

# Off-line Handwriting Text Line Segmentation : A Review

Zaidi Razak<sup>†</sup>, Khansa Zulkiflee<sup>††</sup>, Mohd Yamani Idna Idris<sup>††</sup>, Emran Mohd Tamil<sup>††</sup>, Mohd Noorzaily Mohamed Noor<sup>††</sup>, Rosli Salleh<sup>††</sup>, Mohd Yaakob @ Zulkifli Mohd Yusof<sup>††</sup>, and Mashkuri Yaacob<sup>††</sup>

University of Malaya, Kuala Lumpur, Malaysia

## Summary

Text line segmentation is an essential pre-processing stage for off-line handwriting recognition in many Optical Character Recognition (OCR) systems. It is an important step because inaccurately segmented text lines will cause errors in the recognition stage. Text line segmentation of the handwritten documents is still one of the most complicated problems in developing a reliable OCR. The nature of handwriting makes the process of text line segmentation very challenging. Several techniques to segment handwriting text line have been proposed in the past. This paper seeks to provide a comprehensive review of the methods of off-line handwriting text line segmentation proposed by researchers.

## Key words:

Off-line handwriting recognition, text line segmentation.

## 1. Introduction

In this paper, we review previous work done on text line segmentation in handwritten documents which can be generally categorized into bottom-up and top-down. In the bottom-up approach, the connected components based methods merge neighboring connected components using simple rules on the geometric relationship between neighboring blocks. On the other hand, projection based methods may be one of the most successful top-down algorithms for machine printed documents since the gap between two neighboring text lines in machine printed documents is typically significant, thus the text lines are easily separable. However, these projection based methods cannot be directly used in handwritten documents, unless gaps between lines are significant or handwritten lines are straight.

After a brief description of the characteristics of text line structures in handwritten documents, the rest of the paper is organized as follows. Section 2 describes the challenges in text line segmentation. In Section 3 we review the different approaches to segment a handwritten document into text lines and propose a taxonomy. Section 4 presents an extensive performance evaluation and quantitative comparison of experiment results in the previously

proposed methods. Concluding remarks are given in Section 5.

## 2. Segmentation Challenges

Freestyle and unconstrained handwriting text line segmentation is considered a complex and challenging task due to the following characteristics [1] (Fig. 1):

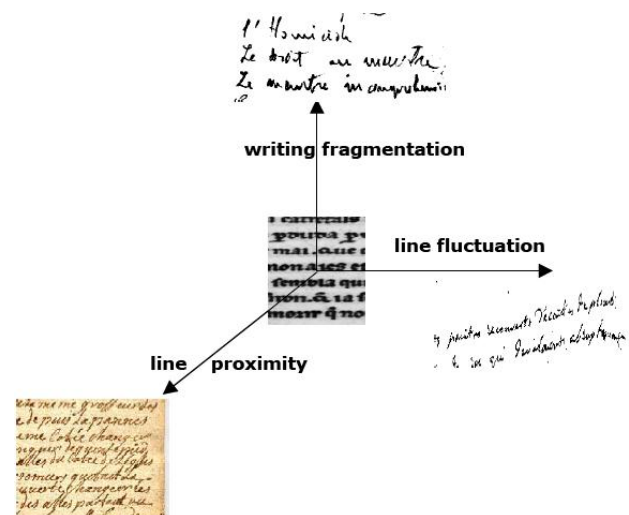


Fig. 1 Complexity in handwritten documents (from Likforman Sulem et.al [1]).

### 2.1 Line Fluctuation

Fluctuating lines or skew variability [2,3]. Lines of text in general are not straight. The inter-line distance variability and inconsistent distance between the components may vary due to writer movement. It may be straight, straight by segments, or curved [1]. According to Okun [4], three types of skew exist in documents:

- A Global skew: all the page blocks have the same orientation,

- multiple skew: unaligned paragraphs or slant is different in different blocks of the page such as the FLAUBERT's drafts [2] which contain several blocks of text arranged in a non linear way, and numerous editorial marks such as erasures and word insertion, and
- non uniform text line skew or varying text line slope: slant is different along the same line of text, for example curvilinear text lines.

## 2.2 Line Proximity

Small gaps between neighboring text lines will cause touching or overlapping of ascenders or descenders [1,3,5,6]. Text lines may be touching or overlapped, when upper-strokes and down-strokes of two consecutive lines are near or touching, or ascenders and descenders of adjacent lines interfere (Fig. 2). Church registers were written with lines of text being close to each other and changing type, size and shape of the handwriting causing text of a given line possibly reaching into adjoining lines and running into each other [7].

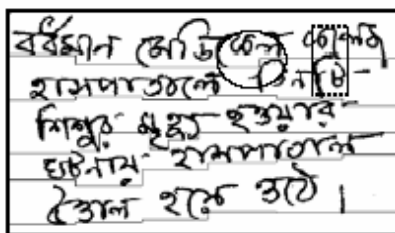


Fig. 2: Overlapping components separated (circle) and touching component separated into two parts (rectangle) in Bangla writing (from Pal and Datta [8]).

## 2.3 Writing Fragmentation

Characters are made up of more than one connected component. This applies to Indian scripts such as Telugu, Tamil, Bangla, and Malayalam and Arabic writing with massive presence of diacritical points.

## 3. Existing Approaches

Handwriting text line segmentation approaches can be categorized according to the different strategies used. These strategies are projection based, smearing, grouping, Hough-based, graph-based and Cut Text Minimization (CTM) approach. Related work can be found in [18-35].

### 3.1 Projection-based approach

In this approach the vertical projection profile is obtained by summing pixel values along the horizontal axis for each y value. From the vertical profile, the vertical gaps between the text lines can be determined. A profile curve can be obtained by projecting black/white transitions or the number of connected components. The profile curve is then analysed to find its maxima and minima. For skewed or moderately fluctuating text lines, the image may be divided into vertical strips and profiles projected in each strip (Zahour *et al.* [3]) using piecewise projections in globally adapting to local fluctuations.

The technique used in [3] divides the text image into columns. In this approach, a partial projection is performed on each column. Then a partial contour following method is used to detect the separating lines, in the direction and opposite direction of the writing. Tripathy and Pal [5] also used a projection based method in text line columns and combined the results of adjacent columns into a longer text line. The document is divided into vertical stripes. Analyzing the heights of the water reservoirs obtained from different components of the document, the width of a stripe is calculated. Stripe-wise horizontal histograms are then computed and the relationship of the histograms peak-valley points is used for line segmentation.

The projection-based algorithm proposed by Arivazhagan *et al.* [7] first obtains an initial set of candidate lines from the piece-wise projection profile of the document (Fig. 3). The lines traverse around any obstructing handwritten connected component by associating it to the line above or below. A decision of associating such a component is made by (i) modeling the lines as bivariate Gaussian densities and evaluating the probability of the component under each Gaussian or (ii) the probability obtained from a distance metric. The proposed method is robust to handle skewed documents and touching lines.

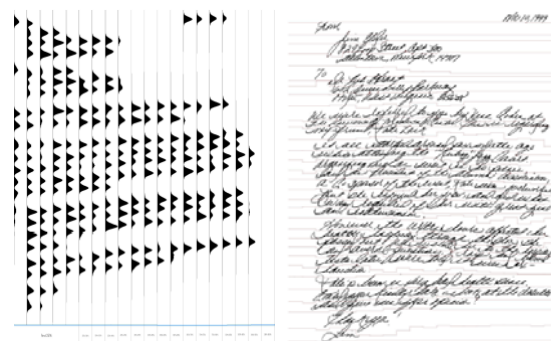


Fig. 3: Piece-wise projection profile (from Arivazhagan *et al.* [7])

Yanikoglu and Sandon [9] first searched for the handwriting text line boundaries and then processed each text line in turn. The boundary between two text lines is not a straight line if the text lines are touching or overlapping. To find the exact boundary between two text lines, they searched for the rough boundary location by analyzing the horizontal pixel density histogram of the line. They then apply a contour following algorithm within that zone to find the exact boundary. The contour following algorithm is modified to operate within a rectangular zone. It is also changed to force a cut at the half-line when necessary, in order to separate text lines that are touching and cannot be separated otherwise.

The algorithm in [10] localized lines by computing the horizontal histograms for the entire image at a couple of relevant skew angles; then the angle and position where the histograms have local minima were chosen as the location between lines. Calculation of the horizontal histograms was done using the traditional histogram calculation executable on DSPs. They refined the line finding algorithm by using a method to blur the words without affecting their location. They computed the pseudo convex hull of each word using the HOLLOW template. The horizontal histogram computed on the pseudo convex hulls is smoothed further via sliding-window. Then the local maximum of the histogram was located since these correspond to the location of the lines. Thresholds were specified to associate all maxima with one line.

### 3.2 Smearing approach

In this technique, consecutive black pixels along the horizontal direction are smeared. If the distance between the white space is within a predefined threshold, it is filled with black pixels. The bounding boxes of the connected components in the smeared image are considered as text lines.

Li et. al [11] proposed a new approach for text line detection by adopting a state-of-the-art image segmentation technique. They first convert a binary image to gray scale using a Gaussian window, which enhances text line structures. Text lines are extracted by evolving an initial estimate using the level set method (Fig. 4). Preliminary experiments show that their method is more robust compared to a bottom-up connected component based approach. Examples show that the method is script independent. This has been qualitatively confirmed by testing it on handwritten documents in different languages, such as Arabic, English, Chinese, Hindi, and Korean. Statistical results show that the algorithm produces

consistent results under reasonable variation of skew angles, character sizes, and noise.

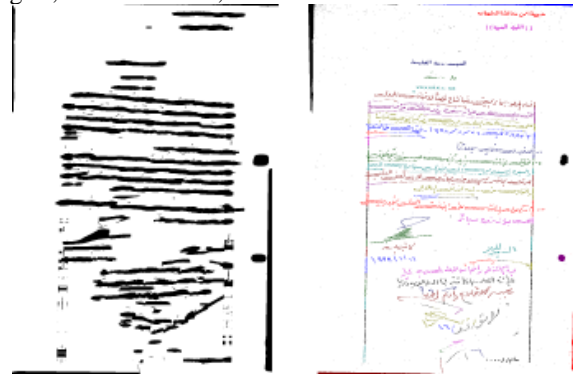


Fig. 4: Smearing using the level set method (from Li et. al [19])

### 3.3 Grouping approach

This method involves building alignments by aggregating units in a bottom-up approach. Units such as pixels, connected components, or blocks are then joined together to form alignments. Likforman-Sulem and Faure [12] proposed an approach based on perceptual grouping of connected components of black pixels. Text lines are iteratively constructed by grouping neighboring connected components based on certain perceptual criteria such as similarity, continuity and proximity. Therefore local constraints on the neighboring components are combined with global quality measures. To handle conflicts, the technique merges a refinement procedure combining a global and a local analysis. According to the authors the proposed technique cannot be used on degraded or poorly structured documents, such as modern authorial manuscripts.

Feldbach and Tönnies [13] proposed a method for line detection and segmentation in historical church registers (Fig. 5). This method is based on local minima detection of connected components and is applied on a chain code representation of the connected components. Line segments are gradually constructed until a unique text line is formed. This algorithm is able to segment text lines closed to each other, touching text lines and fluctuating text lines.

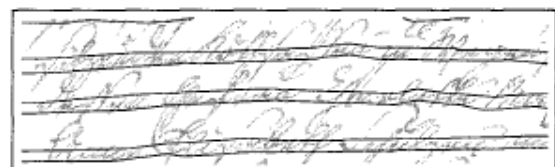


Fig. 5: Line segmentation using a grouping approach (from Feldbach and Tönnies [13])

In [2] the text line extraction problem is seen in the view of artificial intelligence using a production system. The aim is to cluster connected components in the document into homogeneous sets, corresponding to the text lines of the document. To resolve this problem, a search is applied over the graph that is defined by the connected components as vertices and the distances among them as edges.

### 3.4 Hough-based approach

The Hough transform is used for locating straight lines in images. In [14] an iterative hypothesis validation strategy based on Hough transform was proposed. The skew orientation of handwritten text lines is acquired by applying the Hough transform to the center of gravity of each connected component in the document image. If most nearest neighbors of the components in the alignment string belong to the group of components forming the alignment, the alignment has both properties of direction continuity and proximity, and it will be accepted as a text line. This enables the generation of several text line hypotheses. Then a validation is performed to eliminate incorrect alignments between connected components using contextual information such as proximity and direction continuity criteria. Based on the authors, this technique is able to detect text line in handwritten documents which may contain lines oriented in different directions, erasures and annotations between main lines.

Louloudis et al. [15] presented a text line detection method for unconstrained handwritten documents based on a strategy that consists of three distinct steps. The first step comprises preprocessing for image enhancement, connected component extraction and average character height estimation. In the second step, a block-based Hough transform is applied for the detection of potential text lines (Fig. 6) while a third step is used to correct possible false alarms. A grouping method of the remaining connected components which uses the gravity centers of the corresponding blocks is applied. The performance of the proposed method is determined based on a consistent and concrete evaluation technique that relies on the comparison between the text line detection result and the corresponding ground truth annotation.

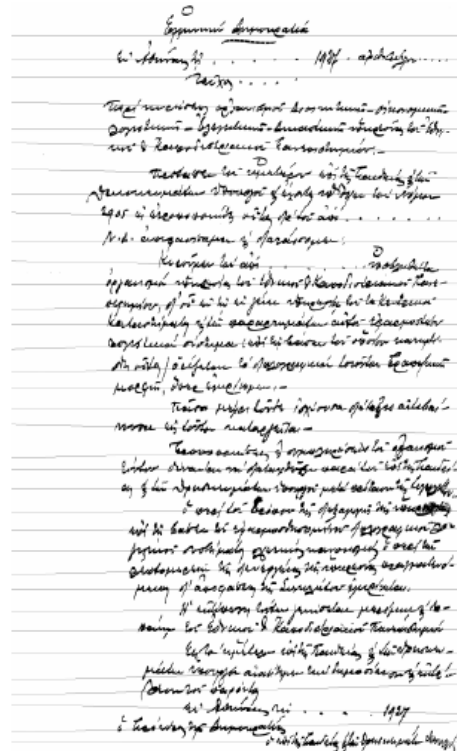


Fig. 6: A block-based Hough transform (from Louloudis et al. [15])

### 3.5 Graph-based approach

A method based on a shortest spanning tree search is presented in [16]. The principle of the method consists of building a graph of main strokes of the document image and searching for the shortest spanning tree of this graph. This method assumes that the distance between the words in a text line is less than the distance between two adjacent text lines.

Sesh Kumar et al. [17] presented a graph cut based framework using a swap algorithm to segment document images containing complex scripts such as in Indian languages. The text block is first segmented into lines using the projection profile approach. The framework enables learning of the spatial distribution of the components of a specific script (Fig. 7) and can adapt to a specific document collection, such as a book. Moreover, they can use both corrections made by the user as well as any segmentation quality metric to improve the segmentation quality.

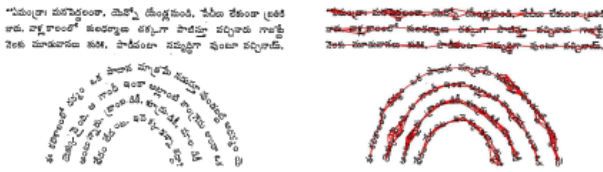


Fig. 7: Graph cut based approach (from Sesh Kumar et al. [17])

### 3.6 CTM Approach

In [6], the CTM method finds a path or cut line in between the text lines to be separated which minimizes the text line pixels cut by the segmentation line, especially descenders from the upper line and ascenders from the lower line. The method attempts to track around ascenders or descenders to avoid cutting them. If the deviation is too great, the segmenter aborts and continues its forward path (Fig. 8). A rough estimate of text line separations were first obtained using vertical projection histograms.

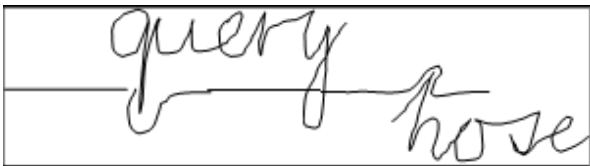


Fig 8: CTM method (from Weliwitage et. al [6])

## 4.0 Performance Evaluation and Comparison

The proposed algorithm described in [3] was tested on 100 samples of text containing 1000 lines. The experiment results in 97% accuracy in correct line segmentation. Errors were caused by baseline-skew variability, and overlaps between characters and diacritical marks in the first column. The association of the diacritical symbols to a text line caused significant errors when the symbol is distant from the separating border line.

Tripathy and Pal [5] used 1627 text lines in single column document pages with different writing styles. Text lines segmentation accuracy is calculated by drawing boundary lines between two consecutive text lines. Then, from the computer's display, line segmentation accuracy was calculated manually. If all text lines are extracted correctly, segmentation accuracy is considered 100%. From 1627 lines tested, 984 lines were segmented correctly.

Experimental results in [7] show that on 720 documents including English, Arabic and children's handwriting consisting of 11,581 lines, 97.31% of the lines were

segmented correctly. On an experiment of over 200 handwritten images with 78,902 connected components, 98.81% of them were associated to the correct lines. Experiments were also conducted on 300 exam essays written on ruled line paper to test the robustness of the algorithm. The algorithm segmented with an accuracy of 96.3%. Furthermore, the segmentation algorithm proved language independence, with accuracy of 98.62% on 120 Arabic handwritten images. Most of the dots above or below a word were associated to the correct lines by the proposed algorithm.

Documents scanned from a Telugu book titled "Aadarsam" containing 256 pages printed in 1973 were used in [17]. The performance of segmentation is calculated using a segmentation quality metric. Segmentation corrections required on a new page are close to zero after adapting to the first 44 pages. The remaining 212 pages were segmented correctly.

The method in [6] was tested on 30 images from NIST special database consisting data in 34 text boxes from 2100 forms scanned at a resolution of 300 pixels / inch and saved in binary format. Text boxes in the form corresponding to a text paragraph of 52 words were abstracted for text line segmentation. The segmentation resulted in 183 text segmenting lines and 213 text lines in the images were correctly segmented into 176 lines with 96% accuracy.

Likforman-Sulem and Faure [12] tested their method on handwritten documents. Fluctuating text lines, sloped annotations or annotations added between main lines were detected. Fluctuation combined with proximity of text lines may cause merged lines. The algorithm in [17] was tested on 10 pages of a handwriting database containing 7000 words from the LOB corpus (Lancaster-Oslo / Bergen) written by a single writer. The experiments were conducted using the MatCNN simulator in Matlab software. The line segmentation algorithm correctly segmented each line in every page.

The algorithm in [11] was tested on more than 10,000 diverse handwritten documents in different scripts, such as Arabic, Hindi, and Chinese. During testing, if a ground-truth line and the corresponding detected line share at least 90% pixels, a text line is considered to be detected correctly. From a total of 2,691 ground-truth lines, their approach correctly detected 2,303 (85.6%) lines. At the text line level, the connected component based method performs significantly worse where only 951 (35.3%) text lines were detected correctly. Errors were caused by two adjacent text lines overlapping significantly, signatures,

the correction in the gap between two lines, and the severe noise introduced during scanning.

The proposed text line detection method in [15] which uses a Block-Based Hough Transform approach was tested on unconstrained handwritten Greek documents using 20 document images taken from the historical archives of the University of Athens for which corresponding text line detection ground truth was manually created. A number of 450 text lines were detected in the images with 96.87% accuracy. Difficulties concerning the extraction of text lines include the variety of accents appearing above or under the text line body and the small difference in the skew angle in the text lines.

In [2], overlapped text lines and text lines where the interline distance is smaller than the intra line distance will cause segmentation errors. Only well separated text lines were correctly segmented. The algorithm in [13] which used a chain code representation was tested on text in church registers ranging 300 years. The algorithm was

applied to images from 61 paragraphs in 7 pages. All 61 paragraphs consisted of 300 lines. The algorithm produced good results with constant values on different handwriting styles. 222 lines (90%) of 246 lines in 49 paragraphs containing six different handwritings were reconstructed correctly. The experimental results are summarized in Table 1.

## 5.0 Conclusion

This paper has provided a comprehensive review of the methods for off-line handwriting text line segmentation previously proposed by researchers. After a brief description of the characteristics of text line structures in handwritten documents, we have describes the challenges in text line segmentation. We also reviewed the different approaches to segment a handwritten document into text lines and proposed a taxonomy. An extensive performance evaluation and quantitative comparison of experiment results in the previously proposed methods was performed.

Table 1: Experimental results

Author	Experiment data	Experiments	Accuracy	Segmentation errors
Zahour et. al [3]	1000 lines in 100 samples	Implemented in a C++ language on a 200MHz PC	97%	Baseline-skew variability, and overlaps between characters and diacritical marks in the first column. Association of diacritic symbols distant from the separating border line, to a text line.
Tripathy and Pal [5]	1627 lines in single column pages with different writing styles	Draw boundary lines between two consecutive text lines. Then, from the computer's display, line segmentation accuracy was calculated manually.	60%	When two consecutive words touch, or distance between two consecutive words is very small.
Arivazhagan et. al [7]	11,581 lines in 720 documents including English, Arabic and children's handwriting	The cut-through accuracy corresponds to the proportion of components correctly classified as overlapping components or cut-through error.	97.31%	A normal component, spanning across two or more lines or lying in between two lines of text
Sesh Kumar et al. [17]	Documents scanned from a Telugu book titled "Aadarsam", printed in 1973 containing 256 pages.	The performance of segmentation is calculated using a segmentation quality metric	Segmentation corrections required on a new page are close to zero after adapting to the first 44 pages. The remaining 212 pages were segmented correctly.	Not stated
Weliwitage et. al [6]	30 images from NIST special database in 34 text boxes from 2100 forms	Text boxes in the form corresponding to a text paragraph of 52 words were abstracted for text line segmentation.	96%	Very short text lines of not longer than one word and text lines not starting from left margin of the image and unclear



				separation of text lines with merging ascenders and descenders may be detected as an extra text line.
Likforman-Sulem et. al [12]	Unconstrained handwritten rough drafts, address blocks, letters and manuscripts.	Alignments found as text lines are crossed by a line, components belonging to the same line share the same identification number inscribed above their enclosing rectangles. Alignments which were found in the Hough domain, but invalidated in a second stage are crossed by a dashed line. Ambiguous components are inscribed in dashed rectangles.	Not stated	When fluctuation is combined with proximity of text lines, merging lines may appear.
Tímár et. al [10]	10 pages of a handwriting database containing 7000 words from the LOB corpus (Lancaster-Oslo / Bergen) written by a single writer	Conducted using the MatCNN simulator in Matlab software	Correctly segmented each line in every page	Not stated
Li et. al [11]	More than 10,000 diverse handwritten documents Arabic, Hindi, and Chinese script	If a ground-truth line and the corresponding detected line share at least 90% pixels, a text line is considered to be detected correctly.	From 2,691 ground-truth lines, correctly detected 2,303 (85.6%) lines. At the text line level, only 951 (35.3%) text lines were detected correctly.	Errors were caused by two adjacent text lines overlapping significantly, signatures, the correction in the gap between two lines, and the severe noise introduced during scanning.
Louloudis et. al [15]	Unconstrained handwritten Greek documents using 20 document images taken from the historical archives of the University of Athens	Block-Based Hough Transform approach for which corresponding text line detection ground truth was manually created	450 text lines were detected in the images with 96.87% accuracy	Various accents above or under the text line body and the small difference in the skew angle in the text lines
Nicolas et. al [2]	Collection of handwritten French novelist Gustave Flaubert drafts		Only well separated text lines were correctly segmented.	Overlapped text lines and text lines where the interline distance is smaller than the intra line distance will cause errors.
Feldbach and Tönnies [13]	Text in church registers from 61 paragraphs consisting of 300 lines in 7 pages ranging 300 years.	Chain code representation	222 lines (90%) of 246 lines in 49 paragraphs containing six different handwritings were reconstructed correctly.	Serious errors occur when the difference between the reconstructed base and centre lines and the real lines was higher than the script size, if a text line was not found, or if an extra line was found at a wrong place.

## 6.0 References

- [1] L. Likforman Sulem, A. Zahour, B. Taconet, "Text line segmentation of historical documents: a survey", *IJDAR*, Vol. 9, No. 2-4, 2007, pp. 123-138.
- [2] S. Nicolas, T. Paquet, L. Heutte, "Text Line Segmentation in Handwritten Document Using a Production System", *Proceedings of the 9th IWFHR*, Tokyo, Japan, 2004, pp. 245-250.
- [3] A. Zahour, B. Taconet, P. Mercy, and S. Ramdane, "Arabic Hand-written Text-line Extraction", in *Proceedings of the Sixth International Conference on Document Analysis and Recognition, ICDAR 2001*, Seattle, USA, September 10-13 2001, pp. 281-285.
- [4] O. Okun, M. Pietikainen, and J. Sauvola, "Document skew estimation without angle range restriction," *IJDAR* 2, pp. 132 - 144, 1999.
- [5] N. Tripathy and U. Pal. "Handwriting Segmentation of Unconstrained Oriya Text," in *International Workshop on Frontiers in Handwriting Recognition*, 2004, pp. 306-311.
- [6] C. Wellitige, A. L. Harvey, A. B. Jennings, "Handwritten Document Offline Text Line Segmentation", in *Proceedings of Digital Imaging Computing: Techniques and Applications*, 2005, pp. 184-187.
- [7] M. Arivazhagan, H. Srinivasan, S. N. Srihari, "A Statistical Approach to Handwritten Line Segmentation", in *Proceedings of SPIE Document Recognition and Retrieval XIV*, San Jose, CA, February 2007.
- [8] Pal U., Datta S. (2003), Segmentation of Bangla unconstrained handwritten text, *Proceedings of Seventh International Conference on Document Analysis and Recognition*, pp 1128 – 1132.
- [9] B. Yanikoglu and P. A. Sandon, "Segmentation of Off-line Cursive Handwriting using Linear Programming", *Pattern Recognition*, Vol. 31, No. 12, 1998, pp. 1825-1833.
- [10] G. Tímár, K. Karacs, Cs. Rekeczky, "Analogic Preprocessing and Segmentation Algorithms For Offline Handwriting Recognition", *Proceedings of IEEE CNNA'02*, World Scientific 2002, pp.407-414.
- [11] Y. Li, Y. Zheng, D. Doermann, and S. Jaeger, "A new algorithm for detecting text line in handwritten documents," in *International Workshop on Frontiers in Handwriting Recognition*, 2006, pp. 35-40.
- [12] L. Likforman-Sulem, C. Faure, "Extracting text lines in handwritten documents by perceptual grouping", *Advances in handwriting and drawing : a multidisciplinary approach*, C. Faure, P. Keuss, G. Lorette and A. Winter Eds, Europia, Paris, 1994, pp. 117-135.
- [13] M. Feldbach, K. D. Tönnies, "Line Detection and Segmentation in Historical Church Registers", *Sixth International Conference on Document Analysis and Recognition*, September, 2001. pp. 743-747.
- [14] L. Likforman-Sulem, A. Hanimyan, C. Faure, "A Hough based algorithm for extracting text lines in handwritten documents", *Third International Conference on Document Analysis and Recognition*, Vol. 2, August 1995, pp. 774-777.
- [15] G. Louloudis, B. Gatos, I. Pratikakis, K. Halatsis, "A Block-Based Hough Transform Mapping for Text Line Detection in Handwritten Documents", *Proceedings of the Tenth International Workshop on Frontiers in Handwriting Recognition*, La Baule, Oct. 2006.
- [16] I.S.I. Abuhaiba, S. Datta, M.J.J. Holt, "Line Extraction and Stroke Ordering of Text Pages", *Proceedings of the Third International Conference on Document Analysis and Recognition*, Montreal, Canada, 1995, pp. 390-393.
- [17] K.S. Sesh Kumar, A. M. Namboodiri, C.V. Jawahar, "Learning Segmentation of Documents with Complex Scripts", in *Fifth Indian Conference on Computer Vision, Graphics and Image Processing*, Madurai, India, LNCS 4338, 2006, pp.749-760.
- [18] Calabretto S., Bozzi A.(1998) The Philological Workstation BAMBI (Better Access to Manuscripts and Browsing of Images). *International Journal of Digital Information (JoDI)*. 1 (3):1-17. ISSN: 1368-7506.
- [19] He J., Downton, A.C. (2003). User-Assisted Archive Document Image Analysis for Digital Library Construction, *Seventh International Conference on Document Analysis and Recognition*, Edinburgh.
- [20] Oztop E., Mulayim A. Y., Atalay V., Yarman-Vural F. (1999), Repulsive attractive network for baseline extraction on document images, *Signal Processing*, 75:1-10.
- [21] Pu Y., Shi Z. (1998) A natural learning algorithm based on Hough transform for text lines extraction in handwritten documents. In *Proceedings of the 6 Intl. Workshop on Frontiers in Handwriting Recognition*, Taejon, Korea, pp. 637- 646
- [22] Shapiro V., Gluhchev G., Sgurev V. (1993), Handwritten document image segmentation and analysis, *Pattern recognition Letters*, 14:71-78.
- [23] Shi Z., Govindaraju V. (2004) Line Separation for Complex Document Images Using Fuzzy Runlength, *Proc. of the Int. Workshop on Document Image Analysis for Libraries*, Palo, Alto, CA, January 23-24.
- [24] Shi Z., V. Govindaraju (2004), Historical Document Image Enhancement using background light intensity normalization, *ICPR 2004*, Cambridge.
- [25] Tseng Y.H., Lee H.J. (1999), Recognition-based handwritten Chinese character segmentation using a probabilistic Viterbi algorithm, *Pattern Recognition Letters*, 20( 8):791-806.
- [26] A. El-Yacoubi, M. Gilloux, R. Sabourin, and C. Y. Suen. An HMM-Based Approach for Off-Line Unconstrained Handwritten Word Modeling and Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21 (8):752-760, Aug. 1999.
- [27] G. Kim, V. Govindaraju, and S. N. Srihari. An Architecture for Handwritten Text Recognition Systems. *International Journal on Document Analysis and Recognition*, 2(1):374, Feb. 1999.
- [28] S. Madhvanath, V. Govindaraju: Local reference lines for handwritten phrase recognition. *Pattern Recognition* 32(12): 2021-2028 (1999)
- [29] S. Madhvanath, E. Kleinberg, and V. Govindaraju. Holistic Verification of Handwritten Phrases. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2 1 (1 2): 1344-1356, Dec. 1999.
- [30] Y. Pu and Z. Shi. A Natural Learning Algorithm Based on Hough Transform for Text Lines Extraction in Handwritten Documents. In *Proceedings of the Sixth International Workshop on Frontiers of Handwriting Recognition (IWFHR VI)*, Taejon, Korea, pages 637-646, 1998.
- [31] M. Shridar and F. Kimura. Segmentation-Based Cursive Handwriting Recognition. In H. Bunke and P. S. P. Wang, editors, *Handbook of Character Recognition and Document Image Analysis*, pages 123-156. World Scientific, Feb. 1997.



- [32] P. Steiner. Zwei ausgewählte Probleme zur Offline-Erkennung von Handschrift. Diploma thesis, University of Bern, Aug. 1995.
- [33] L. O’Gorman, “The document spectrum for page layout analysis,” *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 15, no. 11, pp. 1162–1173, 1993.
- [34] S. Jaeger, G. Zhu, D. Doermann, K. Chen, and S. Sampat, “DOCLIB: A software library for document processing,” in *Document Recognition and Retrieval XIII, Proc. of SPIE* vol. 6067, 2006, pp. 63–71.
- [35] G. Nagy, S. Seth, and M. Viswanathan, “A prototype document image analysis system for technical journals,” *Computer*, vol. 25, no. 7, pp. 10–22, 1992.



**Zaidi Razak** obtained his Bachelor’s degree in Computer Science and Master in Chip Design from University of Malaya in 2000. Currently, he is a lecturer at the Faculty of Computer Science and Information Technology, University of Malaya. His research areas include image processing, Jawi character recognition and System on Chip (SoC)

design. He has published a number of papers related to these areas.



**Khansa Zulkiflee** obtained her Bachelor’s degree in Computer Science from University Technology Malaysia. Currently she is a research assistant at the Faculty of Computer Science and Information Technology, University of Malaya. Her research areas include image processing, Jawi character recognition and System on Chip (SoC) design.



**Mohd Yamani Idna Idris** obtained his Bachelor’s degree in Electrical Engineering and Master’s degree in Computer Science from University of Malaya. Currently, he is a lecturer at the Faculty of Computer Science and Information Technology, University of Malaya. His research areas include System on Chip, SCADA, image

processing.



**Emran Mohd Tamil** obtained his Bachelor’s degree in Electrical-Robotic Engineering from University Technology Malaysia in and Master’s degree in Information Technology from University Technology MARA. Currently, he is a lecturer at the Faculty of Computer Science and Information Technology, University of Malaya. His specializations

are system and network and current research interests are embedded systems, network security, SCADA and chip design.



**Mohd Noorzaily Mohamed Noor** obtained his Bachelor’s and Master degree in Computer Science from University of Malaya. Currently, he is a lecturer at the Faculty of Computer Science and Information Technology, University of Malaya. His research areas include arithmetic and logic structures and detection and estimation. He has published a number of papers related to these areas.



**Rosli Salleh** obtained his Bachelor’s degree in Computer Science from University of Malaya. Then he obtained his Master’s degree in Data Communication Networking and PhD in Computer Science majoring in Virtual Reality, Tele-surgery and Networking both from University of Salford Manchester United Kingdom. He has also obtained CCNA professional qualifications.

Currently, he is a lecturer at the Faculty of Computer Science and Information Technology, University of Malaya. His specializations include Bluetooth network, networking, and Virtual Reality. His current research interests are Bluetooth Scatternet Formation, Public Key Infrastructure, Network Security, Authentication Server, Virtual Simulation, Virtual Reality and Laparoscopic Surgical Training.



**Mohd Yaakob @ Zulkifli Mohd Yusof** obtained his Bachelor’s degree, Master’s degree and PhD from University of Malaya, University of Jordan and University of Wales, respectively. Currently, he is a lecturer at the Academy of Islamic Studies, University of Malaya. His research areas include Quranic Exegesis, Quranic Studies and

Methodology Of Quranic Exegesis.



**Mashkuri Yaacob** is currently the third vice-chancellor of Universiti Tenaga Nasional (Uniten). Mashkuri, who obtained his Bachelor’s degree in Electrical Engineering from the University of New South Wales, Sydney, also holds master’s and doctoral degrees in electronics and computer engineering from the University of Manchester in the

United Kingdom, was named deputy dean of the Engineering Faculty in 1983 and later dean of the Computer Science and Information Technology Faculty in 1994. Mashkuri’s career at UM culminated with his appointment as deputy vice-chancellor (academic) from June 2000 to May 2003. Mashkuri’s main research interests are discrete wavelet, operating system and integrated information.