# ECE 415 Project First Draft

# Optical Character Recognition

## Group Partners:

- Yunjuan Wang <661 180 568>
- Usama Muneeb <656 392 492>

## Section 1: Problem Statement

The goal is to develop an OCR (Optical Character Recognition) application that can take a scanned document image as an input and converts it to digital text with a high degree of accuracy. Traditionally, documents scanned into a computer are saved as PDFs and can be read only on a computer. OCR not only allows a user to scan documents and save them to a computer, but also to be able to edit the documents. Vividly, the application is such like a combination of CamScanner and Adobe Acrobat Pro. It can be widely used since it often occurs when someone tries to scan some paper work into editable files.

The final application we develop is web based (written in NodeJS and Python at backend).

## Section 2: Implementation Details

We are using Python along with the following libraries:

- OpenCV for regular image processing procedures
- TensorFlow for implementing LeNet-5 neural network to recognize letters
- PDFRW for PDF write operations to generate a PDF file.

Also, we are using NodeJS for displaying a web-based user interface.

## Section 3: Block Diagram of the System

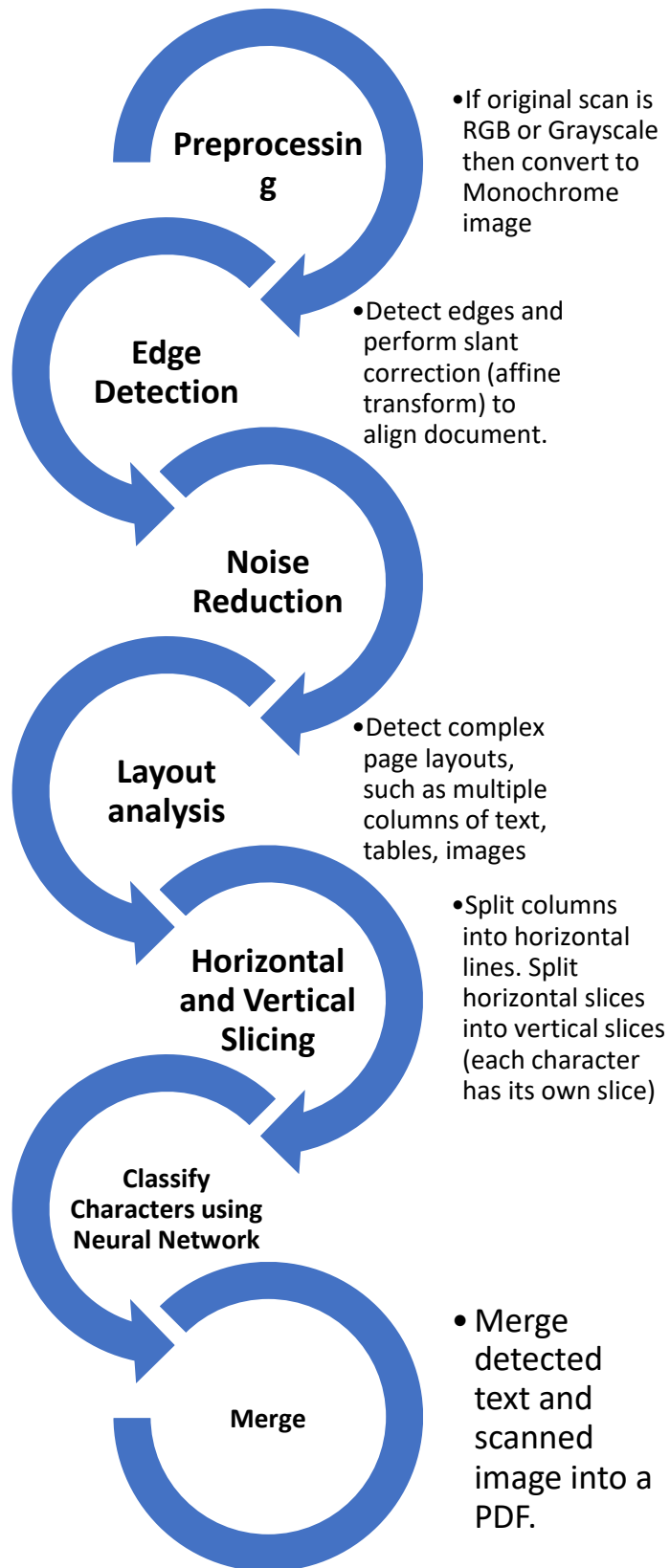The following block diagram summarizes our steps for achieving the OCR application.

**Preprocessing**

- If original scan is RGB or Grayscale then convert to Monochrome image

**Edge Detection**

- Detect edges and perform slant correction (affine transform) to align document.

**Noise Reduction**

- Detect complex page layouts, such as multiple columns of text, tables, images

**Layout analysis**

**Horizontal and Vertical Slicing**

- Split columns into horizontal lines. Split horizontal slices into vertical slices (each character has its own slice)

**Classify Characters using Neural Network**

**Merge**

- Merge detected text and scanned image into a PDF.

*Figure 1 Block Diagram*

## Section 4: Inputs and Outputs

Our input is a JPG, JPEG, PNG or TIF image of a document scanned in at least 300 DPI.

The output is a PDF document which has the **original image in the background**, along with **hidden OCR text overlaid which can be selected and copied**.

## Section 5: Walkthrough of the process

The whole process of OCR can be broken down into six main steps. Each step with its sub steps is summarized below.

1. **Edge Detection** – find the largest object in the image (which should be the paper)
   a. We convert the RGB image into **grayscale**.
   b. We choose a standard size for the image to perform the projective transform. We **resize** the image to height 800 pixels and the width is chosen to retain the aspect ratio.
   c. We perform **bilateral filtering** to reduce noise and smoothen the image, while preserving edges.
   d. We perform **adaptive thresholding** to create a monochrome image. Adaptive thresholding calculates the optimal threshold before applying it.
   e. We perform **median filtering** to get rid of small details.
   f. We **find the contours** of **ALL objects** inside the image. Each object is represented by a sequence of points which represent the contour. We choose four points for each object using OpenCV's polygon approximation feature.
   g. We find the **object with the largest area**. The area is computed using the four coordinates for every (approximated) polygon.
   h. We choose the object with the largest area. This object is considered to be the paper which we need to perform OCR on.
   i. We use these four **source coordinates** and map them to the four corners of the target image, which are (0,0), (0,WIDTH), (HEIGHT,0) and (HEIGHT,WIDTH). These are the **target coordinates**.
   j. This is done by computing a projective transform that maps these **source coordinates** to **target coordinates**.
   k. This results in a good enough approximation of the paper, as shown in the screenshot below.
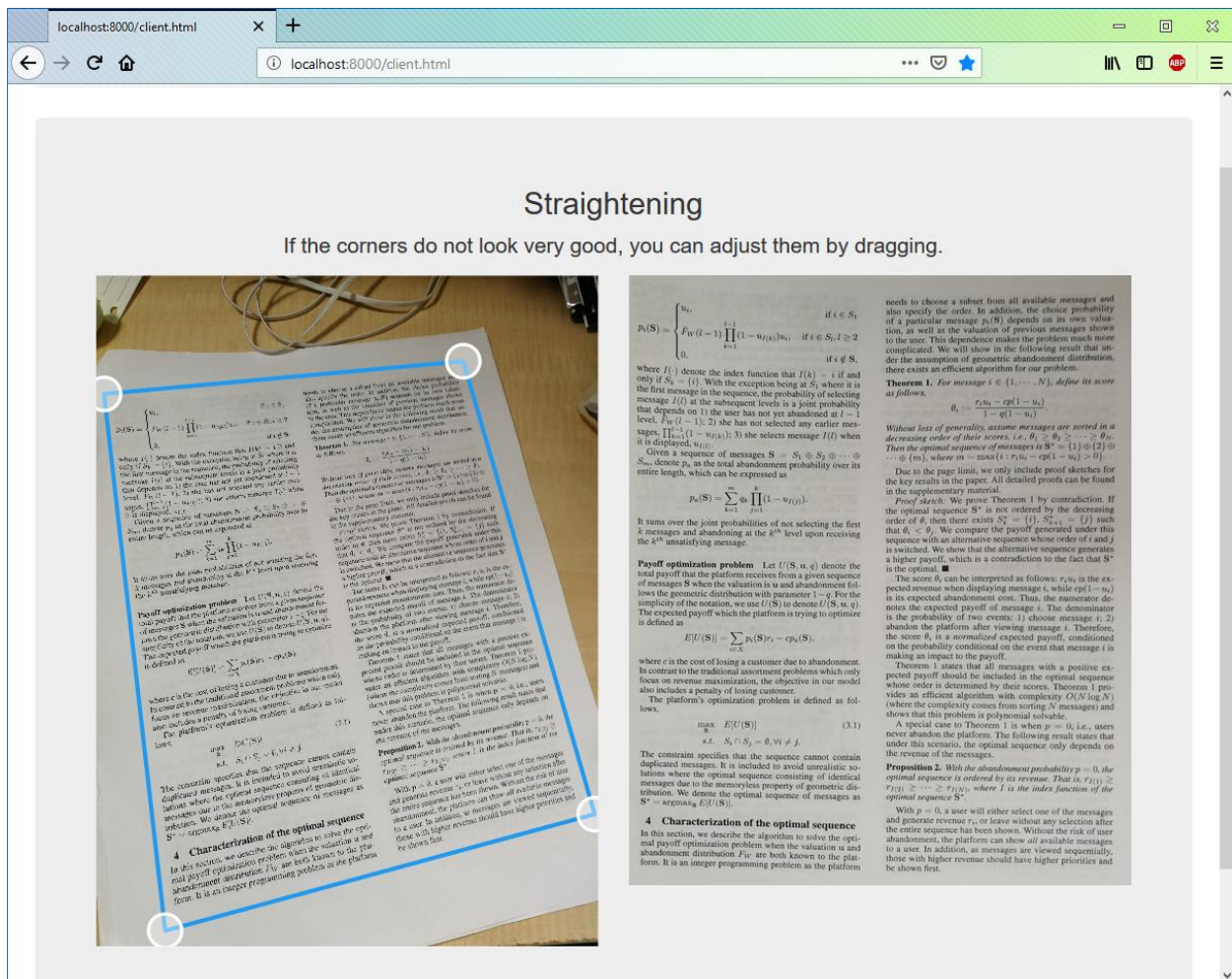   l. In case the approximation is not good, the user can adjust the corners and obtain better straightening.

*Figure 2 Screenshot - Straightening*

2. **Layout Analysis** – breaking the page down into individual objects.
    a. There are three algorithms for doing layout analysis. We integrate two of them and let the user choose whichever layout analysis he thinks performs better. The layout analysis algorithms are described in Section 6. We also give the user a choice to add undetected regions and remove unnecessary regions from the layout analysis.
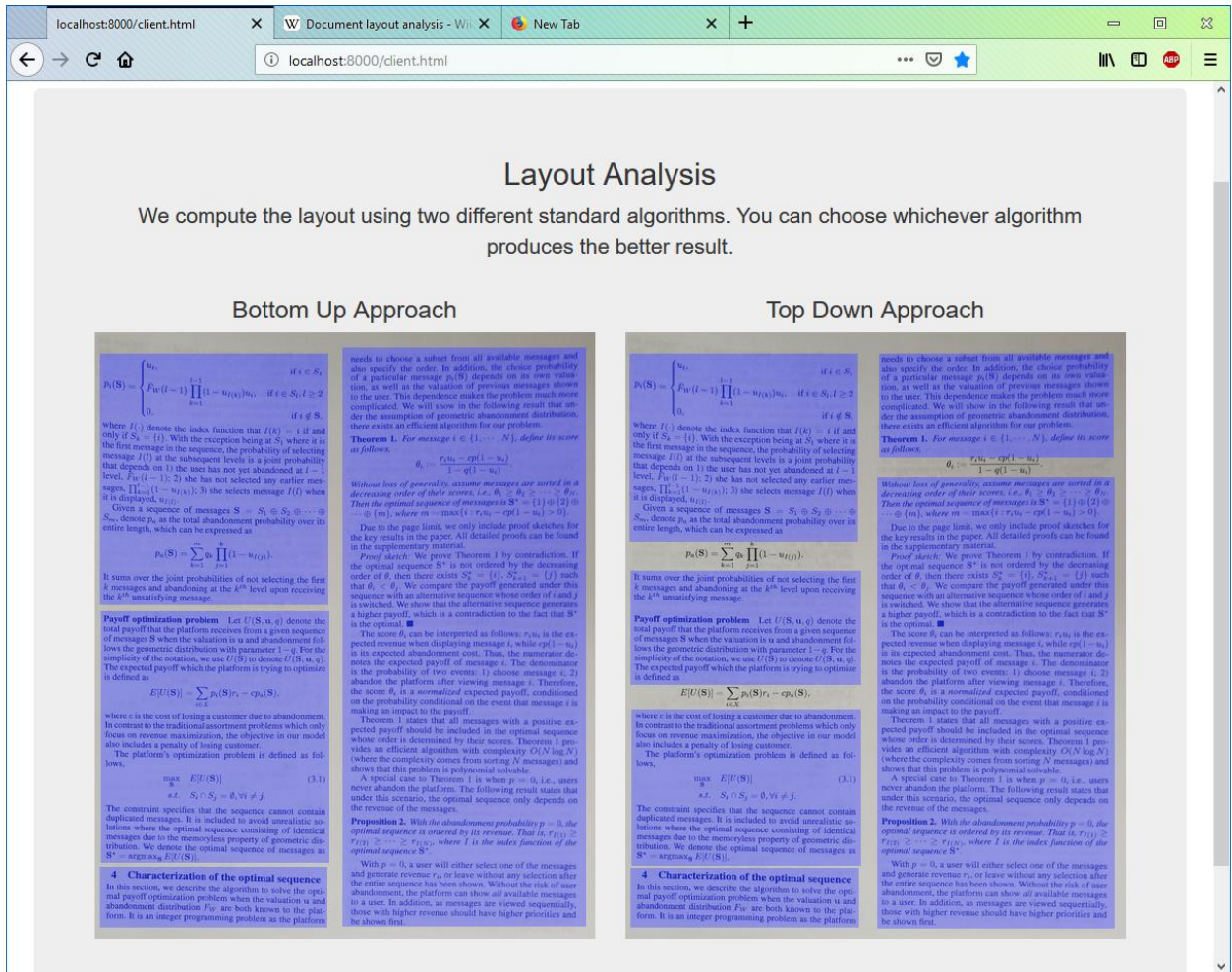
*Figure 3 Screenshot - Layout Analysis*

3. **Line Splitting** – divide blocks of text into individual lines.
    a. Iterate over each block of text given by the Layout Analysis algorithm.
    b. For each block, scan from the top most row of pixels to the bottom most row of pixels.
    c. Calculate the average pixel value of each row.
    d. Apply a threshold on the average pixel value.
    e. Whenever the average value goes below the threshold, create a horizontal line

*Figure 4 Screenshot - Line Splitting*

4. **Character Splitting**
   a. Work in process
5. **Character recognition**
   a. Work in process
6. **Consolidation into a PDF file.**
   a. Not started

# Section 6 - Layout Analysis

Layout analysis is the process of identifying and labeling the regions of interest in the scanned image of a text document.

We need to identify the textual zones from non-textual ones and their arrangement in their correct reading order, so that the copy paste operations copy them in the right sequence.

For layout analysis, we use three algorithms to implement and try to find which one is better. All of our input is binary image which the text is black and background is white.

## Algorithm 1:

This is the simplest one, albeit the results are not as good.

1. Set a threshold (240 for example).
2. Iterate the width, compute the average pixel value of each column, if the average pixel value is greater than 240, then we draw a split line. After that, we can get several columns divided by vertical split line.
3. Do the same thing- iterate the height. Compute the average pixel value of each row, find the average pixel value greater than threshold then draw a horizontal split line. However, before drawing a horizontal split line, make sure it is at least 25 pixels from the previous line above the threshold. This makes sure that this algorithm only separates blocks of text, not individual lines of text.

**Pros: low computation time**

**Cons: Only satisfactory performance**

## Algorithm 2:

This is the same as algorithm 1, however, instead of using a threshold on the average value of pixels, we compute the FFT (Fast Fourier Transform) of a column or a row depending on whether we are doing a horizontal pass (i.e. Step 2 of algorithm 1) or vertical pass  (i.e. Step 3 of algorithm 2). Then we apply a threshold on the magnitude of the high frequencies of the FFT. This detects large variation between the pixels which would be the case if we are dealing with text (there will be a lot of variation between black and white).

**Pros: better performance, not very high computation time (FFT can be done very fast in modern day computers)**

**Cons: still not as good as algorithm 3**

## Algorithm 3:

The purpose of the algorithm is to group contiguous pixels under one label. Each block of text will be assigned a label.

First few steps are preprocessing.

1. Invert the image such that the text is white and background is black.
2. Apply dilation with a suitable kernel size to thicken the text.
3. Blur the image to remove little noise.
4. Convert the image to monochrome by thresholding.

Next, we move onto actual layout analysis.

5. Set two arrays both of size *WIDTH x HEIGHT*:
6. One should be boolean (0 or 1), showing whether the pixel has been grouped (labeled) or not. Let's call this *LABELED_OR_NOT*.
7. Another should not be binary. It should be integer. It should save the actual label (group) of the pixel. Let's call this *PIXEL_GROUP*. The label is also the label of a block of text.
8. Run the following algorithm.

Initialize *i* =1

Loop until *LABELED_OR_NOT(x,y)* contains as many ones as the number of white pixels in the preprocessed image.

> Iterate from the top left pixel *(x,y)* in any direction*
> *see Section 8 for the actual implementation details*

> > If *LABELED_OR_NOT(x,y)*==0:

> > > Flood fill the pixels with label *i*
> > > Update all values *PIXEL_GROUP(x,y)* to *i*
> > > For all these pixels, also set *LABELED_OR_NOT(x,y)* to 1

> Update *i* = *i* + 1 (so that in the next iteration for the next object, PIXEL_GROUP(x,y) will be assigned the new value)

9. Find the 4 points of each block. We iterate the pixel in each block to find the maximum and minimum of x axis and y axis. Then *(min_x, min_y), (min_x, max_y), (max_x, min_y), (max_x, max_y)* is 4 corners of each block.
   The output of the algorithm is obtained the sets of 4 coordinates each for all possible blocks of text.

**Pros: Excellent performance**

**Cons: High running time**

We are working to reduce the running time of Algorithm 3 and use that in the final implementation alongside Algorithm 2 (providing the user a choice between the two).

## Section 7 - Testing and Performance Measurement

We will test whether our algorithm can identify text in different formats (font, margin, line space, text direction), whether our algorithm can distinguish between words and pictures. We will test in two column and three column layouts.

The performance will be measured by the accuracy of text recognition and conversion. The task should be accomplished in a reasonable time frame (no more than a few seconds per page).

## Section 8 – Challenges

In Algorithm 3 for layout analysis, the running complexity was large because for each round of iteration, we had to start from the first pixel. We took an alternate approach. We decided to randomly generate a pair of coordinates, and keep generating them, until we land at an unlabeled pixel.

It takes maximum 3 or 4 random guesses to land at an ungrouped pixel for the large blocks and maximum 20-30 guesses for the smaller blocks. It is hence more efficient than the iteration technique.

It takes significant trial and testing to devise the thresholds for line splitting and layout analysis. We choose thresholds that would be sufficient for major paper sizes like A4, letter, etc. scanned in 300 DPI to 600 DPI (the most common scanning resolutions).

Specifying other parameters like the size of dilation kernel, etc. also takes a lot of calibration.

## Section 9 - Limitations

We assume that our input satisfies the following requirements:

- Paper should be placed on a flat surface. A simple projective transform should be enough to straighten the page.
- Should contain the entire content of the page (no missing parts)

## References

[1] https://opencv.org

[2] https://www.tensorflow.org

[3] http://homepages.cae.wisc.edu/~ece539/project/f03/sarkar-ppt.pdf

[4] http://yann.lecun.com/exdb/lenet

[5] https://en.wikipedia.org/wiki/Document_layout_analysis