# A Technique for Character Segmentation in Middle zone of Handwritten Hindi words using Hybrid Approach

Preeti Sharma [#1], Manoj Kumar Sachan [*2]

[1] Mtech Scholar , Department of Computer Science and Engineering
Sant Longowal Institute of Engineering and Technology, Punjab
[1] *Preeti.niem@gmail.com*

[2] Associate Professor , Department of Computer Science and Engineering
Sant Longowal Institute of Engineering and Technology, Punjab
[2] *manojsachan@gmail.com*

**Abstract—** India is a country where people talk in multilingual and write in multi-script. Devanagari is one of the most popular scripts in India, which is used to write Hindi, Sanskrit, Sindhi, Marathi and Nepali Languages. This research work is performed on Hindi language. A large number of precious and essential documents are available in handwritten form, which needs to be converted into editable form. The existence of Optical Character Recognition (OCR) makes this task easier to convert handwritten text in editable form. Character segmentation is an important phase of OCR, which segment the characters from handwritten words. This enhances the accuracy of OCR system. In this paper a hybrid approach is used to segment the characters that contain single and multiple touching characters within a word. The proposed system is tested on a dataset of various handwritten words written by different writers. The dataset of proposed system contains more than 300 handwritten words in Hindi language. Accuracy of the proposed hybrid system is evaluated to 96% which is better than that of existing techniques.

**Keyword-** *OCR, Segmentation, Character segmentation, Hindi language, Touching character*

\*\*\*\*\*

## I. INTRODUCTION

Natural Language Processing is an interdisciplinary field of artificial intelligence, computer science and computational linguistics. It deals with the interactions between computers and human languages. Every aspect of NLP is used in script recognition[1,3], optical character recognition, sentiment analysis[4], [5] part of speech tagging [6], information extraction , social media analysis[7] etc. In the past years, there was small amount of data available for computation on desktop but as the time passes huge data is generated through various sources such as banks, business application, records of government, bill processing system etc, which needs processing to be converted into editable form so that they can be easily uploaded on internet or stored in computer. The existence of OCR makes this task easier by converting handwritten text to editable form. The OCR provides the facility to electronically search, edit and maintain documents. The general strategy of OCR involves techniques like preprocessing, segmentation and recognition [22,23]. Out of these, the segmentation philosophy is the spine of the general OCR process[8]. Fig 1 demonstrates various phases of OCR. Segmentation is a process that extract the printed and handwritten constituents (characters/ words) of a text image[24]. For example, while applying automatic mail-sorting process, the address must be identified and separated from other data on the envelope like stamps or company logos. The basic purpose of optical character recognition algorithms is to segment the words into isolated characters, so that they can be recognized individually.

Character segmentation is an operation that looks to divide an image of words into individual character images as shown in fig1. Handwritten Character Segmentation is one of the major issue in the field of Optical Character Recognition process, because there are a number of challenges in the segmentation of handwritten text, due to large variety of writing styles or pen-types[9].The proposed work focuses on handwritten character segmentation. The segmentation is performed on a database consisting of various types of words, carrying irregularities in writing style. In our work, hybrid approach is use to segment the handwritten words with focused on the middle region of the Hindi words. Structural features is identified to deal with the input containing the isolated, touching and broken words. Hybrid approach has proved to be an efficient approach to deal with the collected database.

The rest of the paper is arranged as follows : section 2 discusses the related work, section 3 explains the properties of Hindi text with zone description and describes the character segmentation along with its types, section 5 presents the proposed methodology for the handwritten character segmentation of Hindi words, section 6 present the details of experimental results and discussion.

_____



**Fig 1: Components of OCR**

## II. LITERATURE REVIEW

**Dipak and sharvari, 2012** [18] proposed a new technique for segmentation of touching characters based on joint point algorithm. This algorithm has used all foreground and background details for proper segmentation. The position of vertical bars was detected and then joint characters were segmented out. The focus was also given to remove the protruding points; these points were bulging pixels which unnecessarily generate joint points after thinning where no joint actually exists. These must be removed to generate accurate joint points where joints between lines actually exist. The given algorithm works for touching characters only.

**Nakesh kumar and lakhwinder,** *et al.,* **2010 [17],** have done segmentation based on structural approach for handwritten Hindi text had been implemented. All experiments were conducted on database constructed by taking handwritten data from 15 writers. For character separation, the vertical projection method was used after header line detection. The header line was identified using the horizontal projection profile. The line with maximum number of black pixels in upper 10% part of the word is considered as the header line. After that upper and lower modifiers were identified and segmented and then vertical projections were used to extract the characters.

**Naresh kumar and lakhwinder, et al., (2011) [16]** have introduced an algorithm for segmentation of half characters from consonants. The given algorithm uses the structural property of the script. It made an assumption of threshold value to identify the half characters in the word. It scanned the no of pixel successively column wise to segment the conjunct. The algorithm was tested on both handwritten and printed words and provides an accuracy of 83% and 87.5%respectively.the algorithm does not work very well.

**Vikas and Vijay, 2011**[15] have proposed, Bounded Box method for segmentation of document lines, words and characters was proposed and it was based on Pixel Histogram Approach. The global horizontal projection method was used to compute sum of all white pixels on every row and construct corresponding histogram. Both horizontal & vertical histograms for character segmentation has been performed. Unconnected vertical lines in the words are recognized as separate symbol by the algorithm that is being used. This work can be extended for compound letters that are connected at various places and handwritten unconnected compound letter segmentation.

**Sandip and Megha , 2011 [14]** The proposed work have carried out the work on scanned Devanagari handwritten data and segmenting top, bottom modifiers and fused characters. Header line is detected and removed using morphological operations like erosion, dilation, cropping etc. Upper and lower modifiers are extracted using these morphological operations. Highlighted point detection, bottom-line modifier and vertical line detection algorithms are used to identify the points of segmentation. Top modifiers were separated using Moore Neighbor Tracing technique and fused characters were identified and segmented by calculating average width. 52% accuracy had been achieved for segmenting simple and fused characters.

**Saiprakash and Renu,** *et al.,* **2012**[13] proposed new line and character segmentation method had been proposed for handwritten Hindi text. This method basically work on calculating the structural properties of the characters like height, width, size etc. An algorithm had been defined to convert it into straight line by subtracting the expected and actual header line. After that upper, lower modifiers and consonants were separated using horizontal projection profile. This technique was tested on 29 pages and achieved 93.6% accuracy for line segmentation, 98.6% for word and 89.90% for character segmentation. The proposed method fails to segment character in case of overlapped and touching lines, broken parts and touching characters.

**Vaishali and Chhaya ,(2014)**[12] proposed contour based technique for handwritten Hindi text which had been implemented. This paper mainly focuses on the segmentation of upper and lower modifier with the help of counter. Counter is used to return the position of connected point depending on the neighboring pixels to extract the upper and lower modifier, and to overcome the

_____

_____

uncertainty issue, windows are used to identify whether the character was joint or simple. The proposed method worked only in case of variable text sizes and different resolution images.

**Soumen and Ankit, (2015)**[11] In this paper, a new strategy to segment the upper and lower modifier had been proposed. It is divided into three stages In 1[st] stage, the header line is removed to focus on  upper-strip. In 2[nd] stage, the segmentation of upper modifier is performed with the help of statistical information. In 3rd stage same statistical information is used to choose the components on which segmentation is to be performed. And further  the lower modifier are separated from the middle zone. This method was tested on handwritten Hindi word images and the results were found with an average accuracy rate of 96.93 %. In case of shadow characters this method couldn't work well.

**Ankita and Neha., (2016)** [10] proposed a technique to segment the lines, words and characters based on the horizontal projection and bounding box using two MATLAB function i.e. region props and rectangle. The work was carried out on scanned document. Noise removal was done by using medfilt2 and the conversion process from grayscale to binary image was done by using Otsu's global thresholding method [7]. The method segmented the lines, words, isolated characters and modifiers but fails to segment the connected component. 90% accuracy had been achieved for segmenting the isolated characters.

It has been found that although some algorithms exists for segmentation of handwritten hindi words but following limitations exist as mentioned below :

1. These algorithms do not work on all three zones.

2. Existing algorithms do not handle multiple touching characters..

3. Existing systems does not segment words having broken, touching and overlapping characters.

### III. PROPERTIES OF HINDI TEXT

Hindi is considered as the official language of India which is used for writing the devnagari script. Devanagari is two dimensional scripts as consonants are arranged in many ways from top, bottom, left or right to make a meaningful letter. In Devanagari, text is written from left-to-right.

The proposed work is carried out on Hindi language, because it is foremost language in India. Hindi language has a rich character set having 11 Vowels, 33 consonants & 14 modifiers [19,20] . The figure 2 and figure 2.1 shows the consonants, vowels and modifiers in the set of Hindi language:

| अ | आ | इ | ई | उ | ऊ | ऋ | ए | ऐ | ओ | औ |
|---|---|---|---|---|---|---|---|---|---|---|
|   | ा | ि | ी | ु | ू | ृ | े | ै | ो | ौ |

Fig 2: Vowels and their modifiers [6]

Vowels could be formed as free or by using different diacritical imprints, i.e. modifiers or 'matras', which are created above, below, before or after the consonant or vowel. Two or more consonants join together to produce an alternate shape this is called as compound character. One exceptional feature of Devanagari script is horizontal flat line on top of all characters, called header line or 'Shirorekha'. The following table represents the consonant set of Hindi language.

| क | ख | ग | घ | ङ | च | छ | ज | ह |
|---|---|---|---|---|---|---|---|---|
| झ | ञ | ट | ठ | ड | ढ | ण | त | क्ष |
| थ | द | ध | न | प | फ | ब | भ | त्र |
| म | य | र | ल | व | श | ष | स | ज्ञ |

Fig 2.1: Consonants of Hindi language

#### A.  Zones in Hindi words

Header line of one character joins together with the header line of previous and next character to form a single word .A single word or character in Hindi language could be vertically differentiated into three sections: the upper area, the middle area and the lower area [14]. Middle area holds vowels, consonants or synthesis of both vowels and consonant and the upper & lower areas hold vowels, modifiers and their parts. The middle region is the crucial zone of a Hindi word, because it contains vowels, consonants and also a part of some modifiers. An example of handwritten word with different regions is shown below in figure 3:

_____

Fig 3:  Common matras, trips and regions in Hindi word.

*B.  Character  Segmentation*

Character Segmentation partitions the words into individual characters. The precision of character recognition relies on the segmentation. The principle challenge in the division of manually hand written text is due to wide variety of styles or pen-types [17].

According to the mode of data acquisition, character segmentation methodologies are categorized as follows [19]:

- **Online Character Segmentation** is the procedure of fragmenting handwriting that is recorded with a digitizer, according to the pen directions. It catches the data of the pen trajectory. While writing text with digitizer or light pen, the written text is moved to machine at the same moment and segmentation is being performed on that text online. Online here means the moment data is written and entered to the machine and side by side segmentation is performed.

- **Offline Character Segmentation** is the procedure of changing over the image of document into touch design by an optically digitizing unit, for example, optical scanner or camera. The handwritten text is scanned first and converted into digital image.  The segmentation is carried out on this bit-pattern information for machine-printed or written by hand content. According to the input received by the system ,content can be divided into two types  [21]

    **Machine Printed Text**: consists of text contained in books, article papers, magazines,etc.Machine printed characters are static and uniform in height and width size assuming the same font and size are used. An example of machine printed word in Devanagari  script *is shown in figure 4:*

Fig 4: Machine Printed word in Devanagari

    **Handwritten Text:** It can be further divided into two categories: cursive and hand printed script. Characters are non uniform and can vary greatly in size and style. Even the text data written by the same person can vary considerably. In this, the location of characters is not predictable. fig 5 shows the handwritten word.

Fig 5: Handwritten word in Devanagari

There are lots of irregularities in handwritten words because of variation in handwritings and presence of overlapping and incorrect characters. Figure 5 which is given above is an example of handwritten word in Devanagari.

IV. PROPOSED METHODOLOGY

To overcome the problematic issues in character segmentation and to improve the efficiency of existing segmentation strategies, a new algorithm has been implemented using MatLab tool. This hybrid technique is developed to segment the characters from Hindi text consists of isolated, touching, and broken characters. In the proposed system a hybrid approach combining the techniques of horizontal profile, vertical profile and clustering technique has been developed.

A. *Experimental database*: The proposed hybrid algorithm has been tested on a dataset of 300 handwritten words containing isolated, touching, multiple touching and broken characters from different writers

1.*Isolated characters:* Words holding the characters with legitimate spaces in between. The segmentation of these words is relatively less difficult than other sort of words. The characters in these types of words are separated by detecting the gaps within the different characters. An examples of isolated characters are shown in fig 6 below:

_____



fig6: Words containing isolated characters

*2.words with single and multiple touchings:* Touching characters may be of two types: inter-touching and intra-touching. The proposed approach will considere word for inter-touching characters only i.e the end of one characters touches the start point of other character at single or multiple ponit. When two or more characters touch each other, some pixels of one characters gets intermixed with other one, thus making cluster. The main focus of this system is to identify the cluster, the only cluster which is made by two touching characters.



fig7: single and multiple touching characters

*3.words Containing broken characters:* Some time due to style of writers and paper or pen quality characters can be broken from some parts. A few segments of the characters may be missing which represents to a single character as more than one character. The characters may be broken up in horizontal or vertical form.



fig8: broken characters

*B. Hybrid Segmentation Algorithm:*
1.Scan the image at 300 dpi using optical scanner from which word is to be used for input    purpose and set
   the threshold value of the scanned image to 200.
2.Identify the header line with the help of horizontal projection profile.
3.Convert the header line along with two rows above and below the shieorekha into white pixels.



Fig 9: removal of header line

*C. Segmentation of isolated characters:*
**Step1**:Estimate the height of the word, two parameters are used i.e height starting point (hspoint) and height
        last point (hlpoint).

$$Height = -1 * (hlpoint - hspoint)$$

**Step 2:** Compute the width of the word using the formula defined as:

$$width = 1 * (lpoint - spoint)$$ ,Span of the word is identified.
        where lpoint is last point and spoint is starting point

**Step 3:** Detect the bottom line by scanning from bottom.  Fig 9 and 9.1 represent the base line.



Fig 10: detected bottom line or base line          Fig 10.1: base line dissect the lower modifier

**Step 4**: Midpoint procedure is used to segment the simple or isolated handwritten words.
   (I)      Vacant space index value computed on the constrained such as height, bottom line and width of the word.
   (II)     Check the next and previous two pixels (i+1, i+2, i-1, i-2) to store the column values and check the presence of
           broken character.

_____

_____

**Step 5**: Vacant space between the characters is stored in array_ index.

(I)    Find the middle value of each vacant space successively between the characters till the end of the word.

$$MID = (start_{index} + last_{index}/2$$

(II)    This middle value is considered as the segmentation points for the isolated character



Fig11and fig 11.1 shows the segmentation of isolated characters

**Step 7**: Total no of characters in a word is computed to identify the presence of touching characters.

**Step 8**: Find the total number of characters (TOC) by the ratio of height and width.

$$TOC = \frac{Height}{Width}$$

**Step 9**: Compare the total no of characters with the segmentation points.

$$Number\ of\ MID = TOC + 1$$

**Step 10**: If touching character is present then number of segmentation points does not exceed the total no of characters in a single word.

**Step9:** Successively compute the gap between the middle values, if the gap exceeds the 120% (1.2times)of height , there must be a presence of touching character either it may be single or multiple touching .

**Step10:** Apply the clustering to detect the cluster in identified area of interest of the character in the middle portion.

**Step 11**: find the region of interest (cluster) between $(MID1 + 10) - (MID2 + 10)$ to find the heap of pixel i.e. cluster.

(I)    Scan each column to identify the cluster, if pixel count is found to be 10 then it is considered as touching point of the character.

(II)    Segment the touching character by leaving three columns successively.

(III)    Extract the new segmentation points

(IV)    Segment the word from all the segmentation points.

(V)    Display the output to the user.



Fig 12 and 12.1 segmentation of single and multiple touching characters



Fig 13 and 13.1 segmentation of words containing upper and lower modifier

_____

**International Journal on Future Revolution in Computer Science & Communication Engineering**
**Volume: 3 Issue: 7**

**ISSN: 2454-4248**
**01 – 10**
_____

*E. Flowchart of proposed methodology*



## V. RESULTS AND DISCUSSION

The Hybrid Segmentation Algorithm (HSA) implemented in current research work is tested on various documents containing handwritten Hindi language text. Different words are tested under four main categories as:

- Isolated Characters

_____

- Touching Characters

- Broken Characters

- Multiple Touching Characters

HSA , Explain the results in detail  as mentioned in figures including the comparison .

The following is the graph showing the efficiency of outputs on different types of input words:



Table 1 depicts the overall accuracy of the system in each category of inputs:

| Type of Input | Total No. of Words | Correctly Segmented Words | Incorrectly segmented words | Accuracy (%age) |
|---|---|---|---|---|
| Isolated | 150 | 147 | 3 | 98% |
| Touching | 50 | 48 | 2 | 96% |
| Broken | 50 | 49 | 1 | 98% |
| Multiple Touching | 50 | 44 | 6 | 88% |
| Total | 300 | 288 | 19 | 96% |

Table 2: Comparison of the existing and proposed system

| Ref. | Technique used | Type of input | Average Accuracy |
|---|---|---|---|
| B. Thakhral et. Al[1] | Cluster Identification Method | Handwritten Touching, conjunct & overlapping, characters | 94% |
| M. Kumar et. al [2] | Water Reservoir Principle, Projection profiles | Isolated and Touching characters in Gurumukhi | 93.5% |
| Proposed Work | Hybrid approach | Touching,  Multiple Touching Characters ,broken characters | 96% |

_____

_____

## VI. CONCLUSION AND FUTURE SCOPE

Segmentation of Devanagari handwritten text is complex task due to large variations in handwriting of different people. Though these challenges made segmentation very difficult but the proposed algorithm has been developed to overcome these difficulties. Issues like touching, broken and isolated characters are attempted for segmentation. Proposed hybrid approach efficiently works for the input dataset containing around 300 words of different categories and gives overall accuracy of 96% for handwritten character segmentation. Algorithm provides better results than existing methods Algorithm yields good results in each category of Hindi words. Approximately 830 consonants and about 768 consonants are extracted effectively. This system gives promising results for isolated, conjuncts and touching characters.. Some of the words from input database which are not been correctly segmented are words having skewed or uneven header lines, overlapping characters.

Table 3: Words having skewed header line and corresponding result



*Table 4: visualize the results of the words containing isolated, touching and broken words*

| Input image | Resultant image | Input image | Resultant image |
|---|---|---|---|
| **Results for segmentation of isolated characters** | | | |
|  | | | |
| **Results for segmentation of touching & multiple touching characters** | | | |
|  | | | |
| **Results for segmentation of broken characters** | | | |
|  | | | |

## REFRENCES

[1] G. Singh and M. K. Sachan, "Data capturing process for online Gurmukhi script recognition system," *IEEE Int. Conf. Comput. Intell. Comput. Res.*, pp. 518–521, 2015.

[2] G. Singh and S. Manoj, "Offline Gurmukhi Script Recognition using Knowledge Based Approach & Multi-Layered Perceptron Neural Network," pp. 266–271, 2015.

[3] G. Singh and M. Sachan, "Multi-layer perceptron (MLP) neural network technique for offline handwritten Gurmukhi character recognition," *2014 IEEE Int. Conf. Comput. Intell. Comput. Res. IEEE ICCIC 2014*, 2015.

_____

_____

[4] S. K. Singh, S. Paul, D. Kumar, and H. Arfi, "Sentiment Analysis of Twitter Data Set: Survey," *Int. J. Appl. Eng. Res.*, vol. 9, no. 22, pp. 13925–13936, 2014.

[5] S. K. Singh and S. Paul, "Sentiment analysis of social issues and sentiment score calculation of negative prefixes," *Int. J. Appl. Eng. Res.*, vol. 10, no. 55, pp. 1694–1699, 2015.

[6] P. K. Singh, S. K. Singh, and S. Paul, "Sentiment classification of social issues using contextual valence shifters," *Int. J. Eng. Technol.*, vol. 7, no. 4, 2015.

[7] S. K. Singh and K. S. Manoj, "Importance and Challenges of Social Media Text," *Int. J. Adv. Res. Comput. Sci.*, vol. 8, no. 3, pp. 2015–2018, 2017.

[8] R. G. Casey and E. Lecolinet, "A Survey of methods and strategies in character segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 18, no. 7, pp. 690–706, 1996.

[9] Sharma, Preeti, and Manoj Kumar Sachan. "A Review on Character Segmentation of Touching and Half Character in Handwritten Hindi Text." *International Journal of Advanced Research in  Computer Science* 8.3 (2017).

[10] A. Srivastav and N. Sahu, "Segmentation of Devanagari Handwritten Characters," Int. J. Comput. Appl., vol. 142 – No.1, 2016.

[11] Bag and A. Krishna, "Character segmentation of hindi unconstrained handwritten words," Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics), vol. 9448, pp. 247–260, 2015.

[12] M. V. G. Bhujade and M. C. M. Meshram,  "A Technique for Segmentation of  Handwritten Hindi          Text," Int. J. Eng. Res.Technol., vol. 3, no. 2, pp. 1491–1495, 2014

[13] S. Palakollu and R. Rani, "Handwritten     Hindi Text Segmentation Techniques for Lines and Characters," vol. I, 2012.

[14] Kamble, Sandip N., and Megha Kamble. "Morphological Approach for Segmentation of Scanned  Handwritten Devnagari Text 1."

[15] J. Dongre, vikas and H. mankar, vijay, "Devanagri Document Segmentation Using Histogram approach," Int. J. Comput. Sci. Eng. Inf. Technol., vol. 1, No.3, 2011.

[16] N. K. Garg, L. Kaur, and M. K. Jindal, "The segmentation of half characters in handwritten Hindi text," Commun. Comput. Inf. Sci., vol. 139 CCIS, pp. 48–53, 2011.

[17] Garg, Naresh Kumar, Lakhwinder Kaur, and M. K. Jindal. "Segmentation of handwritten hindi text." International Journal of Computer Applications (IJCA) 1.4 (2010): 22-26.

[18] Koshti, Dipak K., and Sharvari Govilkar. "Segmentation of touching characters in handwritten devanagari script." International Journal of Computer Science and its Applications 2.2 (2012): 83-87.

[19] V. Bansal and R. Sinha, "Segmentation of Touching characters in Devanagari," Proc. CVGIP, Delhi, vol. 2, no. 2, pp. 83–87, 1998.

[20] Ma, Huanfeng, and David Doermann. "Adaptive Hindi OCR using generalized Hausdorff image comparison." ACM  Transactions on Asian Language Information Processing (TALIP) 2.3 (2003): 193-218.

[21] Pathak, Rajesh, and Ravi Kumar Tewari. "Distinction between machine printed text and handwritten Text in a document." International Journal of Scientific Engineering and Research (IJSER) 3.7 (2015): 13-17.

[22] Sachan, M.K., Lehal, G.S., Jain, V.K. (2011) 'A Novel Method to Segment Online Gurmukhi Script', Proceedings of International Conference on Information Systems for Indian Languages, ICISIL 2011, Patiala, Communications in Computer and Information Science, Springer-Verlag Berlin Heidelberg, Germany,Vol. 139, pp. 1-8.

[23] Sachan, M.K., Lehal, G.S., Jain, V.K. (2011), 'A System for Online Gurmukhi Script Recognition', Proceedings of International Conference on Information Systems for Indian Languages, ICISIL 2011, Patiala, Communications in Computer and Information Science, Springer-Verlag Berlin Heidelberg, Germany, Vol. 139, pp. 294-295.

[24] Singh, G., and M. Sachan. "A framework of online handwritten gurmukhi script recognition." International Journal of Computer Science and Technology (IJCST) 6 (2015): 52-56.

[25] Sachan, M.K., 'Analysis of Shape Variations in Recognition of Online Gurmukhi Script", Proceedings of International Conference on Electrical and Electronics: Techniques and applications (EETA-2015), 2015,pp. 312-316

_____