

# Modified Method of Document Text Extraction from Document Images Using Haar DWT

Navjot Kaur

**Abstract**— This paper extends the technique used for Document Text Extraction from Images using 2-D Haar Wavelet. The discrete wavelet transform is a very useful tool for signal analysis and image processing, especially in multi-resolution representation. It can decompose signal into different components in the frequency domain. Two-dimensional discrete wavelet transform (2-D DWT) decomposes an input image into four sub-bands, one average component (LL) and three detail components (LH, HL, HH). The multi-resolution of 2-D DWT has been employed to detect edges of an original image. We select an appropriate threshold value and preliminarily remove the non-text edges in the detail component sub-bands. Then we use the logical AND operator to further removes the non-text regions. Another idea of removing the large size area in the image is merged with this idea to eliminate the non-text region from Document Images.

**Keywords**— Average component, Detail components, Document text, DWT, Multi-resolution of 2-D DWT, Non-Text Edges, Sub-band images, Text extraction, 2-D Haar Wavelet

## 1 INTRODUCTION

Large amounts of information are embedded in images which are often required to be automatically recognized and processed. Texts in images contain useful information which can be useful to fully understand images. Text recognition from document images receives a growing attention because of potential applications in content based indexing, archiving documents.

The term *document* is no longer confined to scanned pages and any camera based image can be subject to operations like text information extraction (TIE) for applications such as optical character recognition (OCR), image/video indexing, mobile reading system for visually challenged persons etc.

Data capture of documents by optical scanning or by digital video yields a file of picture elements, or pixels, that is the

raw input to document analysis.

Some examples of Document Images are as shown in the figure:



fig 1. Document Image example 1

- Navjot Kaur has completed masters degree program in computer engineering from Punjabi University, Patiala (Punjab), India. E-mail: jyoti\_navi2k@yahoo.co.in

## Corporation takes over Victoria Public Hall

Proposes to restore the 125-year-old building after the elections

By S. H. Ramakrishnan

**CHENNAI:** With effect from April 1, 2009, the Chennai Corporation has taken over the Victoria Public Hall that is situated near the Ripon Building.

With the Board of the Victoria Public Hall Trust, resulting to hand over the hall and other related property to the civic body at a meeting held last month, it has paved the way for the Corporation to renovate the property.

Corporation sources said that the board members who met on March 11 unanimously resolved to hand over the property including assets, savings and deposits to the civic body.

The civic body proposes to restore the over 125-year-old building, which was designed as a town hall. The civic body would call for tenders after the election and the restoration work is expected to commence in July.

The Municipal Corporation of the city of Madras had leased out about 0.14 acres in the People's Park for 99 years to the Trustees of the Victoria Public Hall Trust



A view of the Victoria Public Hall in Chennai. - PHOTO: M. VEDHAN

beginning April 1, 1886, for the hall. The hall, built in 1887, hosted lectures, balls and stage performances. The lease rent was eight annas a ground or 16.28 for the property. The lease expired in 1982 and as the civic body did not want to extend the lease, a legal battle ensued. Now as a compromise, has been reached in the matter, a petition to that effect would be submitted in the court. The Corporation is also likely to take action against sub-leases.

fig 2. Document Image example 2

Text extraction is a critical and essential step as it sets up the quality of the final recognition result. It aims at segmenting text from background, meaning isolated text pixels from those of background. A text extraction system usually assumes that text is the major input contributor, but it also has to be robust against variations in the detected text's bounding box size. A very efficient text extraction method could enable the use of commercial OCR without any other modifications.

## OBJECTIVES

To extract text from Document images using 2-D Haar Wavelet and by eliminating large size areas in the image. In case we have larger area components in the image, we can get better result.

## RESEARCH METHODOLOGY

The edges detection is accomplished by using 2-D Haar DWT and some of the non-text edges are removed using thresholding. Afterward, we connect the isolated candidate text edges in each detail component sub-band of the binary image. Although the color component may differ in a text region, the information about colors does not help extracting texts from images. If the input image is a gray-level image, the image is processed directly starting at discrete wavelet transform. If the input image is colored, its RGB components are combined to give an intensity image Y as follows:

$$Y = 0.299R + 0.587G + 0.114B \quad (1)$$

Image Y is then processed with discrete wavelet transform and the whole extraction algorithm afterward. If the input image itself is stored in the DWT compressed form, DWT operation can be omitted in the proposed algorithm.

To remove the large size area in the image, the outline of the suggested idea (MATLAB code) is as follows:

```
[CC NOB] = bwlabel(final, 8);
S = regionprops(CC, 'Area');
Binary=final;
Stats = regionprops(CC, 'Area', 'pixelidlist');
Cleaned = binary;
For p=1:length(stats)
    Ara(p)=stats(p).Area;
End
Area=sort(ara,'descend');
MIN_AREA=area(2)+10;
For region = 1 : length(stats)
    If stats(region).Area > MIN_AREA
        Cleaned(stats(region).pixelidlist) = 0;
    End
End
Cl=imresize(cleaned,[512 512]);
B=[0 1 0; 1 1 1; 0 1 0];
K=imdilate(IM2,B);
J = immultiply(cl,k);
B=[1 1 1; 1 1 1; 1 1 1];
K=imdilate(J,B);
J = immultiply(IM2,k);
Final = bwareaopen(final,20);
```

where

- `bwlabel(final,8)` returns matrix CC, of the same size as final, containing labels for the connected objects in final having 8-connected objects.
- `regionprops` measures a set of properties for each connected component (object) in CC, which is a structure returned by `bwconncomp`.
- `length(stats)` finds number of elements along the largest dimension of an array.
- `B` is the dilation operator.
- `imdilate(IM2, B)` dilates the gray-scale, binary, or packed binary image IM2.
- `Immultiply(cl, k)` multiplies each element in array cl by the corresponding element in array k.
- `bwareaopen(final, 20)` removes from a binary image all CCs that have fewer than 20 pixels, providing another binary image.

## Results & Discussion

The experimentation of the proposed algorithm was carried out on a data set consisting of different document images. Currently the data set consists of 10 images (All images are given in the Appendix A). We tried the implemented technique on a set of test images and get the results as follows:



fig3. Original image



fig4. Resultant image obtained after applying AND operation on the dilated horizontal, vertical and diagonal sub-band images

We will check the performance of the implemented technique using the following Statistical measures of the performance.

- 1) Sensitivity/ Recall rate: Sensitivity relates to the test's ability to identify positive results.

$$\text{sensitivity} = \frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$$

- 2) Specificity: Specificity relates to the ability of the test to identify negative results.

$$\text{specificity} = \frac{\text{number of true negatives}}{\text{number of true negatives} + \text{number of false positives}}$$

- 3) Precision: Precision is defined as the proportion of the true positives against all the positive results (both true positives and false positives)

$$\text{precision} = \frac{\text{number of true positives}}{\text{number of true positives} + \text{false positives}}$$

- 4) F-measure: The F-measure can be used as a single measure of performance of the test. The F-measure is the harmonic mean of precision and recall.

$$F = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- 5) Accuracy: The accuracy is the proportion of true results (both true positives and true negatives).

$$\text{accuracy} = \frac{\text{number of true positives} + \text{number of true negatives}}{\text{number of true positives} + \text{false positives} + \text{false negatives} + \text{true negatives}}$$

where

- True positive: wrong is correctly diagnosed as wrong
- False positive: right incorrectly identified as wrong
- True negative: right is correctly identified as right
- False negative: wrong incorrectly identified as right

## FINDINGS

Measures are obtained in case of the test images (given in the Appendix-A). The average of those measures is as follows.

**Table1. Various measures obtained in case of 10 test images (All images are given in Appendix A)**

Image/ Measures	Recall- rate (%)	Specificity (%)	Precision (%)	Accuracy (%)
1	72	88.67	86.36	80.56
2	52	74.81	27.65	71.25
3	38	75	52.85	59.64
4	51.5	88.21	53.96	80.44
5	65.9	87.68	63.04	82.41
6	56.5	60.75	37.58	59.52
7	75	86.80	57.53	84.53
8	41.05	88.96	70.90	70
9	56.25	58.39	24	57.98
10	41.09	78.75	59.40	62.53

Taking all test images into consideration (All images are given in Appendix-A),

- Average Recall Rate = 54.87%
- Average Specificity= 78.80%
- Average Precision Rate= 53.327%
- Average Accuracy= 70.886%
- Average F-measure=54.090%

## RECOMMENDATIONS/SUGGESTIONS

For the procedure to be effective, a priori knowledge about the structure of the page is necessary. This technique is therefore particularly useful when the layout is constrained, such as is often the case when considering pages from scientific journals.

## CONCLUSIONS

We have implemented an effective document text extraction method based on the fact that in text regions, horizontal edges, vertical edges and diagonal edges are mingled together while they are distributed separately in non-text regions. Larger areas are detected to ease the method and try removing the non-text regions which are left even after the above processings.

Bottom-up technique merge evidence at increasing scales to form, e.g., words from characters, lines from words. Actually the processing of document image segmentation and classification comes under an OCR pre-processor step. The text blocks that are detected by this technique are used as an input to the OCR system.

## SCOPE FOR FURTHER RESEARCH

Despite the many efforts spent on the subject, there is still much room for improvement in document segmentation techniques, which is the key factor to improve the overall perfor-



mance of an automatic reading/processing system.

## APPENDIX A

### List of Test images



fig5



fig6



fig7



fig8



fig9



fig10



fig11



fig12



fig13



fig14

## List of Resultant images



fig15. Result of fig.5



fig16. Result of fig.6



fig17. Result of fig.7

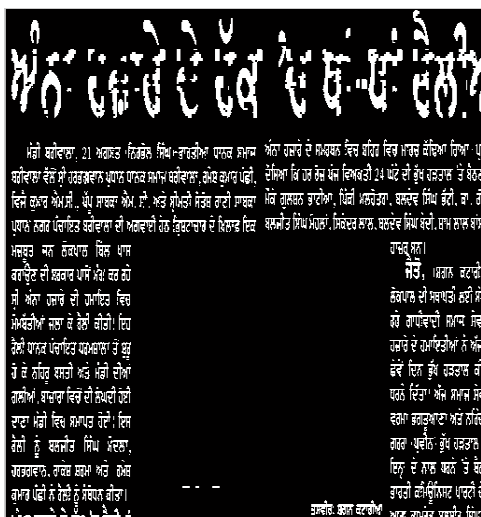


fig18. Result of fig.8

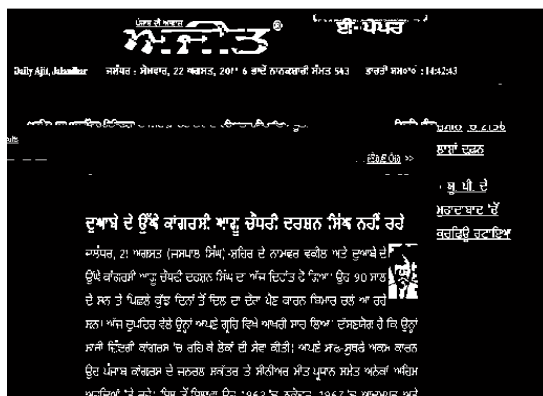


fig19. Result of fig.9



fig20. Result of fig.10

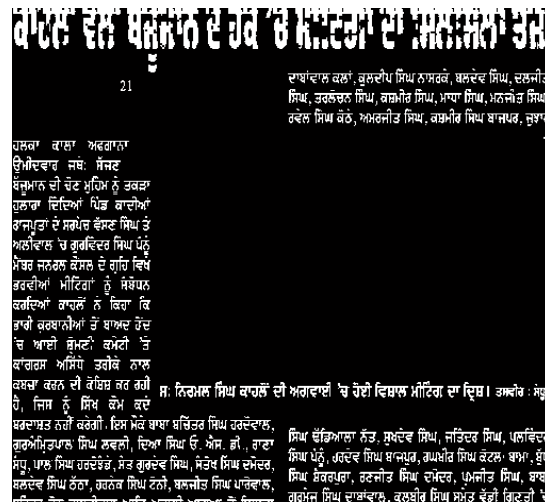


fig21. Result of fig.11

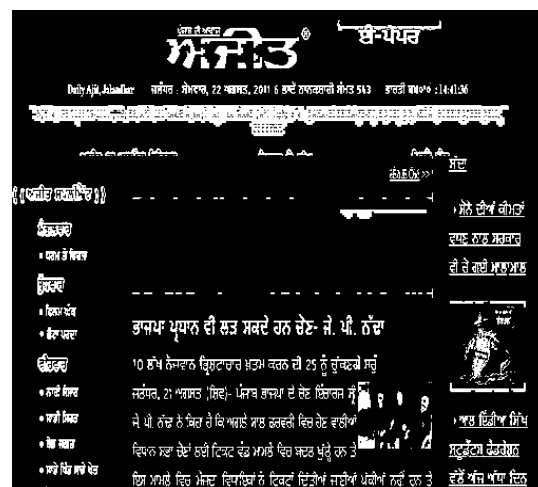


fig 22. Result of fig.12



fig23. Result of fig.13



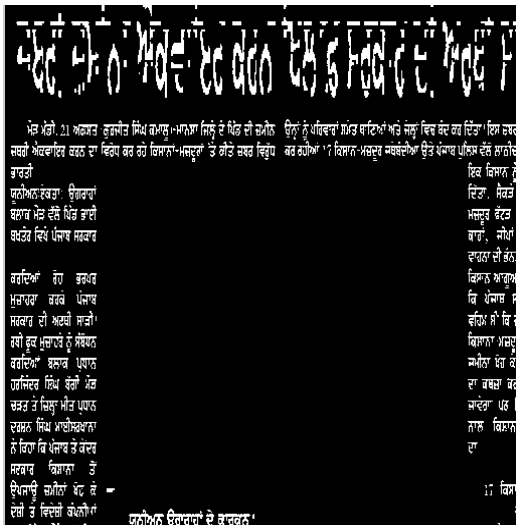


fig24. Result of fig.14

## ACKNOWLEDGMENTS

First of all, I would like to express my deep sense of respect and gratitude towards my guide Dr. Rajesh Kumar Bawa, Professor, Department of Computer Science, Punjabi University, Patiala, who has been the guiding force behind this work. I am greatly indebted to him for his constant guidance, useful suggestion and sustained encouragement throughout the work.

I also wish to acknowledge valuable interaction i've had with my other teachers of the department. Thanks are also due to all of my lab mates, from whom I learned a lot.

Finally, my parents... I am endlessly grateful to my parents, for giving me the opportunity to open my eyes in one of the most beautiful planets I have ever known. I would like to express my sincere thanks to the almighty who kept me motivated to do some purposeful work.

## REFERENCES

[1] S.Audithan, RM. Chandrasekaran (2009), "DOCUMENT TEXT EXTRACTION FROM DOCUMENT IMAGES USING HAAR DISCRETE WAVELET TRANSFORM", European Journal of Scientific Research ISSN 1450-216X Vol.36 No.4 (2009),

pp.502-512.

[2] Shyama Prosad Chowdhury, Soumyadeep Dhar, Amit Kumar Das, Bhabatosh Chanda, Karen mcmenemy (2009),"ROBUST EXTRACTION OF TEXT FROM CAMERA IMAGES", ICDAR '09 Proceedings of the 2009 10<sup>th</sup> International Conference on Document Analysis and Recognition.

[3] Ujjwal Bhattacharya, Swapan Kumar Parui, Srikanta Mondal (2009), "DEVANAGARI AND BANGLA TEXT EXTRACTION FROM NATURAL SCENE IMAGES," icdar, pp.171-175, 2009 10th International Conference on Document Analysis and Recognition.

[4] Keechul Jung, Kwang In Kim and Anil K. Jain(2004), "TEXT INFORMATION EXTRACTION IN IMAGES AND VIDEOS: A SURVEY", The journal of the Pattern Recognition society.

[5] G. Rama Mohan Babu, P. Srimaiyee, 3A. Srikrishna(2005-2010), "TEXT EXTRACTION FROMHETEROGENOUS IMAGES USING MATHEMATICAL MORPHOLOGY", Journal of Theoretical and Applied Information Technology.

[6] S. A. Angadi, M. M. Kodabagi," A TEXTURE BASED METHODOLOGY FOR TEXT REGION EXTRACTION FROM LOW RESOLUTION NATURAL SCENE IMAGES" International Journal of Image Processing (IJIP) Volume(3), Issue(5)

[7] H. Tran, A lux, H.L. Nguyen T. And A. Boucher(2005)," A NOVEL APPROACH FOR TEXT DETECTION IN IMAGES USING STRUCTURAL FEATURES", The 3rd International Conference on Advances in Pattern Recognition, LNCS Vol. 3686, pp. 627-63.

[8] X. Liu, H. Fu and Y. Jia.(2008)," GAUSSIAN MIXTURE MODELING AND LEARNNG OF NEIGHBOR CHARACTERS FOR MULTILINGUAL TEXT EXTRACTION IN IMAGES", Pattern Recognition, Vol. 41, pp. 484-493.

[9] P. Dubey(2006)," EDGE BASED TEXT DETECTION FOR MULTI-PURPOSE APPLICATION", Proceedings of International Conference Signal Processing, IEEE, Vol. 4.

[10] K. Subramanian, P. Natajajan, M. Decerbo, and D. Casta-



- non(2007)," CHARACTER-STROKE DETECTION FOR TEXT-LOCALIZATION AND EXTRACTION", *Proceedings of Ninth International Conference on Document Analysis and Recognition, IEEE*, pp. 33-37.
- [11] C. Mancas-Thilou, B. Gosselin(2006)," SPATIAL AND COLOR SPACES COMBINATION FOR NATURAL SCENE TEXT EXTRACTION", *Proceedings of IEEE International Conference on Image Processing*, pp. 985-988.
- [12] W. M. Pan, T. D. Bui, and C. Y. Suen(2007),"TEXT SEGMENTATION FROM COMPLEX BACKGROUND USING SPARSE REPRESENTATIONS", *Proceedings of Ninth International Conference on Document Analysis and Recognition, IEEE*, pp. 412-416.
- [13] J. Liang, D. Doermann, and H. P. Li.(2005)," CAMERA-BASED ANALYSIS OF TEXT AND DOCUMENTS: A SURVEY". *Int'l J. Document Analysis and Recognition*, 7(2-3):84-104.
- [14] D. F. Dunn and N. E. Mathew(2000), "EXTRACTING COLOUR HALFTONES FROM PRINTED DOCUMENTS USING TEXTURE ANALYSIS," *Pattern Recognition*, vol. 33, no. 3, pp. 445-463.
- [15] M. I. C. Murgui(1998), "DOCUMENT SEGMENTATION USING TEXTURE VARIANCE AND LOW RESOLUTION IMAGES," in *Proceedings of IEEE Southwest Symposium on Image Analysis and Interpretation*, Tucson, Arizona, USA, pp.164-167.
- [16] L. Clieque, L. Lombardi, and G. Mazini(1998), "A MULTI-RESTORATION APPROACH FOR PAGE SEGMENTATION," *Pattern Recognition Letters*, vol. 19, no. 2, pp. 217-225,.
- [17] K. Etemad, D. S. Doermann, and R. Chellappa(1998), "MULTISCALE SEGMENTATION OF UNSTRUCTURED DOCUMENT PAGES USING SOFT DECISION INTEGRATION," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 19, no. 1, pp. 92-96.
- [18] A. K. Jain and Y. Zhong(1996), "PAGE SEGMENTATION USING TEXTURE ANALYSIS," *Pattern Recognition*, vol. 23, no. 2, pp. 743-770.
- [19] Y. K. Ham, M. S. Kang, H. K. Chung, and R. H. Park(1995)," RECOGNITION OF RAISED CHARACTERS FOR AUTOMATIC CLASSIFICATION OF RUBBER TIRES", *Opt. Eng.*, Vol. 34, pp.102-108.
- [20] T. Sato, T. Kanade, E. K. Hughes, and M. A. Smith(1998)," VIDEO OCR FOR DIGITAL NEWS ARCHIVE", *Proc. Of IEEE Workshop on Content based Access of Image and Video Databases*, pp. 52-60.
- [21] B. Shahraray and D. C. Gibbon(1995)," AUTOMATIC GENERATION OF PICTORIAL TRANSCRIPTS OF VIDEO PROGRAMS", *Proc. Of SPIE*, Vol. 2417.
- [22] Julinda Gllavata, Ralph Ewerth and Bernd Freisleben (2003), "A ROBUST ALGORITHM FOR TEXT DETECTION IN IMAGES", *Proceedings of the 3rd international symposium on Image and Signal Processing and Analysis*.