# Optical Character Recognisation

Supriya kale[1], Ajinkya Shinde[2], Somsundarji Shinde[3], Prof. Pallavi Chandake[4]

[1,2,3,4]*(MarathwadaMitra Mandals Institute of technology pune), (department of E&TC)*

***ABSTRACT:*** *Optical Character Recognition by using Template Matching is a system which is useful to recognize the character or alphabets in the given text by comparing two images of the alphabet. The objectives of this system prototype are to develop a program for the Optical Character Recognition (OCR) system by using the Template Matching algorithm . This system has its own scopes which are using Template Matching as the algorithm that applied to recognize the characters, which are in both in capitals and in small (A – Z),and the numbers (0 -9) used with courier new font type, using bitmap image format with 240 x 240 image size and recognizing the alphabet by comparing between images which are already stored in our database is already . The purpose of this system prototype is to solve the problems of blind peoples who are not able to read , in recognizing the character which is before that it is difficult to recognize the character without using any techniques and Template Matching is as one of the solution to overcome the problem.*

***Keywords:*** *Matlab*

## I. INTRODUCTION

The Optical Character Recognition is a mobile application.which is implemented on MATLAB and it requires only MATLAB supportive laptop or pc.also It can implement on smart mobile phones of android platform. This paper combines the functionality of Optical Character Recognition and speech synthesizer. The objective is to develop user friendly application which performs image to speech conversion system using MATLAB or android phones. The OCR takes image as the input, gets text from that image and then converts it into speech. This system can be useful in various applications like banking, legal industry, other industries, and home and office automation. It mainly designed for people who are unable to read any type of text documents. In this paper, the character recognition method is presented by using OCR technology and android phone with higher quality camera. OCR is the acronym for Optical Character Recognition. This technology allows to automatically recognizing characters through an optical mechanism. In case of human beings, our eyes are optical mechanism. The image seen by eyes is input for brain. The ability to understand these inputs varies in each person according to many factors. OCR is a technology that functions like human ability of reading. Although OCR is not able to compete with human reading capabilities.
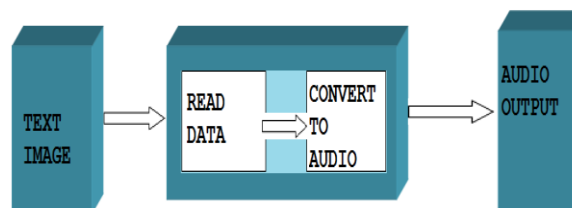


FIG. TYPICAL BLOCK DIAGRAM OF OCR

## II. Objective

i) To develop a prototype of Optical Character Recognition (OCR) system. ii) To apply a Template Matching approach in recognizing character.
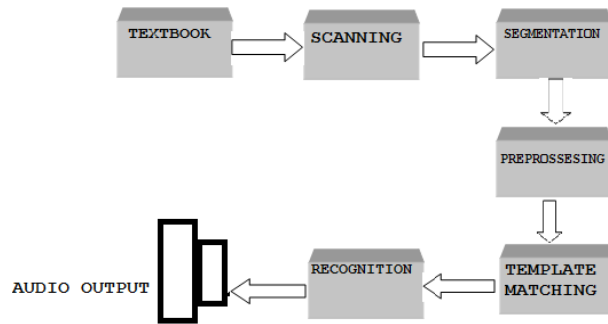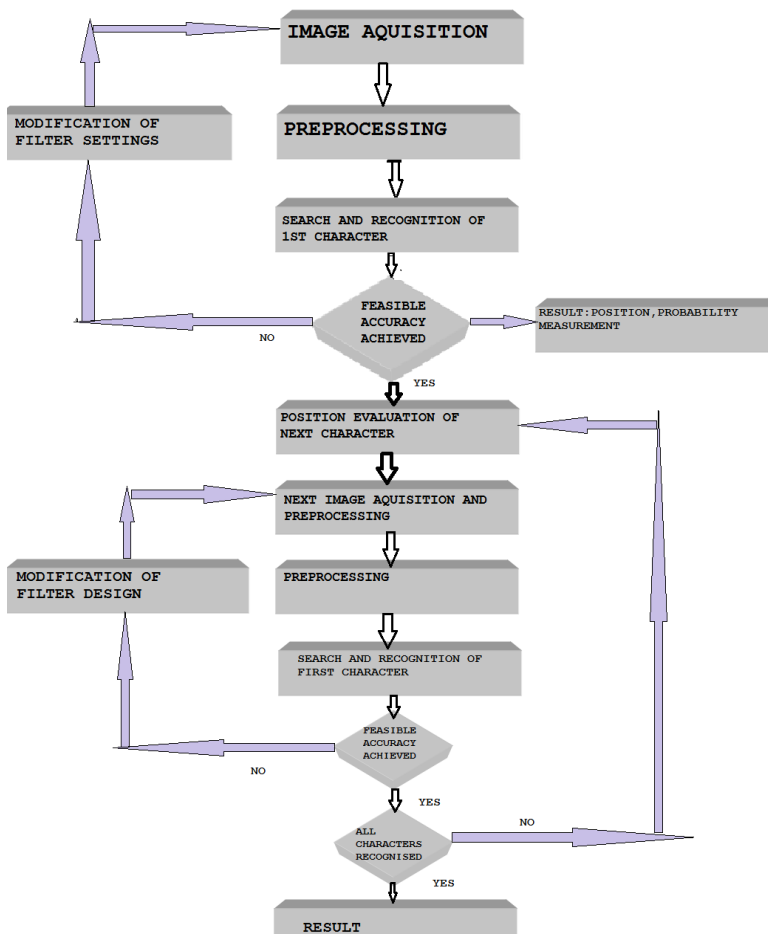
## III. Block Diagram



FIG: OPTIVAL CHARACTER RECOGNITION

## IV. Block Diagram Discription

IN this system first we scan the data which we want to read through scanner .after the scanning data is segmented in different small parts .after segmentation it is easy to recognize for the machine .segmentation makes each letter separate .then segmentized data is send for the preprocessing .in this preprocessing the size of all letters makes equal ,because in given image the font size of a letters may be different .so we adjust the size and the font of the letters according to machine which is 24 x 48.then after preprocessing  complete word is formed and this word is read and converted to audio.

## V. Flowchart

**Scanning :** in scanning process we scan the given image by an scanner .for the scnning we use the camera which havinh high quality of resolution.if the quality of the camera bad then the audio we get is not as good as we want.after scanning we send the scnned image for the segmentation.

**Pre-processing**

       OCR software often "pre-processes" images to improve the chances of successful recognition. Techniques for the preprocessing includes:

- De-skew – If the document was not aligned properly when scanned, it may need to be tilted a few degrees clockwise or counterclockwise in order to make lines of text perfectly horizontal or vertical.
- Despeckle – remove positive and negative spots, smoothing edges
- Binarization – Convert an image from color or grey scale to black-and-white (called a "binary image" because there are two colours). In some cases, this is necessary for the character recognition algorithm; in other cases, the algorithm performs better on the original image and so this step is skipped. rgb2gray
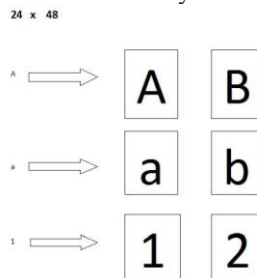- Convert RGB image or colormap to grayscale
  Syntax
  I = rgb2gray(RGB)
  newmap = rgb2gray(map)
  Description
  I = rgb2gray(RGB) converts the true colour image RGB to the grayscale intensity image I. rgb2gray converts RGB images to grayscale by eliminating the hue and saturation information while retaining the luminance.
  newmap = rgb2gray(map) returns a grayscale colormap equivalent to map
- Line removal – Cleans up non-glyph boxes and lines
- Layout analysis or "zoning" – Identifies columns, paragraphs, captions, etc. as distinct blocks. Especially important in multi-column layouts and tables.
- For Line and word detection Establishes baseline for word and character shapes, separates words if necessary.
- For Script recognition In multilingual documents, the script may change at the level of the words and hence, identification of the script is necessary before the right OCR can be invoked to handle the specific script.
- Character isolation or "segmentation" – For per-character OCR, multiple characters that are connected due to image artifacts must be separated; single characters that are broken into multiple pieces due to artifacts must be connected.
- Then the resizing of the each character is done by normalizing the aspect ratio and scale to 24 x 48.



       Segmentation of fixed-pitch fonts is accomplished relatively simply by aligning the image to a uniform grid based on where vertical grid lines will least often intersect black areas. For proportional fonts, more sophisticated techniques are needed because whitespace between letters can sometimes be greater than that between words, and vertical lines can intersect more than one character. Preproccesing makes it easy to recognize the given data

## VI. Template Matching Algorithm

       The template matching algorithm has been fully implemented and tested. An input character is Size normalized to a 24 x 48 grid and compared by distance to a set of size-normalized prototype determined. if the letter from the image is matched with the letters stored in database up to maximum percentage the that letter is get selected . Up to 18,000 prototypes have been used at one time with this technique. Experiments have shown that performance steadily improves by adding more number of fonts and data sizes in specified library .

to recognize the large number of variations in hand printed text are better represented as some of the more obscure prototypes are added to the training data.

The performance of this algorithm has been determined with a training set of 18,000 characters and a test set of 2,000 characters. The results of this analysis are shown in Table 1. The percentage correct is shown for a given number of classes. The error rate is 100 minus this value. It is seen that the technique is 95.8 percent correct at guessing the input is among four of the 40 classes.

The M classes or the M' prototypes that most closely match an input character are then The data base that is currently being used for these experiments consists of about 10,000 in printed characters. These were extracted from about 2000 handwritten postal addresses that were the existing methods and systems for optical character recognition provide more reliability of the recognition of the given texts with high and medium print quality. A small number of errors in long texts is usually not a serious problem one does not notice them at all or corrects them easily. However such systems are not always able to cope with the task of characters recognition in industrial systems, for example, while recognising serial numbers and inscriptions on components, products, packing etc. The main requirements in this class of problems are reliability and stability, since even single errors in recognition of relatively short inscriptions may produce a serious problem. Algorithms which are used in industrial systems should be stable to different kinds of defects that originate from displacement or deformation of the object, distortion of the image acquired from the camera, image In the below discussion the learned set is in xml format. This learned set is basically coordinates related information which will be explained in below article. Noises, colour changes under different lighting or pollution etc. In these cases the algorithms for recognition of the printed text give quite poor results.in some time we got errors while detecting the characters whichare looks same as with the other characters like 'o'is looks same as number zero '0'etc.

OCR is generally an "offline" process, which analyzes a static document. Handwriting movement analysis can be used as input to text recognition. Instead of merely using the shapes of glyphs and words, this technique is able to capture givenimage of data, and convert it to sound by using simple technique like template matching.s character recognition", "dynamic character recognition", "real-time character recognition", and "intelligent character recognition".

## REFERENCES

[1] Kailash S. Sharma,A. R. Karwankar, Dr. A.S.Bhalchandra," Devnagari Character Recognition Using Self Organizing Maps" ICCCCT'10
[2] http://www.heatonresearch.com/articles/series/1
[3] R.M.K. Sinha, and Veena Bansal, "On Automating trainer for construction of prototypes for Devnagari text recognition", Technical report TRCS-95-232, IIT Kanpur, India 1995.
[4] http://en.wikipedia.org/wiki/Handwriting_recognition
[5] R.M.K. Sinha, and Veena Bansal, "On Devanagari documentation processing", IEEE International Conference on Systems, Man and Cybernetics, Vancouver, Canada 1995.