

A Search Engine for Handwritten Documents

Sargur Srihari, Chen Huang, Harish Srinivasan

Center of Excellence for Document Analysis and Recognition(CEDAR)
University at Buffalo, State University of New York, Buffalo, USA
{srihari,chuang5,hs32}@cedar.buffalo.edu

ABSTRACT

The design and functionality of a versatile search engine on handwritten documents is described. Documents are indexed using global image features, e.g., stroke width, slant, word gaps, as well local features that describe shapes of characters and words. Image indexing is done automatically using page analysis, page segmentation, line separation, word segmentation and recognition of characters and words. Several types of searches are permitted: (i) Word / Phrase Spotting; (ii) Text to Image Search; (iii) Plain Text Search; (iv) Word Recognition from Lexicon. The words in the document are characterized by various features and it forms the basis for the different searching techniques. The system was implemented in Microsoft Visual C++. The paper reports on the functional capabilities of the various search techniques and their performance.

1. INTRODUCTION

The need for searching archives of scanned handwritten documents are in applications such as historical manuscripts, forensics, personal records, as well as criminal records. In each of these applications there is a need for indexing and retrieval based on textual content as well as user-indexed terms. In the forensic application, there is a need for searching a large database of handwritten documents not only for textual information but also for visual content such as writer characteristics. There exists many text books¹⁻⁵ describing the methodology employed by forensic document examiners. This paper describes the various possible searches that one can attempt and to provide for a complete search engine for a digital library of handwritten documents. Search engines in general are the most popular applications that are a part of the word wide web that provides for a complete textual information retrieval system. On the other hand scanned handwritten documents provide images to search on rather than text. Combining the two, results in having the query as well as the results to be either text or image. A handwriting document management system, known as CEDAR-FOX^{6,7} has been developed at CEDAR, whose search aspects are the subject of this paper.

1.1. CEDAR-FOX System

As a document management system for handwritten documents, CEDAR-FOX^{6,7} provides several functionalities: interactive document analysis, searching a digital library of handwritten documents, image indexing for content-based retrieval. For the purpose of indexing based on writer characteristics, features are automatically extracted after several image processing functions followed by character recognition. The features of the documents are extracted at various levels, document level, word level and character level. The indexing includes segmenting a given document automatically into lines and then into words. The user can then use interactive graphical tools to assist in obtaining more accurate features, for e.g., providing the transcript of the document. The unique aspect of CEDAR-FOX is that image matching is driven by probabilistic writer verification, i.e., whether the query and the document were likely to be written by the same writer.

The search capabilities of the system come with interactive graphical user interfaces that provide for searching (i) within a single document; (ii) from another active document and (iii) from a collection of preprocessed indexed documents. The searching can also be done on a region of interest (ROI) of a document.

1.2. Organization of paper

Section 2 describes the indexing aspects of CEDAR-FOX which has two parts: (i) Preprocessing and (ii) Feature Extraction. Section 3 is divided into two sections (i) Similarity Measures and (ii) Various Search Techniques. Section 4 describes the performance and evaluation and Section 5 is the conclusion remarks.

2. INDEXING

Indexing of handwritten document images does not only involve textual information like in the IR counterpart but includes visual image features of words and characters.

2.1. Line and Word Segmentation

Indexing of documents begins at the image processing stage. The raw image of the document is subjected to a number of image processing functions that automatically segment the document into lines and then into words. Character components from these words are then recognized. After the preprocessing step, each document is now characterized by their line, word and character images. Information such as number of lines, words and characters, their position in the document are stored along with the document image. Figure 1 shows the line and word segmentation.

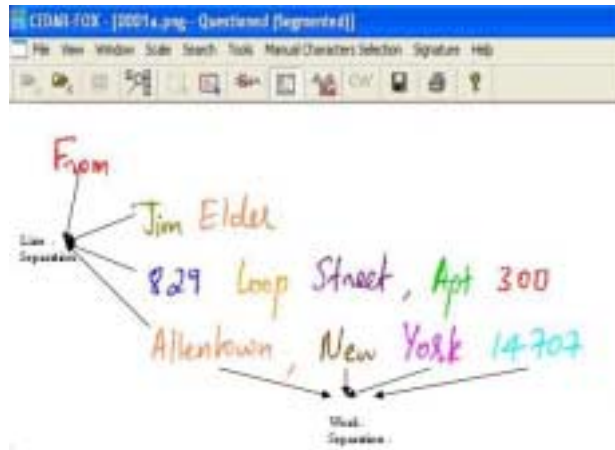


Figure 1. Example of line and word segmentation.

2.2. Feature Extraction

Image features are computed at the document word and the character levels. The document level features are called macro features whereas the word and character level features are called the micro features. The initial 12 macro features reported in Ref. 6,7 were truly at the document level, meaning they were computed either directly from the entire document (no. of black pixels, threshold, etc.) or on a line-by-line basis and applied to the entire document (no. of exterior contours, no. of interior contours, slope, etc.). Macro features are used for matching documents based on similarity of writership which is not the subject of this paper.

Word features extracted from each word of the document form the key index for image search. The Word features used in CEDAR-FOX consist of 1024 bits corresponding to gradient (384 bits), structural (384 bits) and concavity (256 bits) features. Each of these three sets of features relies on dividing the scanned image of the word into 8×4 region. The gradient features capture the stroke flow orientation and its variations using the frequency of gradient directions, as obtained by convolving the image with a Sobel edge operator, in each of 12 directions and then thresholding the resultant values to yield a 384-bit vector. The structural features representing the coarser shape of the word capture the presence of corners, diagonal lines, and vertical and horizontal lines in the gradient image, as determined by 12 rules. The concavity features capture the major topological and geometrical features including direction of bays, presence of holes, and large vertical and horizontal strokes. All the 1024 binary features are converted from the original floating point number computations.

3. SEARCH TECHNIQUES

Several different types of searches can be done on a handwritten document based on fact that the query and the results can be either text or image. This primarily requires two kinds of matching: (i) Word image to Word image and (ii) Text to Word image.

3.1. Similarity Measures

Two types of similarity measures are central to matching a query with image data: image-to-image and text-image. These are described below.

3.1.1. Image-Image similarity

Each word image is first converted into a feature vector. This is done by enclosing the image in a 4 x 8 grid and then computing gradient, structural and concavity features for each cell of the grid. The result is a binary vector of length 1024. The success of using this feature vector in word spotting was reported in Ref..⁸

A distance measure between the vectors is used as the similarity between the word images. The *correlation measure*⁹ is used to compare two binary vectors. Given S_{ij} , $i, j \in 0, 1$, the number of occurrences of matches with i in the first feature vector and j in the second feature vector at the corresponding positions, the dissimilarity D between two feature vectors X and Y is given by equation 1.

$$D(X, Y) = \frac{1 - S_{11}S_{00} + S_{10}S_{01}}{2\sqrt{(S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10})}} \quad (1)$$

Two similar word images have a smaller value of D than two dissimilar images. A smaller value of D would result in that image having a higher rank during the searching process.

3.1.2. Image-Text score

For matching a text string to a word image it is necessary to have a measure of confidence that a given word image is composed of the characters in the given text. In word recognition it is a distance measure between the word image and the members of a lexicon. The smaller the value of the score, the more likely the given word image is composed of the characters in the word. To compute this score, a model of word recognition called *word model recognition (WMR)* is used. WMR has been used successfully in various applications such as postal street name recognition.

The key steps of WMR are (i) segmentation: this returns the segmentation points to be used for grouping one or more segment(s) to form meaningful characters; the segmentation points are determined using a combination of ligatures and concavity features on the contour, (ii) feature extraction: generates feature vectors for combination of segments that are hypothesized to be characters, (iii) recognition: uses the segmentation statistics such as how often a particular character is split into how many segments (1.4), character cluster centroids, a lexicon and feature vector derived from the test image. A dynamic matching scheme is employed that computes the distance between the image and the word.¹⁰⁻¹² Figure 2 shows an example of a word image and the possible segmentation points. The matching score computed with one lexicon entry is shown in Figure 2.

3.2. Search Techniques

Each of the four types of search are described here. These are: the four combinations, together with commonly used nomenclature, are

1. Image query to image result: *word/phrase spotting*
2. Image query to text result: *word recognition*
3. Text query to image result: document image search using text word
4. Text query to text result: text search using transcript

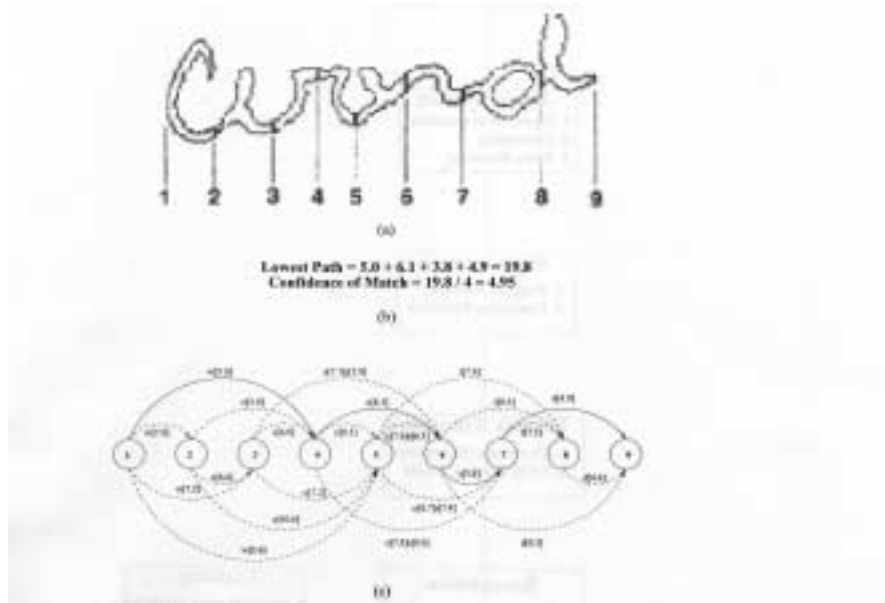


Figure 2. Matching between a sample image and one lexicon entry “word”. (a) Segmentation points for the image, (b) Score of the match, (c) Matching paths and Scores

3.2.1. Word / Phrase Spotting

Here, both the query and result are word/phrase images. The user is given a graphical interface to select a query word / phrase image from a document. A phrase can be two or more words in one line. The query image can now be searched for (i) from within the same document, (ii) in another active document that is currently open or (iii) from a set of preprocessed documents. The Word Feature Similarity discussed in Section 3.1.1 is used to match the query image against all target images. For the first two kinds of searches, the result set consists of all the words/phrase images in the current or the other document ranked from top to bottom. The top rank image is the closest in similarity to the query image and last rank image is the most dissimilar image to that of the query. For the third search wherein, the query image is searched for from a set of preprocessed documents, the top ten images from each of the preprocessed documents are selected and then resultant set is sorted top to bottom by rank. Figure 3 shows a sample screen shot of the CEDAR-FOX system after word spotting.



Figure 3. Screen shot of CEDAR-FOX system with word spotting the image of a word.

3.2.2. Text to Image Search

In this search mode, the query to the search is a word text. The task of this search mode is to return all the word images that are likely to be composed of the characters in the given text. The images that are retrieved can be from (i) current document, (ii) from another open document or (iii) a set of preprocessed documents. The Text vs. Word image similarity measure (WMR) discussed in Section 3.1.2 is used to match all the images against the one query text. The result of the search is a list of word images sorted top to bottom in rank. Each of these retrieved word images form a link to the document from which they came from and hence the corresponding documents can also be retrieved. Figure 4 shows a sample screen shot of the CEDAR-FOX system after Text to Image Search.

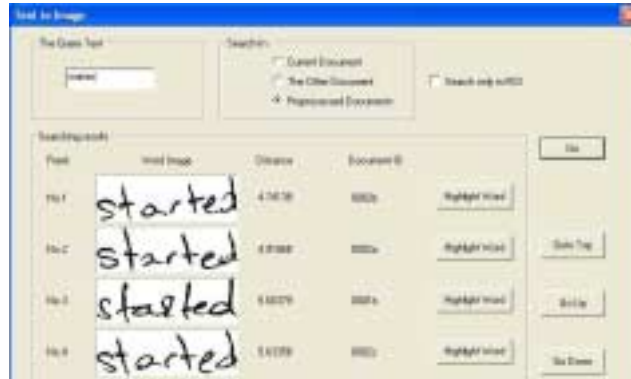


Figure 4. Screen shot of CEDAR-FOX system with the results of Text to Image Search

3.2.3. Word Recognition(Image to Text Search)

In this search mode, the query is a user selected word image from a document. The task of this search is to recognize the word. The user gives a Lexicon of words to search from and the result is a sorted order of the Lexicon. For a given input image, the lexicon entries are sorted based on the distance from the lexicon to the image. The distance measure used here is text vs. Image similarity (WMR). The top ranked lexicon or text represents the most likely content for the query word image. The Text vs. word image similarity measure (WMR) discussed in Section 3.1.2 is used to compute the score between each word of the lexicon to the given query word image. Figure 5 shows a sample screen shot of the CEDAR-FOX system after Word Recognition.



Figure 5. Screen shot of CEDAR-FOX system with the results of Word Recognition

3.2.4. Plain Text Retrieval

The last search option works very similar to the IR counterpart. The query to the search is a text. This search requires that all the documents to be content identified i.e., the entire transcript of the document is already provided with each of the document. The task of this search is to retrieve all the word images that contains the text as that of the query text and rank them. The matching considers the priority of the information represented by the words. The distance measure is edit distance based. Figure 6 shows a sample screen shot of the CEDAR-FOX system after Plain Text Search.



Figure 6. Screen shot of CEDAR-FOX system with the results of Plain Text Search

4. PERFORMANCE

More than 1000 people wrote three samples of a preformatted letter, each. A complete system for handwritten document management, known as CEDAR-FOX, has been developed for document analysis, identification, verification and document retrieval. As discussed, a database for handwritten document can be built using the systems database functionality. The system design stresses the importance of the individuality in handwritten documents since it not only handles meta-data but also stores the discriminating elements of handwriting in the database to be used as indexes during retrieval. The system extracts the individuality represented by the document features. The document identification can then be done using the features.

The system is implemented in a Windows environment. It has a number of interactive features that makes it useful as a document management tool: The system automatic does line and word segmentation, and automatic feature extractors at the document, word and character level. For creation of a document database, the system provides tools to save the entire preprocessed and indexed images. Document retrieval uses several query models. Retrieval can be done using any word in a document or text that user type in. Document retrieval provides user an efficient way to spot the occurrence of a word image, recognize a word from a lexicon list, and search for text using image as well as transcript

4.1. Text to Image Search: Evaluation

To evaluate this search mode, a corpus of 100 document images from different writers consisting of $100 \times 150 = 15,000$ word images was taken. The top 1,000 results were displayed and precision recall curves were averaged using 10 different query text. The query texts used were such that there is exactly one occurrence of that word in every document. Hence in the corpus of 100 documents, there are a total of 100 correct matches i.e., one correct match in each document. The total corpus size is 15,000 words.

Therefore we have $Recall = \frac{x}{100}$ and $Precision = \frac{x}{N}$ where N is the number of documents considered for Precision and x is the number of correctly retrieved word images. Figure 7 shows the Precision Recall Curve.

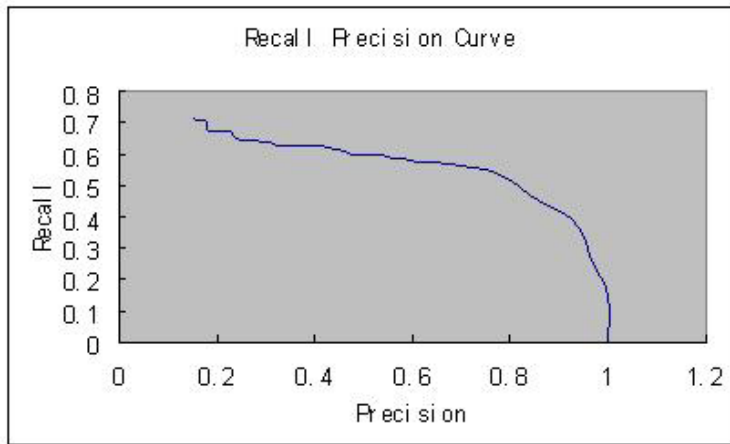


Figure 7. Recall Precision Curve for text to image search.

4.2. Word Recognition: Evaluation

To evaluate this search technique, a lexicon of 150 words were taken. A total of 100 query word images used. Each query word image has exactly one correct match in the Lexicon. The results are tabulated as a percentage of the number of times the correct match was returned as (i) Top rank, (ii) Within top 5 ranks, (iii) Within top 10 ranks. Table 1 shows the evaluation result.

Table 1. Performance of Word Recognition

Rank of The Correct Lexicon for that query word image	Percentage ($\frac{\text{Number of times}}{\text{total no of queries}} \times 100$)
1	89 %
< 5	99 %
< 10	100 %

4.3. Word Spotting: Evaluation

To evaluate this search technique, 100 query word images were used. Each query word image was searched in one another document written by the same writer. The results are tabulated as a percentage of the number of times the correct match was returned as (i) Top rank, (ii) Within top 5 ranks, (iii) Within top 10, (iv) Within top 20, (v) Within top 50. Table 2 shows the evaluation result.

5. CONCLUSION

We have described different kinds of search capabilities of a document analysis and management system for handwritten documents. Image as well as text can be used for retrieval of data. The methods are based on image-to-image similarity measure and a text-to-image score. The performance of the system in the various types of searches has been presented. Future work on searching handwritten documents will involve the use of named entities to overcome ambiguity inherent in handwritten documents.

ACKNOWLEDGMENTS

CEDAR-FOX is the result of effort of many students and researchers at CEDAR. Others who have contributed to the system and its testing are Meenakshi Kalera, Siyuan Chen, Vinay Shah and Vivek Shah.

Table 2. Performance of Word Spotting

Rank of The Correct Image for that query word image	Percentage ($\frac{\text{Number of times}}{\text{total no of queries}} \times 100$)
1	80 %
< 5	81 %
< 10	85 %
< 20	89 %
< 50	100 %

REFERENCES

1. A. S. Osborn, *Questioned Documents, 2nd*, Boyd Print. Co., 1929.
2. E. W. Robertson, *Fundamentals of Document Examination*, Nelson-Hall, 1991.
3. R. R. Bradford and R. Bradford, *Introduction to Handwriting Examination and Identification*, Nelson-Hall, 1992.
4. R. Huber and A. Headrick, *Handwriting Identification: Facts and Fundamentals*, CRC Press, 1999.
5. O. Hilton, *Scientific examination of questioned documents*, CRC Press, 1993.
6. S. N. Srihari, S.-H. Cha, H. Arora, and S. Lee, "Individuality of handwriting," in *Journal of Forensic Sciences*, **47**(4), pp. 856–872, 2002.
7. S. N. Srihari, B. Zhang, C. Tomai, S. Lee, Z. Shi, and Y. C. Shin, "A system for handwriting matching and recognition," in *Proceedings of the Symposium on Document Image Understanding Technology (SDIUT 03)*, Greenbelt, MD, 2003.
8. S. N. S. B.Zhang and C.Huang, "Word image retrieval using binary features," in *Document Recognition and Retrieval XI, SPIE, Bellingham, WA*, **5296**, pp. 45–53, 2004.
9. B. Zhang and S. N. Srihari, "Binary vector dissimilarity measures for handwriting identification," in *Document Recognition and Retrieval X, SPIE, Bellingham, WA*, **5010**, pp. 28–38, 2003.
10. H. Sakoe, "Two-level dp-matching - a dynamic programming-based pattern matching algorithm for connected word recognition," in *IEEE Trans. Acoustics, Speech, and Signal Processing*, **27**, pp. 588–595, Dec. 1979.
11. C. S. Myers and L. R. Rabiner, "A level building dynamic time warping algorithm for connected word recognition," in *IEEE Trans. Acoustics, Speech, and Signal Processing*, **29**, no.2, pp. 284–297, 1981.
12. H. Ney, "A comparative study of two search strategies for connected word recognition: Dynamic programming and heuristic search," in *IEEE Trans. Pattern Analysis and Machine Intelligence*, **14**, no. 5, pp. 586–595, 1992.