portion of a reconstructed Document.

## Chapter 9

# Result Analysis

The methodology we improved for the decomposition part was mainly tested on various pages from newspaper. A big portion of those testing dataset was taken from “Prothom-Alo”. All the scan copy was in 300dpi. We also increased our dataset to test the result. The basic of our analysis is, we took a document count how many images there are, how many paragraphs also how many headlines. Then we compare is the output is correct or not.

### 9.1 Accuracy for Current Decomposition Methodology

CONFUSION MATRIX		P R E D I C T E D				
A C T U A L		Image	Table	Headline	Sub-Headline	Col- umn
	Image	112	1	1	0	0
	Table	2	27	0	0	0
	Headline	1	0	122	2	0
	Sub-Headline	0	0	1	61	1
	Column	0	0	0	4	300

Figure 8.1: Result analysis of current methodology

	TP	TN	FP	FN	Accuracy
Image	112	518	3	2	98.25%
Table	27	615	1	2	93.10%
Headline	122	508	2	3	97.6%
Sub-Headline	61	566	6	2	96.83%
Column	300	334	1	4	98.68%

## 9.2 Accuracy Comparison with Previous (Batch 2010) Decomposition Methodology

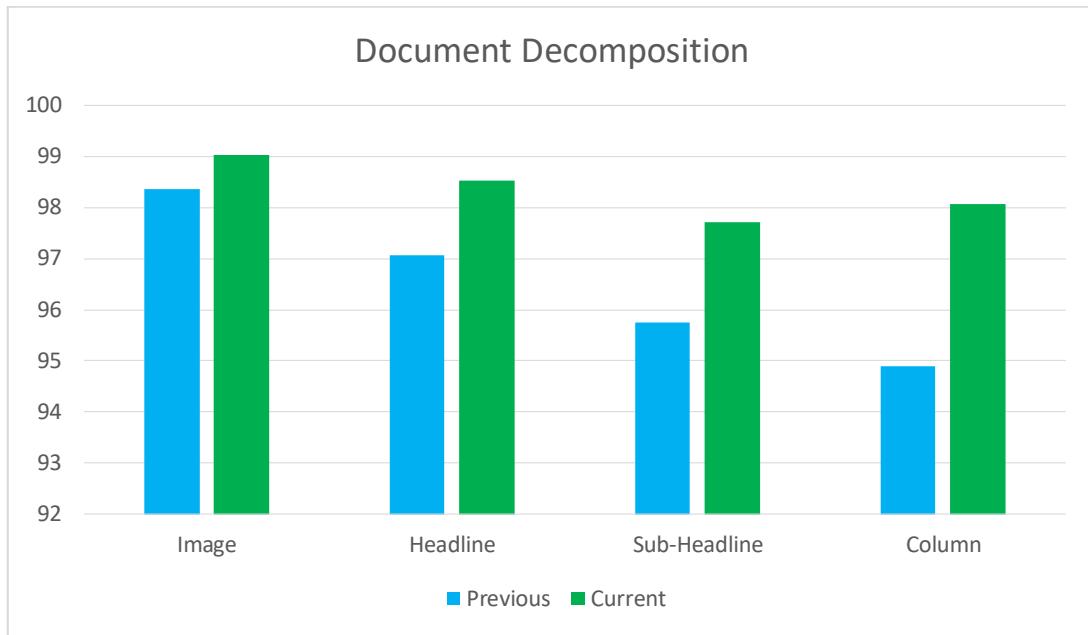


Figure 9.1: Result comparing with previous methodology.

# Chapter 10

## Conclusion

In this last chapter will discuss about how the whole could be more improved also what necessary steps should be taken.

### 10.1 Findings in our Developed Methodology

Though we wanted to solve most of the problems for segmenting a Bangla Document, there are still some lacking which should be eliminated soon. Some of these are minor issues, some should be should by another individual team. These cases are discussed below:

- i. Current system doesn't have a dynamic range for decomposition. Dynamic range will solve decomposing any kind of documents perfectly. Current system was analyzed and studied by observing the newspaper most of the example was "Prothom Alo" and "Jonokantha".
- ii. For the DE Skew Methodology, there will be some problem if the Contours slop is not inside the range of,  $-45 \geq \text{slop} \geq 45$ . If we exclude this range the output will not be correct and the system will failed to segment. Also our system cannot detect that, was the given document was upside down or not.
- iii. In our Document Reconstruction system, we just append the result by the serial of the component from their position in the document. But this should be reconstructed as the document. Our observation was that, this part should be researched by another team. Because for reconstructing the document format of the output should be well studied as it can always differ for the layout of the document.

## References

- [1] L. Eikvil, Optical Character Recognition, December 1993.
- [2] Wikipedia, "Optical character recognition," [Online]. Available: [https://en.wikipedia.org/wiki/Optical\\_character\\_recognition](https://en.wikipedia.org/wiki/Optical_character_recognition). [Accessed 8 April 2016].
- [3] blogDNA, "Deskew Scanned Images of Documents & Text Pages with Solway Deskew," [Online]. Available: <http://www.blogsdna.com/18083/deskew-scanned-images-of-documents-text-pages-with-solway-deskew.htm>. [Accessed 9 April 2016].
- [4] C. Project, "Detect image skew angle and deskew image," [Online]. Available: <http://www.codeproject.com/Articles/104248/Detect-image-skew-angle-and-deskew-image>. [Accessed 9 April 2016].
- [5] B. S. Mehmet Sezgin, "Survey over image thresholding techniques and quantitative performance evaluation," p. 23, January 2004.
- [6] K. Summers, "Automatic Discovery of Logical Document Structure," *Ph. D. Thesis*, 1998.
- [7] C. Strouthopoulos, "Text identification for document image analysis using a neural network," *Image and Vision Computation*, vol. 16, no. 12-13, 1998.
- [8] Wikipedia, "International Conference on Document Analysis and Recognition," www.wikipedia.org, 4 February 2015. [Online]. Available: [https://en.wikipedia.org/wiki/International\\_Conference\\_on\\_Document\\_Analysis\\_and\\_Recognition](https://en.wikipedia.org/wiki/International_Conference_on_Document_Analysis_and_Recognition).
- [9] B. Gatos, S. Mantzaris, K. Chandrios, A. Tigris and S. J. Perantonis, "Integrated algorithms for newspaper page decomposition and article tracking," in *ICDAR*, Bangalore, India, 1999.
- [10] C. Strouthopoulos and N. Papamarkos, "Text identification for document image analysis using a neural network," *Image and Vision Computing*, vol. 16, no. 12-13, 1998.
- [11] Wahl, M. F., Wong, Y. K., Casey and R. G., Block Segmentation and Text Extraction in Mixed Text/ Image Documents, 1982.

- [12] B. Kong, S. Chen, Haralick, M. R., Phillips and T. I., Automatic Line Detection in Document Images Using Recursive Morphological Transforms, 1995.
- [13] K. Hadjar, O. Hitz and R. Ingold, "Newspaper Page Decomposition using a Split and Merge Approach," in *International Conference on Document Analysis and Recognition*, Seattle, 2011.
- [14] P. Saikrishna and A. G. Ramakrishnan, "Script independent detection of bold words in multi font-size documents.,," in *In proceeding of National Conference on Computer Vision, Pattern Recognition Image Processing and Graphics (NCVPRIPG-2013)*, 2013.
- [15] D. S. B. F. R. Chen and L. D. Wilcox, "Detection and location of multicharacter sequences in lines of imaged text," *Journal of Electronic Imaging*, vol. 5, p. 37 – 49, January 1996.
- [16] D. B. K. R. McConnell and R. Schaphorst, *FAX: Digital Facsimile Technology and Applications*, Artech House, 1989.
- [17] P. N. M. Javed and B. Chaudhuri, "Automatic detection of font size straight from run length compressed text documents," *International Journal of Computer Science and Information Technologies*, vol. 5 (1), p. 818–825, 2014.
- [18] A. Z. L. Likforman-Sulem and . B. Taconet, "Text line segmentation of historical documents: a survey," *International Journal on Document Analysis and Recognition*, vol. 9, pp. 123-138, 2007.
- [19] J. Lu and D. Jiang, "Survey on the technology of image processing based on dct compressed domain," in *International Conference on Multimedia Technology (ICMT)*,, 2011.
- [20] P. N. M. Javed and B. Chaudhuri, "Extraction of projection profile, run-histogram and entropy features straight from run-length compressed documents.,," in *Proceedings of Second IAPR Asian Conference on Pattern Recognition (ACPR13)*, Okinawa, Japan, 2013.
- [21] U. Pal and B. B. Chaudhuri, "OCR in Bangla: an Indo-Bangladeshi language," in *Conference B: Computer Vision & Image Processing., Proceedings of the 12th IAPR International*, 1994.
- [22] U. Pal and B. B. Chaudhuri, "A complete printed Bangla OCR system," vol. 31, no. 5, p. 531–549, 1 March 1998.

- [23] Wikipedia, "Tesseract (software)," [Online]. Available: [https://en.wikipedia.org/wiki/Tesseract\\_\(software\)](https://en.wikipedia.org/wiki/Tesseract_(software)). [Accessed 11 April 2016].
- [24] H. Yan, "Skew Correction of Document Images Using Interline Cross-Correlation," *Graphical Models and Image Processing*, vol. 55, no. 6, pp. 538-543, 1993.
- [25] H.-F. Jiang, . C. . c. Han and K.-C. Fan, "A Fast approach to the detection and correction of skew documents," vol. 18, pp. 675-686, 1887.
- [26] C. Singh, N. Bhatia and A. Kaur, "Hough Transform Based Fast Skew Detection and Accurate Skew Correction Methods," *elsevier*, vol. 41, pp. 3528-3546, 2008.
- [27] R. Szeliski, Computer Vision: Algorithms and Applications, 2010.
- [28] Rich and Knight, Artificial intelligence, 2009.
- [29] systransoft, "Machine Translation," [Online]. Available: <http://www.systransoft.com/systran/corporate-profile/translation-technology/what-is-machine-translation/>. [Accessed 8 April 2016].
- [30] elsevier, "Pattern Recognition," [Online]. Available: <http://www.journals.elsevier.com/pattern-recognition/>. [Accessed 8 April 2016].