

A Methodology For Enhancement And Reconstruction of Document Using Haar Wavelet Technique

Miss. Ankita S. Nathe¹, Dr. Chetan J. Shelke²

^{1,2}Dept of Computer Science and Engineering

^{1,2}P. R. Pote College of Engineering and Management, Amravati, India

Abstract- Character perception is an unique way to diagnose the unidentifiable character into understandable text. Nowadays, with the increasing urge of getting information digital media has been roped in because of its convenient use. Hence the printed media needs to get into digital format but the printed text needs to get recognized by the computer this difficulty is solved by the pattern matching technique used in this methodology this uses haar wavelet method to recognize the character ,decompose it and frame it into understandable format keeping its prototype in mind this method can turn boon for the unidentifiable text also it can overcome the problems of printed text character perception system help human ease and reduce their jobs of manually handling and processing of document.

Keywords- Character perception, haar wavelet, pattern matching, printed text.

I. INTRODUCTION

Recognizing and analyzing character is one of most difficult task. Character recognition methods aims at identifying segmented characters in images of printed text. The objective is to develop an offline character recognizing system to recognize Machine Printed Characters from the document images, which comprises of text in machine printed format, which would be convert editable form.

Documents are valuable resource of information. Lots of documents are available in libraries in the form of hard copy of any book, magazine, newspaper which is having great amount of valuable information. So it is needed to convert this documents into which will help to keep it forever that is nothing but transforming it into digital form. Many documents suffer from several factors such as low paper quality, degradation, lack of standard alphabets, stains, noise, dense and arbitrary layout, typesetting imperfections, low print contrast and fonts etc.

These factors deny the operation of stereotyped image reorganization method to documents. So predominantly it is necessary to reduce the noise and improve the quality[1]. To remove the noise from the documents [2] appropriate method is applied. Further, there is a need to apply appropriate image enhancement techniques [3] to enhance the quality of these documents. Image Enhancement stage improves the clarity of images for human viewing, removes blurring and reduces noise and increases contrast for revealing more details [4]. To maintain the original document consistently[5,6] it is essential that these documents are transformed into digital media. One of the objectives of this paper is to preserve the documents forever.

II. LITERATURE SURVEY

Enhancement of old images and documents by digital image processing techniques were proposed by, Preeti et al., 2015. Binarization techniques applied to remove the noise and improve the quality of the documents. Specialized processing is required to these document images for removing background noise in order to become more legible. A hybrid binarization approach is proposed in this paper for improving the quality for the old documents. Combination of global and local thresholding techniques are used for the same. Initially, a technique named global thresholding is applied to the whole image. The image area that still has background noise is detected and the technique is again re-applied to each area separately [7].

Text Line Detection in Corrupted and Damaged Historical Manuscripts is proposed by Rabaevet al.2013. Method grouped text lines by analyzing the evolution maps of connected components. A sweep line moved from left to right is further used to check whether elements lie in the same line. However, the method can only detect lines of equal-size texts which are chosen in their dataset. Method is found to be powerful to detected characters in torn and damaged Manuscripts [8].

Binarization-Free Text Line Segmentation for Historical Documents Based on Interest Point Clustering proposed by Angelika et al., 2012 . Interest points representing parts of characters are extracted from gray-scale images. Next, word clusters are identified in high-density regions and touching components such as ascenders and descenders are separated using seam carving. Finally, text lines are generated by concatenating neighboring word clusters, where neighborhood is defined by the prevailing orientation of the words in the document [9].

III. PROPOSED SYSTEM

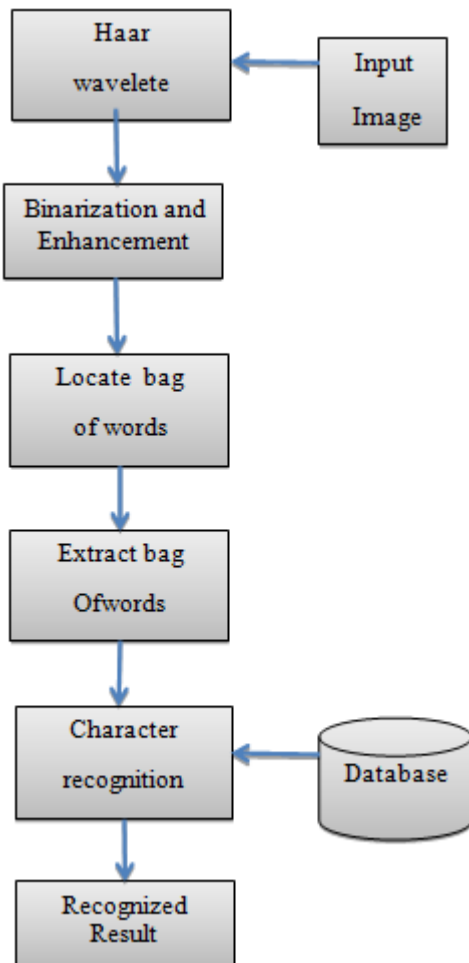
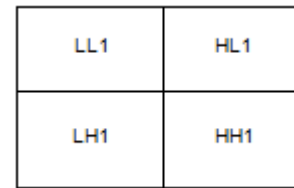


Fig 1:- Block diagram of character recognition system

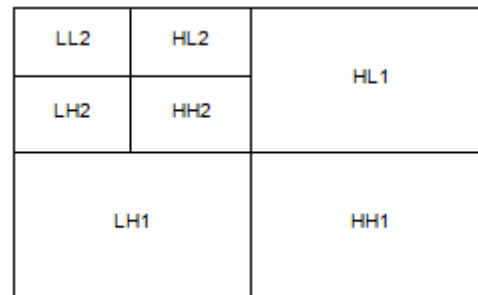
1. Image Acquisition: First of all raw unprocessed image is provide as input, it can be any text that contain preferred language. Initially images are taken from internet. Which are later available in folder of the system that to apply the following procedures.

2. Image Decomposition(Hear wavelete): -The wavelet decomposition is nothing but to transform the original raw image into several components with single low-resolution

component as shown in Figure (a) and (b). Figure (a) Shows first level decomposition and figure (b) shows second level decomposition. A higher level of decomposition is obtain by repeating the same pattern of splitting it into number of rows and columns.



(a) First level Decomposition



(b) Second level Decomposition

3. Image Binarization and Enhancement: This consists of the processes for Binarization and Enhancement. If given images are in RGB scale. before further processing convert it into a grayscale image and then into a binary image with appropriate thresholding will get an image without loss of information.

In binarization, it takes the original scan image matrix as input and processes each pixel of it in such a way that, if the pixel value is greater than 200 it assumes it as a white pixel and change it to value 1. Similarly, if the pixel value is less than 200 it assumes it as a black pixel and changes it to 0.

Table 1. Some important pre-processing operations Processes Description

Process	Description
Binarization	It is use to separates image pixels as text or background.
Noise Reduction	Noise Reduction Better improvements of image acquisition
Thresholding	Separating information from its background.

Enhancement takes the binary image matrix as input, scan it from top row to bottom row for every pixel of each row. If it finds any pixel in a row then it assumes it as the top of the image. After that, it again scans all pixel of every row and if it

finds any pixel then it will increment the row number. At the end of the image, it will save the last row where it last found a black pixel in the image, and assume that row as the bottom of the image. Similarly, by the same process, it will detect the right and left the edge of that image. It will scan it from left most column to right most column for every pixel of each column. If it finds any pixel in a column then it assumes it as the leftmost edge of the image. After that, it again scans all pixel of every column and if it finds any pixel it will increment the column number. At the rightmost end of an image, it will save the last column number where it found the last black pixel in the image, and assume that column as the rightmost edge of the image

4. Locate bag of words:-In the previous stage all images are decomposed. As character area are shown with white pixel and background with black. If character area of white pixel is greater than 90% then consider that there is character otherwise not.

5. Extract Image: The left edge, right edge, top edge and bottom edge of the entire image are detected and cropped.

6. Character Recognition: The difficult stage is to recognize the character, here the matching operation is going to perform on the image which is provided as input and another one is the image from the dataset. Before matching, each image from the dataset is taken and binarized. Then both the acquired input character image and the binarized dataset image is resized in a particular size. Then it performs a pixel by pixel matching operation between the input image and all the stored image in the dataset. It does the same for all the dataset image and finally, the one with the maximum number of matched pixel is selected and accordingly the appropriate character is detected from the data.

7. Character Printing: The character with the maximum value of the correlation coefficient is then printed on the screen in the form of an image.

Algorithm (Haar Wavelet)

1. Start
2. Input an Image I
3. Input Decomposition Level L
4. Resize an Image with size of Width & Height as power of 2.
5. Initialize $i=1$, $Width_{new}=0$, $Height_{new}=0$.
6. While(1)

$Width_{new_end}=Width_{new}+Width/L;$

$Height_{new_end}=Height_{new}+Height/L;$

$Di =$ Decompose an Image I ($Width_{new}, Width_{new_end}, Height_{new}, Height_{new_end}$)
 Save Di
 $i=i+1;$
 $Width_{new} = Width_{new_end};$
 $Height_{new} = Height_{new_end};$
 End

7. Save Result

8. Stop

IV. CONCLUSION

As an overall scenario of the technique is for solving the problem of offline character recognition. This system will be developed by using the technique which is Haar wavelet based pattern matching approach to recognize the character image. Additionally interface of system is user friendly because of the steps provided in the system for character recognition. The future scope includes recognition of cursive handwriting and online character recognition with increase accuracy.

REFERENCES

- [1] Shenbagavadivu S and Devi R, "An investigation of noise removing techniques used in spatial domain image processing", *International Journal of Computer Science and Information Technology*, vol. 2, no. 7, pp. 198–203, 2013.
- [2] Prabhdeep S and Arora A, "Analytical analysis of image filtering techniques", *International Journal of Engineering and Innovative Technology (IJEIT)*, vol. 3, no. 4, pp. 234–237, 2013.
- [3] Shakair K and Mahmud J., "Salt and pepper noise detection and removal by tolerance based selective arithmetic mean filtering technique for image restoration", vol. 8, no. 6, pp. 1234–1239, 2008.
- [4] P. Janani, J. Premaladha and K. and S. Ravichandran, "noise removal techniques: A review", *Indian Journal of Science and Technology*, vol 8, No. 22, pp. 1-5, 2015.
- [5] Papiya C, "Histogram equalization by cumulative frequency distribution", *International Journal of Scientific and Research Publications*. vol. 2, no. 7, pp. 1–4, 2012.
- [6] Ramkumar M and Karthikeyan B, "A survey on image enhancement methods", *International Journal of Engineering and Technology (IJET)*, vol. 5, no. 2, pp. 960–1012, 2013.
- [7] P. Kale, G. Phade, S. Gandhe, and P. A. Dhulekar, "Enhancement of old images and documents by digital image processing techniques," *2015 International*

Conference on Communication, Information & Computing Technology (ICCICT), 2015.

- [8] I. Rabaev, O. Biller, J. El-Sana, K. Kedem, and I. Dinstein, "Text Line Detection in Corrupted and Damaged Historical Manuscripts," *2013 12th International Conference on Document Analysis and Recognition*, 2013.
- [9] A. Garz, A. Fischer, R. Sablatnig, and H. Bunke, "Binarization-Free Text Line Segmentation for Historical Documents Based on Interest Point Clustering," *2012 10th IAPR International Workshop on Document Analysis Systems*, 2012.