

## A Review on the Text Segmentation Techniques

*Harpreet mann, Chetan Marwaha*

M.Tech Scholar Department of Computer Engineering and Technology, Guru Nanak Dev University, Amritsar  
mann.harpreet25@yahoo.in

Senior programmer, Department of Computer Engineering and Technology, Guru Nanak Dev University, Amritsar  
cmarwaha@rediffmail.com

### ABSTRACT

This paper represents the text segmentation is an important process of image processing and understanding. Basically it is defined as the process of dividing the image into different parts of homogeneity. The aim of text segmentation is to simplify the representation of an image into something that is more meaningful and easier to understand. So the overall objective is to represent the nearest neighbor criteria for grouping components in the same line which results in clusters as well as large dataset containing variety of images. The Sobel and Laplacian for enhancing degraded low contrast pixels in it.

### Keywords

Text segmentation, Text enhancement, Text line segmentation.

### 1. INTRODUCTION

Text segmentation has been enhanced appreciably; nice of aged pieces of software including Indus continues to complicated due to the complexity. Indus files consist of signs that appear to be including attractive throughout images. Generally, all these signs are etched manually on infrequent types of surface for example pebbles while in the ancient time. For that reason, Indus piece of software discovered in close up type that was applied by people with regards to connection inside the past. Determine 1 displays some images of them files, where texting are related to animal-like images in different types such as a solo horn as well as horns. This difficulty makes the segmentation dilemma harder in addition to interesting. Due to the large variety for these files in addition to the lack of scholars in the field of epigraphy, it is difficult in order to read every one of the website programs manually the way it uses a lot of time. So that you can minimize guide book work, there is certainly an excuse for the particular digitization with the website programs in order to preserve crucial information with regard to foreseeable future study. The vast majority of figure functions placed can not be relevant to fresh photographs; rather, objects has to be already available. Also, many of the functions need which specific characters always be explained specifically a person concept in an effort to discriminate between figure in addition to non-character objects.

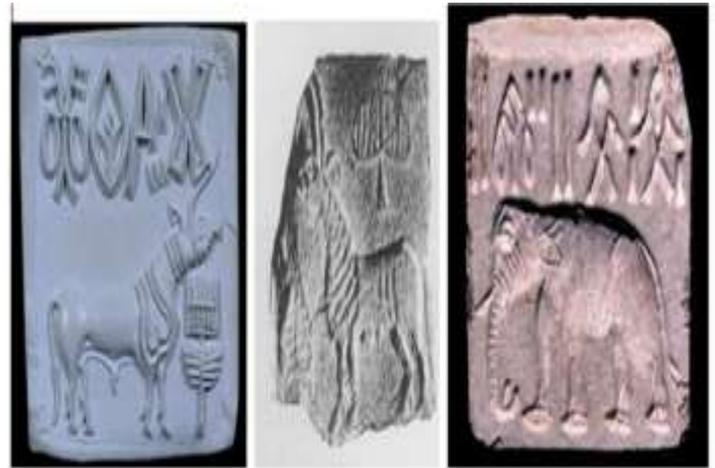


Figure 1: Sample Indus document images.

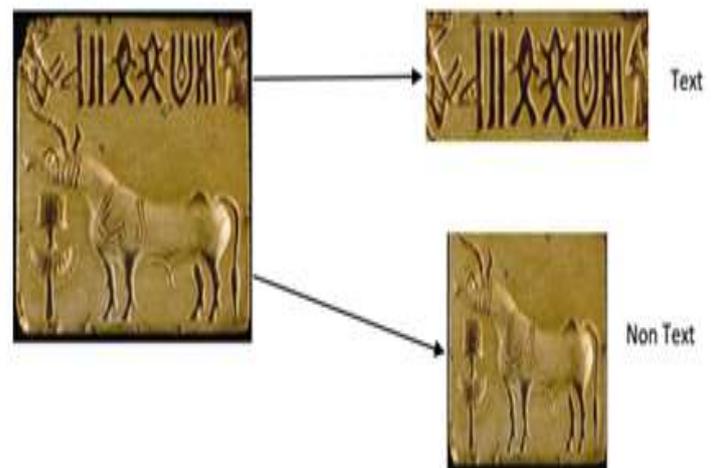


Figure 2: Illustrating text and non-text components in Indus documents.

## 1.1 CONVERSION OF RAW SCRIPT DATTO DIGITAL DATA

Producing an automatic criteria pertaining to transforming live program files so that you can a digital files includes some ways, specifically:

- 1.2.1 Text line segmentation
- 1.2.2 Word segmentation,
- 1.2.3 Character segmentation, and
- 1.2.4 Character recognition.

## 1.2 TEXT LINE SEGMENTATION

Text line segmentation is a vital stage mainly because it helps some other ways to obtain superior acceptance rates. In addition, word line segmentation is actually difficult for your record just like Indus a result of the unusual properties connected with word ingredients as well as unpredictable historical past variations. Consequently, within this work, word line segmentation through Indus texts can be on target. The vast majority of solutions are generally formulated dependant on geometrical options like feature proportion as well as measurements to get word line segmentation. Consequently, these techniques probably are not appropriate for word line segmentation through Indus record images, where by one particular can't expect uniform measurements as well as structure caused by sophisticated background. That's why, it can be come to the conclusion there is an immense setting to get developing a fresh way of segmenting word strains through Indus record images.

### 2. Text enhancement

To improve small contrast wording components. For this function, Laplacian plus Sobel procedures for the enter photograph as they are well recognized gradient procedures to reinforce the information from the photograph are considered. Because Sobel business will be the earliest buy kind, they present very good information for top contrast pixels. Hence, they boosts simply excessive contrast perimeters regarding wording parts but not the sides regarding small contrast parts as in Indus file images. To overpower this challenge, Laplacian business will be offered, which often boosts each small contrast and contrast p since this business entails another buy derivative. Aside from, your Laplacian business introduces industrial noise to get elaborate track record information. In order to maintain increased perimeters plus control track record industrial noise, intersection business from the Sobel and Laplacian business results will be performed. For example, the results of Sobel and Laplacian operations on the input images are respectively shown in Fig. 3(a) and (b), where one can notice Sobel enhances high contrast information, while Laplacian enhances both low and high contrast information along with noises. To take the advantage of both Sobel and Laplacian, we perform an intersection operation as shown in Fig. 3(c), where it is noted that only significant information is highlighted.

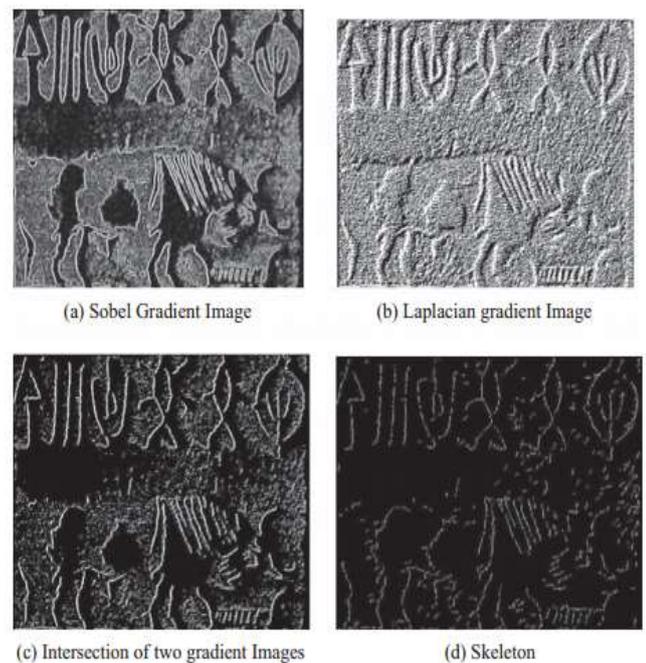


Figure 3: Intermediate results of text enhancement method.

## 3. RELATED WORK

Yuanwang et al. (2017) [1] studied the important advancement has been manufactured inside detecting textual content inside landscape images. Even so, most of state-of-the-art methods cannot effectively work while encountered blurred, low-resolution as well as small-sized texts. Many associated locations are viewed as prospects, that try and capture personality locations as many as possible. On this cardstock, most people propose to your girlfriend your story procedure, that is dependent on thorough segmentation, for you to identify textual content inside landscape images. Palaiahnakote et al. (2017)[7] proposed that Bib number/text detectors and also acceptance with Demonstration pure photos can be demanding on account of unconstrained techniques put together by background and bib range well variations. This kind of papers shows the latest multi-modal way of include body and also text message detectors ways in the novel way to accomplish results compared with the common approaches of which make use of the options associated with text messaging yet not numbers associated with hit-or-miss orientations. In the earliest period, many of us investigate HOG functions alongside through an SVM classifier to get breasts detectors in the insight image. Then the well-known Take hold of Slice method is adapted to get foreground segmentation in the breasts image. Aladhahalli et al. (2016) [3] have studied of which identity segmentation coming from wording lines within changed historical papers graphics will be hard due to intricate background non-availability of standard buildings with wording patterns. That papers proposes the latest approach determined by watershed type pertaining to segmenting personas coming from wording lines within changed historical papers images. This recommended approach filter out there noises pixels simply by looking at Sobel plus Laplacian principles with pixels, which ends up in edges of which signify wording components. A.S. Kavitha et al. (2016) [4] has studied that Text Word

segmentation via changed for the worse Old Indus screenplay pictures allows To make certain that Persona Recognizer (OCR) to attain good popularity fees with regard to Hindus pieces of software; having said that, it's complicated as a result of sophisticated historical past such images. In this paper, many of us current the latest solution to segmenting wording and non-text throughout Indus papers in accordance with the idea that wording factors usually are fewer cursive in comparison with non-text ones. To do this, many of us offer the latest mixture of Sobel and Laplacian with regard to improving changed for the worse low compare pixels. Xiaobing Wang et al. (2015) [5] has studied that text message recognition with healthy picture illustrations or photos is a scorching and also tough injury in routine acknowledgement and also personal computer vision. Along with the difficult situations with healthy picture illustrations or photos, we all offer a strong two-steps strategy with this papers based on multi-layer segmentation and higher purchase conditional arbitrary field (CRF). Offered a great knowledge graphic, the procedure divides text message by reviewing the backdrop by utilizing multi-layer segmentation, that decomposes the knowledge graphic straight into in search of layers. Pradipta Maji et al. (2015) [6] segmentation method, including judiciously the particular merits regarding rough-fuzzy computing plus multiresolution photograph research process, regarding files possessing either wording plus illustrations or photos regions. The idea assumes on of which the written text plus non-text as well as illustrations or photos regions of a certain report are generally thought to be have diverse textural properties. The *M*-band wavelet package research plus rough-fuzzy-possibilistic *h*-means are used for text-graphics segmentation problem. In this connection, a good unsupervised function collection way is travelling to pick out a set of suitable plus non-redundant capabilities regarding text-graphics segmentation problem. Ultimately, the particular rough-fuzzy-possibilistic *h*-means algorithm is familiar with handle the particular doubt dilemma regarding report segmentation. The full approach is invariant underneath the typeface measurement, brand angle, plus screenplay on the text. A effectiveness on the consist of process, plus a comparability using relevant methods, is confirmed about a set of true to life report images. Moayad at al. (2014) [7] best parts any new strategy for on line Arabic textual content acceptance using a multiple Anatomical Algorithm formula (GA) in addition to Balance Search algorithm criteria (HS). The manner is divided in not one but two development: textual content segmentation applying dominating point diagnosis, in addition to recognition-based segmentation applying GA in addition to HS. To begin with,

the pre-segmentation algorithm criteria runs on the revised dominating point diagnosis algorithm criteria in order to level a small range of items which will describes the writing skeleton. A made textual content metal framework coming from this procedure can be indicated seeing that online vector, applying 6-directional style, to minimize the effect associated with character entire body with segmentation process. Antonio et al. (2013) [8] says that wording segmentation within universal features can be recommended by way of studying the best binarization valuations for a commercial eye personality identification (OCR) system. The actual business OCR can be temporarily unveiled in addition to the details which affect your binarization for increasing the category scores. The aim of the work should be to offer the ability to quickly evaluate common textual exhibit data, to make sure that jobs that involve visual individual proof can be carried out without man intervention. The challenge to be fixed should be to recognise wording characters in which be visible on your exhibit, in addition to the color in the characters' forefront along with background. The actual papers introduces that this thresholds will be knowledgeable by way of: (a) selecting lightness and also hue part of a color knowledge cell, (b) improving the bitmaps' quality, along with (c) calculating your segmentation limit range with this cell. Hemant Misra et al. (2011) [9] says that the work involving word segmentation can be greeted originating from a subject matter acting perspective. Most of us look into the usage of not one but two not being watched subject matter designs, hidden Dirichlet portion (LDA) and multinomial blend (MM), to be able to portion your word within semantically coherent parts. Your planned subject matter model dependent approaches constantly outperform a normal baseline method upon quite a few datasets. A serious benefit from this planned LDA dependent tactic is the fact together with the portion limitations, the idea outputs the niche supply connected with every single segment. This information is involving probable use in applications like portion retrieval and discussion analysis. Even so, this planned approaches, especially the LDA dependent method, include large computational requirements. Minhua Li et al. (2010) [10] says that text contained in illustrations or photos plus online video casings supply critical indications pertaining to details listing plus retrieval. Yet it's difficult to help phase wording by illustrations or photos, especially those illustrations or photos together with complicated background. This particular document provides a different conditional arbitrary area technique, by which contextual functions are generally announced in wording segmentation.

#### 4. COMPARISON TABLE

Table 1: Comparison of Various Techniques

Ref No	Authors	Year	Technique	Features	Limitations
[1]	Yuanwang Wei, et al	2017	Exhaustive segmentation	Robustness to noise and better segmentation quality	Threshold determination to segment an image using otsu method cannot produce efficient results when gray levels between

					foreground and background do not change remarkably
[2]	R. Raghavendra	2017	Natural scene text detection	Aircraft Recognition in High-Resolution Satellite Images	Loss of essential details because of presence of noise in satellite images
[3]	Tong Lu, et al	2016	watershed model	High accuracy results and handles the problem of over segmentation efficiently.	NA
[4]	A.S. Kavitha, et al	2016	Object Detection in High-Resolution Remote Sensing Images	Minimal execution time and much better accuracy	The solution is not considered to be efficient
[5]	Xiaobing Wang et al	2015	Selective Search for Object Recognition	Use of the powerful Bag-of-Words model for recognition	Lack of fast and efficient image segmentation for medical images.
[6]	Pradipta Maji, et al	2015	Scene Labeling	A Efficient and much better segmentation results and improved performance	A Meta heuristic technique has not been considered.
[7]	Bestoun S. Ahmed, et al	2014	Optical Character Recognizer (OCR)	the multi-images help improve the detection precision of DBN than using only single-image	Loss of essential details because of presence of noise in satellite images
[8]	X. Sun ,et al	2009	Natural scene text detection	Aircraft Recognition in High-Resolution Satellite Images	Because of some architectural issues the solution is not scalable and efficient.
[9]	Sun, Hao, et al	2009	Rough-fuzzy clustering and multiresolution image analysis	Automatic Target Detection in High-Resolution Remote Sensing Images.	Lack of fast and efficient image segmentation for medical images.
[10]	E.Geoffrey ,et al	2010	An evolutionary harmony search algorithm with dominant point detection	Salient Region Detection	NA
[11]	G.-X. Zhang, et al.	2011	Multi-modal approach	the absolute maximum and minimum amount bounds of information propagation and review performance	Overhead is considered to be ignored
[12]	Xiaojun Du, et al	2009	Display text segmentation	Internet of Things (IOT); Commercial Products; Pattern; Survey;	An adaptive power control method to improve the network energy efficiency is ignored

[13]	G. Louloudis, et al	2009	latent Dirichlet allocation (LDA) and multinomial mixture (MM)	reducing the network traffic load, detect road congestion efficiently with low bandwidth consumption	The solution is not efficient and scalable
------	---------------------	------	--	--	--

## CONCLUSION

Text segmentation from degraded historical Indus script images helps optical character recognizer to achieve good reorganization rates for various scripts and it is challenging due to the complex background images. So in this paper has been discussed for segmenting text lines from degraded historical document images like Indus and the various operations for enhancing the low contrast pixels. It discusses the comparison on various techniques related to text segmentation which focuses on character segmentation from segmented text lines and character recognition. The color images contain maximum information for efficient text segmentation but still there are some issues related to the binarization as well as the consequence of diverse region on text segmentation has been ignored. So to overcome these issues we will propose an ant colony optimization based text segmentation technique for text segmentation algorithm for color images.

## REFERENCES

- [1] Yuanwang Wei, Zhijiang Zhang, Wei Shen, Dan Zeng, Mei Fang, Shifu Zhou, Text detection in scene images based on exhaustive segmentation, *Signal Processing: Image Communication*, Volume 50, February 2017, Pages 1-8
- [2] Palaiahnakote Shivakumara, R. Raghavendra, Longfei Qin, Kiran B. Raja, Tong Lu, Umapada Pal, A new multi-modal approach to bib number/text detection and recognition in Marathon images, *Pattern Recognition*, Volume 61, January 2017.
- [3] Aladhahalli Shivegowda Kavitha, Palaiahnakote Shivakumara, Govindaraj Hemantha Kumar, Tong Lu, A New Watershed Model based System for Character Segmentation in Degraded Text Lines, *AEU - International Journal of Electronics and Communications*, Available online 9 November 2016
- [4] A.S. Kavitha, P. Shivakumara, G.H. Kumar, Tong Lu, Text segmentation in degraded historical document images, *Egyptian Informatics Journal*, Volume 17, Issue 2, July 2016, Pages 189-197,
- [5] Xiaobing Wang, Yonghong Song, Yuanlin Zhang, Jingmin Xin, Natural scene text detection with multi-layer segmentation and higher order conditional random field based analysis, *Pattern Recognition Letters*, Volumes 60–61, 1 August 2015, Pages 41-47,
- [6] Pradipta Maji, Shaswati Roy, Rough-fuzzy clustering and multiresolution image analysis for text-graphics segmentation, *Applied Soft Computing*, Volume 30, May 2015, Pages 705-721,
- [7] Moayad Yousif Potrus, Umi Kalthum Ngah, Bestoun S. Ahmed, An evolutionary harmony search algorithm with dominant point detection for recognition-based segmentation of online Arabic text recognition, *Ain Shams Engineering Journal*, Volume 5, Issue 4, December 2014, Pages 1129-1139,
- [8] Antonio Fernández-Caballero, María T. López, José Carlos Castillo, Display text segmentation after learning best-fitted OCR binarization parameters, *Expert Systems with Applications*, Volume 39, Issue 4, March 2012, Pages 4032-4043,
- [9] Hemant Misra, François Yvon, Olivier Cappé, Joemon Jose, Text segmentation: A topic modeling perspective, *Information Processing & Management*, Volume 47, Issue 4, July 2011.
- [10] Minhua Li, Meng Bai, Chunheng Wang, Baihua Xiao, Conditional random field for text segmentation from images with complex background, *Pattern Recognition Letters*, Volume 31, Issue 14, 15 October 2010.
- [11] Xiaojun Du, Wumo Pan, Tien D. Bui, Text line segmentation in handwritten documents using Mumford–Shah model, *Pattern Recognition*, Volume 42, Issue 12, December 2009, Pages 3136-3145,
- [12] G. Louloudis, B. Gatos, I. Pratikakis, C. Halatsis, Text line and word segmentation of handwritten documents, *Pattern Recognition*, Volume 42, Issue 12, December 2009, Pages 3169-3183.