Research Article

# Indus image segmentation using watershed and histogram projections

## Abstract

Character segmentation is the major step of document image analysis and optical character recognition (OCR). The character segmentation is necessary to detect all the character regions in the image document. The proposed method preprocesses the image document with edge detection techniques to enhance the character edges. Further, the watershed algorithm is implemented to identify the regions of the character. Also, the multiple histogram projections are used to identify the characters. The watershed regions and the multiple histogram projections are compared to analyse the actual character regions improving the accuracy of character recognition. The Proposed method is evaluated on Telugu and Indus images and has extracted the characters accurately.

**Keywords:** watershed model, multiple projections, edge detection

**Dasubabu Talari, Anupama Namburu**
Department of CSE, Acharya Nagarjuna University, India

**Correspondence:** Anupama Namburu, Department of CSE, Acharya Nagarjuna University, India, Email namburianupama@gmail.com

## Introduction

The recognition process in general involves the segmentation of text lines, words, and the characters. The Segmentation of the handwritten document is still one of the most concerned challenging problems due to complexity in handwritten text. The success of recognition thus depends on the result of character segmentation from text lines and words eliminating the background. In this work the scope is limited to segmenting characters from text lines detected by a segmentation method as character segmentation is a challenging step in the recognition process. Conventional character segmentation methods such as projection profile based methods may not work for degraded historical documents from text lines due to the absence of regular spacing between character components. Therefore, there is a need for developing a new method for segmenting such characters. Anupama et al.[1] proposed a method based on projection method for Telugu script document segmentation. The method fails to extract the characters in presence of touching limes. The character segmentation in degraded text lines like indus document images is proposed by Aladhahalli.[2] As a conventional technique for text line segmentation, global horizontal projection analysis of black pixels has been utilized.[3–6] Partial or piece–wise horizontal projection analysis of black pixels as modified global projection technique is employed by many researchers to segment text pages of different languages.[7–9] In this paper to detect the touching characters, first the edge detection techniques Sobel & Prewitts are applied and followed by watershed algorithms to obtain the character region. The multiple projections are applied to these watershed images to obtain the projections of the character regions. The water shed regions and the histogram projections are compared to extract the exact character region. This method eliminating false lines detection of characters in overlapped text lines.

Sobel mask

$$x = \begin{array}{ccc} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{array} \quad y = \begin{array}{ccc} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{array}$$

Prewitt mask

$$x = \begin{array}{ccc} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{array} \quad y = \begin{array}{ccc} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{array}$$
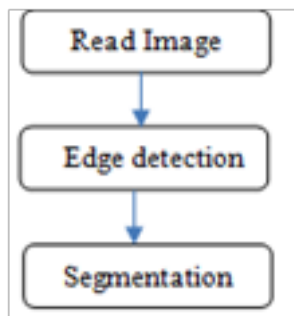
## Related works

In literature available, there are various approaches for character segmentation. However, the Indus documents are shown in Figure 1 has its background distorted and the characters are cursive in nature because of the usage of tools by hands to engrave texts like pictures on hard materials (http://en.wikipedia.org/wiki/Indus_scripts). Therefore, the decipherment of Indus documents in history remains as a research issue in the field of document image analysis. Since there are not many methods on Indus character segmentation in literature, in this section, we review the literature on the segmentation of characters from degraded, historical and handwritten document images. Most of the methods in[1,2,7,10–12] are based on projection profiles, whereas[12] use component grouping. These methods work well for plain and high–resolution texts with clear spaces between characters and watershed algorithm for Indus images. These methods segment lines using horizontal projection profile, and then use vertical projection profile to segment words or characters. Such a method scans vertically for black pixels. As a result, it may not be suitable for Indus documents. Watershed algorithm for Indus document images are considered in literature due to lack of spacing between the characters. There are methods which explore watershed algorithm for segmenting text lines.[12–15] Most of the methods use the results of morphological operations as the input for watershed to segment text lines. It is true that the performance of the morphological operation depends on the size of mask and binary output. Therefore, the methods do not perform well for complex documents such as Indus.



**Figure 1** Indus image.

## Proposed methodology

Here, a new technique which automatically identify and segment the text line region of handwritten documents (Figure 2).



**Figure 2** Steps in segmentation algorithm.

### Edge detection

Pre–processing aims to produce data that are easy for segmentation accurately. The Indus characters images are often contain degraded background. Hence, the background needs to be eliminated from that of the foreground characters. In order to do so, the Sobel edge detection algorithms are used to extract/highlight the characters from the background (Figure 3a–3c).



**Figure 3(a)** Original Images.



**Figure 3(b)** Gray scale image.



**Figure 3(c)** Edges highlighted with sober filter.

### Segmentation

Once the character edges are highlighted the water shed algorithm is applied to obtain the regions of the characters. Morphological watersheds provide a complementary approach to the segmentation of objects. It is especially useful for segmenting objects that are touching one another. To understand the watershed transform, an image $F(x,y)$ is considered as a topological surface, where the intensity values of $F(x,y)$ correspond to heights. To extract this observation, inspired by the characteristics of the watershed algorithm, namely, water flow and volume of collection water, a watershed algorithm to detect spaces between character components is proposed. The watershed algorithm finds water flow and high volume of collection of water where there is a space between two character components. These two properties work well even if any touching exists between character components. In this way, watershed algorithm helps in segmenting characters from Indus text lines by finding non–linear spacing between character components (Figure 4) (Figure 5).



**Figure 4** Image after water shed.



**Figure 5** Image indicated with color.
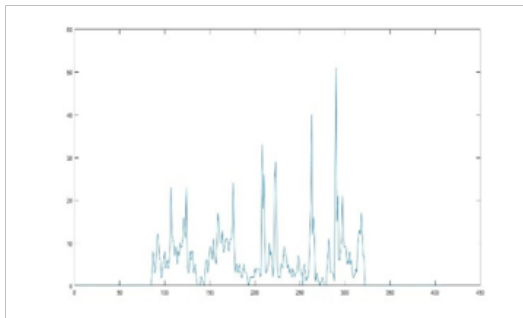
$F(x,y)$ , L= Watershed (F), where L is label matrix.

## Character segmentation

Once the water shed image is obtained the histogram projection are calculated for the watershed image. The procedure to create the histogram projections are indicated in the following steps. As each character in the Indus image after applying watershed can easily be identified, the histogram vertical projections are applied to obtain the regions (Figure 6).

Follow these steps for Word segmentation:

a.  Scan the segmented line image vertically and find the number of ON pixels in each column.

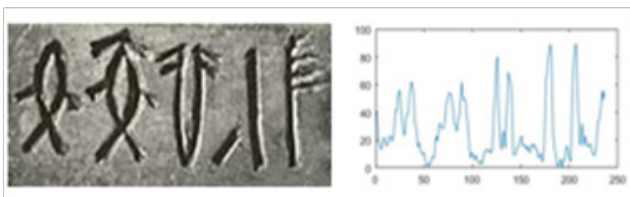b.  Plot the histogram in x direction for the ON pixel count for the image.

c. Scan the histogram projection to find first ON pixel count with zero and remember that x coordinate as x1.

d. Continue scanning the histogram projection then find lots of ON pixel counts to be non zero since the characters would have started.

e. Finally we get the first ON pixel count as zero and remember that x coordinate as x2.

f. Scan the image from x1 to x2 columns and get the segment character and Clear the X1 and X2.

g. Repeat the above steps till the end of the vertical histogram.



**Figure 6** The histogram projections of the Indus image.
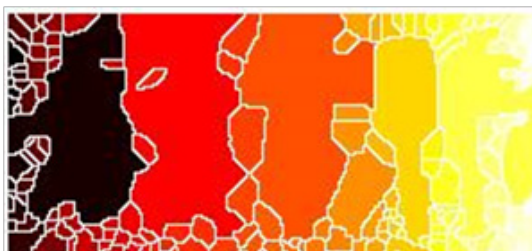
## Experimental results

The experimental results of all the segmentation steps for Indus image are shown in (Figure 7a–7g).
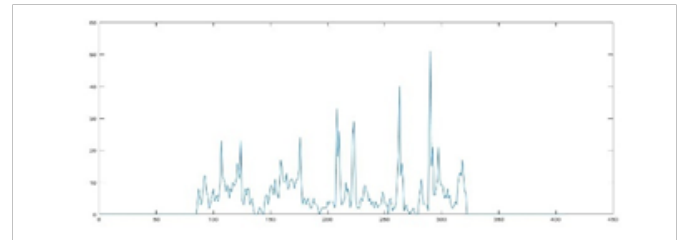


**Figure 7** Indus Image and Projection.



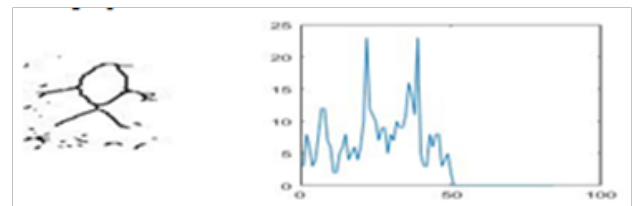**Figure 7(b):** Image after watershed algorithm.



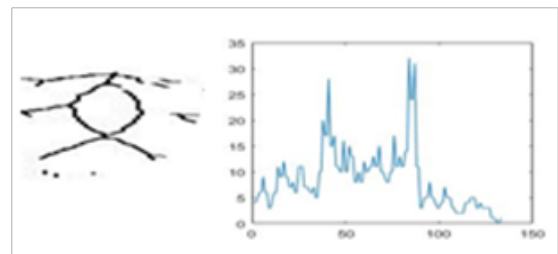**Figure 7(c)** regions of characters watershed algorithm.



**Figure 7(d)** Projections of image.



**Figure 7(e)** Characters extracted using Watershed and projection.



**Figure 7(f)** Segmentation and projection.



**Figure 7(g)** Segmentation and projection.

## Conclusion and future work

In this paper a new method for segmenting characters from text lines and degraded document images like Indus. the proposed algorithm is tested with several document images. Even though this algorithm provides robust results such as detection rate DR (98%) and Recognition Accuracy RA (98%).We have proposed the watershed model for identifying non–linear spacing between characters by exploiting catchment basin and flow of water. Experimental results and the comparisons with the existing methods show that the proposed method outperforms the existing methods in terms of recall and precision. The future work would be extending the same method for blur images and multiple touching character component images in multi scales or multi oriented environments.

## Acknowledgments

None.

## Conflict of interest

Author declares that there is none of the conflicts.

## References

1. Anupama N, Rupa C, Reddy ES. Character Segmentation for Telugu Image Document using Multiple Histogram Projections. *Global Journal of Computer Science and Technology Graphics and Vision*. 2013;13(5):1–7.

2. Aladhahalli Shivegowda Kavitha, PalaiahnakoteShivakumara ,Govindaraj Hemantha Kumar. A New Watershed Model based System for Character Segmentation in Degraded Text Lines. *AEUE–International Journal of Electronics and Communications*. 2017. p. 45–52.

3. Pal U, Chaudhuri BB. Indian script character Recognition: A Survey. *Pattern Recognition*. 2004;37:1887–1899.

4. Chaudhuri BB , Pal U. A complete printed Bangla OCR system. Pattern Recognition. 1998;31:531–549.

5. Vijay Kumar, Pankaj KSenegar. Segmentation of Printed Text in Devnagari Script and Gurmukhi Script. *IJCA: International Journal of Computer Applications*. 2010;3(8):24–29.

6. Pal U, Sagarika Datta. Segmentation of Bangla Unconstrained Handwritten Text. *Proc. 7th Int Conf on Document Analysis and Recognition*. 2003. p. 1128–1132.

7. Wong K, Casey R, Wahl F. Document Analysis System. *IBM J Res Dev*. 1982;26(6):647–656.

8. Likforman–Sulem L, Zahour A, Taconet B. Text line Segmentation of Historical Documents: a Survey. *International Journal on Document Analysis and Recognition*. 2007;9(2):123–138.

9. Pal U, Roy PP. Multi–oriented and curved text lines extraction from Indian documents. *IEEE Trans on Systems, Man and Cybernetics*. 2004;34(4):1676–1684.

10. Phan TV, Zhu B, Nakagawa M. Development of Nom Character Segmentation for Collecting Patterns from Historical Document Pages. *Workshop on Historical Document Imaging and Processing*. 2011. p. 133–139.

11. Silva C, Kariyawasam C. Segmenting Sinhala Handwritten Characters. *International Journal of Conceptions on Computing and Information Technology*. 2014;2(4):22–26.

12. Mathivanan P, Ganesamoorthy B, Maran P. Watershed algorithm based segmentation for handwritten text identification. *ICTACT Journal on Image and Video processing*. 2014;4(3):767–772.

13. Kaur A, Verma A, Ssiet, Derabassi. The Marker–Based Watershed Segmentation–A Review. *IJEIT*. 2013;3(3):171–174.

14. Shadkami P, Bonnier N. Watershed Based Document Image Analysis. *Advanced Concepts for Intelligent Vision Systems*. 2010;6474:114–124.

15. Otsu N. A threshold selection method from gray–level histograms. *IEEE Transactions on Systems, Man, and Cybernetics SMC*. 1979;9(1):1–5.