

## A Review of Character Segmentation Methods

Chirag Patel<sup>A\*</sup>, Atul Patel<sup>A</sup> and Dipti Shah<sup>B</sup><sup>A</sup>Faculty of Computer Science and Applications, Charotar University of Science And Technology (CHARUSAT), Changa, State Gujarat, Country India<sup>B</sup>P.G Department of Computer Science, S.P. University, V. V. Nagar State, Gujarat, Country IndiaAccepted 25 December 2013, Available online 30 December 2013, **Vol.3, No.5 (December 2013)**

### Abstract

Image to text conversion is the vital area of research for many years. Mainly, Optical Character Recognition (OCR) is used to extract characters from the image. Character segmentation is a preprocessing step for an OCR. In this paper, we have discussed different character segment methods used in various domains. Some of the methods are used for handwritten character recognition and some of the methods are for vehicle Number Plate (NP) detection. The major focus of this research is to identify the approaches that can be useful in the vehicle NP detection. After analyzing the existing character segmentation methods, the favored methods for NP detection are discussed in the conclusion section. The paper is concluded by suggesting the future scope of research in this research area.

**Keywords:** OCR, Character Segmentation, Dynamic Programming, Handwritten character recognition

### 1. Introduction

Image to text processing is the topic of research for last several years. The most common method is Optical Character Recognition (OCR) to extract text from the images. In large and complex images, it is essential to segment the image and then extract the characters by using character segmentation method. Then after the segmented character should be sent to OCR engine for the further process. The process is well depicted in Fig.1. As shown in this fig 1 (a), first an image of number plate (NP) is captured which is further processed to find the region of interest. In this figure, the purpose is to extract the characters of captured NP to detect vehicle number.



**Fig.1** Image segmentation and character segmentation (a) Original Image (b) Segmented Image (c) Image with segmented characters

By using an image segmentation process the number plate region is detected as shown in fig 1 (b). In order to identify the vehicle number each character should be clipped from the segmented NP. This task can be accomplished by using character segmentation method. The segmented characters are shown in fig 1 (c).

In the following section existing character segmentation methods are discussed, which is followed by discussion and conclusion section. The paper is concluded by suggesting future scope in the area of character segmentation.

### 2. Literature review

It is quite difficult to separate the touching characters. To segment the touching and fused Devnagri characters (Bansal & Sinha, 2002) proposed a two-pass method. In the first pass the words are easily clipped into characters which are easily separable. Then after statistical information like height and width are calculated to hypothesize whether a character box is composite. In the second pass, the hypothesized are characters are segmented for further processing. The system achieved 85% of segmentation accuracy.

In (Chen & Wu, 2009), multi-plane approach is proposed to extract text from the complex document images. The algorithm works in two stages namely - localized histogram multilevel thresholding and multi-plane region matching and assembling. Then after by using text extraction procedure textual objects are extracted and detected in the respective planes. As per the authors, this method works well in complex images with

\*Corresponding author: Chirag Patel

**Table 1** Various character segmentation approaches

| Ref                          | Method                                       | Tools/Platform                            | Language         | Accuracy (In percentage)   |
|------------------------------|--|---|------------------|----------------------------|
| (Bansal & Sinha, 2002)       | Not reported                                 | Not Reported                              | Devnagari        | 85                         |
| (Chen & Wu, 2009)            | Localized histogram multilevel thresholding  | 2.4GHz Pentium-IV personal computer using | English, Chinese | Not Reported               |
| (Choudhary, et al., 2013)    | Novel  | Not Reported                              | English          | 83.5                       |
| (Grafmüller & Beyerer, 2013) | Bayes theorem, Prior knowledge               | 2.26 GHz Intel Core 2 Duo machine         | Not Reported     | Not Reported               |
| (Lacerda & Mello, 2013)      | Skeletonization, Feature extraction          | Intel Core-i5, 2.67 GHz, 2                | English          | Not Reported               |
| (Rehman & Saba, 2011)        | Geometric features based & neural assistance | Pentium core 2 duo processor              | English          | 88.08                      |
| (Roy, et al., 2012)          | Dynamic programming                          | Not Reported                              | English          | 91.36 (Angel based method) |
| (Tan, et al., 2012)          | Non linear clustering                        | Not Reported                              | English, Chinese | 85.4 (For SegNcut method)  |
| (Verma & Lee, 2011)          | Segment confidence-based binary segmentation | Not Reported                              | Not Reported     | Not Reported               |
| (Zheng, et al., 2012)        | Not reported                                 | C#  | Chinese          | Not reported               |
| (Lee, et al., 1996)          | Multistage Graph Search algorithm            | Not Reported                              | Not Reported     | 98.02                      |
| (Ying, et al., 2010)         | Separator 's Symbols frame of reference      | Not Reported                              | Not Reported     | 99.1                       |
| (Vishwanath, et al., 2012)   | Horizontal and Vertical segmentation         | Not Reported                              | English          | 94                         |

background objects with uneven, gradational, and sharp variations in contrast and illumination conditions.

As per (Choudhary, et al., 2013) character segmentation is most crucial step for OCR. The authors proposed novel segmentation algorithm for recognition of handwritten characters. In this method segmentation points are located after thinning the word image to get the stroke width of a single pixel. The authors mention that it is challenging task to do offline handwritten character segmentation. The authors achieved 83.5% of accuracy among 200 words.

As per (Grafmüller & Beyerer, 2013), character segmentation plays vital role in OCR at various industries like food, pharmaceutical, electronics or automobile. The authors used prior knowledge based method for more reliable character segmentation. As per the authors line and character segmentation is important as wrongly segmented characters can cause lots of classification errors. The authors assume the user already specifies that text region.

Character segmentation becomes more difficult when it is used in recognition of cursive scripts. In (Rehman & Saba, 2011) a novel approach to solve this problem is proposed. Segmentation is performed by considering the analysis of geometric features and ligatures of characters. The authors mention that these features are the strong points for segmentation in cursive handwritten words. The authors used heuristics and pre-processing based four-step method, which includes processes like ligature analysis and character shape analysis. The authors achieved 88.08% of accuracy by using CEDAR (Anon., n.d.) test set of 317 words.

Another handwritten character segmentation method based on nonlinear clustering is proposed by (Tan, et al., 2012). In this approach, the entire text is segmented into strokes, then the similarity matrix of which is computed according to stroke gravities. Then after non-linear clustering is applied on similarity matrix to find out cluster label for these strokes. By using clustered labels, characters are formed by combining the strokes. The authors used two nonlinear clustering methods, namely, spectral clustering based on the normalized cut (Ncut) and kernel clustering based on Conscience On-Line Learning (COLL). The process starts by converting image into a binary image than by using similarity matrix, Ncut/COLL and other methods the image is grouped into labels and finally the characters are segmented.

A novel approach named Segment confidence-based binary segmentation (SCBS) for cursive handwritten word segmentation is proposed by (Verma & Lee, 2011). The approach is a repetitive process of fusion and segmentation of handwritten word images based on a set of suspicious segmentation points (SSPs). A variable called SegMex is maintained to limit maximum number of segments. After completing several repetitive steps mentioned in this paper the characters are successfully segmented.

A dynamic programming based approach for multi-oriented touching text character segmentation is proposed by (Roy, et al., 2012). The algorithm first segments the touching characters into primitives and then finds the best sequence of character shapes based on a dynamic programming approach using these primitive segments. The authors also performed an ad-hoc segmentation

approach of touching characters based on the prior knowledge of the number of characters in the touching string. The authors achieved the best result in this scenario but in the real world it is not possible as it is difficult to have prior knowledge about the number of characters in the image. The authors (Jia, et al., 2007) also proposed prior knowledge based method to segment characters from degraded license plate characters.

A projective method for preliminary character segmentation is proposed by (Zheng, et al., 2012). The authors proposed two methods, one is syncopation based on projection, which is a projection method of being based on pixels of characters in a vertical and horizontal direction. The authors claim that this method is very effective for Individual character cutting, nonetheless, under different conditions of aliasing character or touching symbols or multiple character separation etc. Another method is Syncopation Based on Connected Region, in which connected regions are searched in the binary image based on the algorithm discussed in this paper. The system was developed using C# language.

To extract characters from a license plate, a separator symbol's frame of reference is used by (Ying, et al., 2010). In this method, firstly the separator symbol's frame of reference is constructed based of certain features mentioned in this paper. The the character segmentation process begins with characters pre' segmentation process which uses this symbol. Then after by using vertical projection features of binary plate image the precise interval or distance between the characters is calculated. The authors carried experiment on 219 original sample images adopted in experiment is 219 with the resolution of 320 X 240.

Another hybrid approach for NP character segmentation is proposed by (Vishwanath, et al., 2012). In this approach, the preprocessing method followed by horizontal and vertical segmentation is used to segment the NP characters. The authors used Hough Transform for doing horizontal and vertical segmentation. A study of license plate character segmentation based on vertical projection is well explained by. (Xia & Liao, 2011). A survey on license plate character segmentation of video images (Yutao, et al., 2011) provides useful information about NP Segmentation and character segmentation.

Some of the other useful methods proposed by (Batuwita & Bandara, 2005), (Lee, et al., 1996), (Qi, et al., 2008), (Muralikrishna & Koti Reddy, 2011), (Yoon, et al., 2011), (Palrecha, et al., 2011), (Ban, et al., 2012) and (Syama, et al., 2012) also suggested a new way to doing character segmentation. Different character segmentation methods including tools, platform and language support are presented in Table 1.

## Discussion and Conclusion

Different character segmentation methods are discussed in this paper. There are different methods of character segmentation such as localized histogram multilevel thresholding, Bayes theorem, prior knowledge, feature extraction, dynamic programming, nonlinear clustering,

multistage graph search algorithm, segment confidence-based binary segmentation, separator symbol's frame of reference and horizontal-vertical segmentation. All these methods are very useful as a preprocessing step for the OCR. Some of algorithms based on prior knowledge and separator symbol's frame of reference might not be useful for NP segmentation as it is difficult have prior knowledge regarding vehicle NP in advance. Dynamic programming and Segment confidence-based binary segmentation (SCBS) based methods can be really useful for NP character extraction

## Future Scope

As a preprocessing part of the OCR character segmentation plays a very important part mainly in document analysis and processing and vehicle number extraction from vehicle NP identification. As none of the methods, provide 100% accuracy. The accuracy of OCR depends on the character segment method as the wrongly segmented character can cause misidentification of the character. Researchers are working in the direction to achieve 100% accuracy for character segmentation process in order to improve accuracy in OCR.

## Acknowledgements

The authors thank the Charotar University of Science and Technology for providing necessary infrastructure and resources to accomplish this research.

## References

- Anon., n.d. [Online] Available at: <http://www.cedar.buffalo.edu/handwriting/HRdatabase.html>
- Ban, K.-D., Yoon, H., Kim, J. & Yoon, Y., 2012. Blob detection and filtering for character segmentation of license plates. s.l., s.n., pp. 349-353.
- Bansal, V. & Sinha, R., 2002. Segmentation of touching and fused Devanagari characters. *Pattern Recognition*, April, 35(4), pp. 875-893.
- Batuwita, K. B. M. R. & Bandara, G. E. M. D. C., 2005. An Improved Segmentation Algorithm for Individual Offline Handwritten Character Segmentation. s.l., s.n., pp. 982-988.
- Chen, Y.-L. & Wu, B.-F., 2009. A multi-plane approach for text segmentation of complex document images. *Pattern Recognition*, July, 42(7), pp. 1419-1444.
- Choudhary, A., Rishi, R. & Ahlawat, S., 2013. A New Character Segmentation Approach for Off-Line Cursive Handwritten Words. *Procedia Computer Science*, Volume 17, pp. 88-95.
- Grafmüller, M. & Beyerer, J., 2013. Performance improvement of character recognition in industrial applications using prior knowledge for more reliable segmentation. *Expert Systems with Applications*, 40(7), pp. 6955-6963.
- Jia, X., Wang, X., Li, W. & Wang, H., 2007. A Novel Algorithm for Character Segmentation of Degraded License Plate Based on Prior Knowledge. s.l., s.n., pp. 249-253.
- Lacerda, E. B. & Mello, C. A., 2013. Segmentation of connected handwritten digits using Self-Organizing Maps. *Expert Systems with Applications*, November, 40(15), pp. 5867-5877.
- Lee, S.-W., Lee, D.-J. & Park, H.-S., 1996. A new methodology for gray-scale character segmentation and recognition. *IEEE*

- Transactions on Pattern Analysis and Machine Intelligence*, 18(10), pp. 1045-1050.
- Muralikrishna, M. & Koti Reddy, D., 2011. An OCR-character segmentation using Routing based fast replacement paths in Reach Algorithm. s.l., s.n., pp. 1-7.
- Palrecha, N. et al., 2011. Character segmentation for multi lingual Indic and Roman scripts. s.l., s.n., pp. 45-49.
- Qi, W., Li, X. & Yang, B., 2008. A Character Segmentation Method without Character Verification. s.l., s.n., pp. 581-584.
- Rehman, A. & Saba, T., 2011. Performance analysis of character segmentation approach for cursive script recognition on benchmark database. *Digital Signal Processing*, May, 21(3), pp. 486-490.
- Roy, a. P., Pal, U., Lladós, J. & Delalandre, M., 2012. Multi-oriented touching text character segmentation in graphical documents using dynamic programming. *Pattern Recognition*, May, 45(5), pp. 1972-1983.
- Syama, K. et al., 2012. Performance Study of Active Contour Model Based Character Segmentation with Nonlinear Diffusion. s.l., s.n., pp. 118-121.
- Tan, J. et al., 2012. A new handwritten character segmentation method based on nonlinear clustering. *Neurocomputing*, July, Volume 89, pp. 213-219.
- Verma, B. & Lee, H., 2011. Segment confidence-based binary segmentation (SCBS) for cursive handwritten words. *Expert Systems with Applications*, September, 38(09), pp. s 11167-11175.
- Vishwanath, N., Somasundaram, S., Baburajani, T. & Nallaperumal, N., 2012. A hybrid Indian license plate character segmentation algorithm for automatic license plate recognition system. s.l., s.n., pp. 1-4.
- Xia, H. & Liao, D., 2011. The study of license plate character segmentation algorithm based on vetical projection. s.l., s.n., pp. 4583-4586.
- Ying, H., Song, J. & Ren, X., 2010. Character segmentation for license plate by the separator symbol's frame of reference. s.l., s.n., pp. V1-438,V1-442.
- Yoon, Y., Ban, K.-D., Yoon, H. & Kim, J., 2011. Blob extraction based character segmentation method for automatic license plate recognition system. s.l., s.n., pp. 2192-2196.
- Yutao, W., Ruixia, T., Ling, M. & Gang, Y., 2011. License plate character segmentation from video images: A survey. s.l., s.n., pp. 25-30.
- Zheng, Z. et al., 2012. Character Segmentation System Based on C# Design and Implementation. *Procedia Engineering*, Volume 29, pp. 4073-4078.
- Chirag Patel** received Bachelor degree in computer applications (B.C.A) degree from Dharmsinh Desai University Nadiad, Gujarat, India in 2002 and Masters Degree in Computer Applications (M.C.A) from Gujarat University, Gujarat, India in 2005. He is pursuing PhD in Computer Science and Applications from Charotar University of Science and Technology (CHARUSAT). He is with MCA Department at Smt Chandaben Mohanbhai Patel Institute of Computer Applications, Charotar University of Science and Technology (CHARUSAT), Changa, Gujarat, India. His research interests include Information Retrieval from image/video, Image Processing and Service Oriented Architecture.
- Dr. Dipti Shah** received Bachelor degree in Science; B.Sc.(Maths), M.C.A. Degree from S.P. University, Gujarat, India. She has also received Ph.D in Computer Science, degree from S.P. University, Gujarat, India. Now she is Professor at G.H.Patel Department of Computer Science, S.P. University, Anand, Gujarat, India. Her Research interests include Computer Graphics, Image Processing, Multimedia and Medical Informatics.
- Dr. Atul Patel** received Bachelor in Science B.Sc (Electronics), M.C.A. Degree from Gujarat University, India. M.Phil. (Computer Science) Degree from Madurai Kamraj University, India. He has received his Ph.D degree from S. P. University. Now he is Professor and Dean, Smt Chandaben Mohanbhai Patel Institute of Computer Applications, Charotar University of Science And Technology (CHARUSAT) – Changa, India. His main research areas are wireless communication and Network Security