

# Spotting Words in Latin, Devanagari and Arabic Scripts

Sargur N. Srihari, Harish Srinivasan, Chen Huang and Shravya Shetty

{srihari,hs32,chuang5,sshetty}@cedar.buffalo.edu

Center of Excellence for Document Analysis and Recognition (CEDAR), Buffalo, NY, USA

## Abstract

*A system for spotting words in scanned document images in three scripts, Devanagari, Arabic and Latin is described. Three main components of the system are a word segmenter, a shape based matcher for words and a search interface. The user gives a query which can be either a word image or text. The candidate words that are searched in the documents are retrieved and ranked, where the ranking criterion is a similarity score between the query and the candidate words based on global word shape features. This renders the word spotting technique to be independent of the script used. The performance of system is seen to be better for printed text as compared to handwritten. For handwritten English, a precision of 60% was obtained at a recall of 50%. An alternate approach comprising of prototype selection and word matching, that yields a better performance for handwritten documents is also discussed. For printed Sanskrit documents, a precision as high as 90% was obtained at a Recall of 50%.*

## 1 Introduction

Word spotting is a content-based information retrieval task to find relevant words within a repository of scanned document images. The retrieval of matching words finds applications in several areas like forensics, historical documents, personal records, etc. In these applications it is of interest to retrieve words from a large database of documents based on the visual and

the textual content.

Spotting handwritten words in documents written in the Latin alphabet has received considerable attention [1],[2] and [3]. [1] use pseudo 2-D hidden Markov models to represent keywords, [2] use hierarchical shape models and [3] extract profile-based holistic shape features from a line or word image and use dynamic time warping (DTW). A comparison of seven methods showed that the highest precision was obtained by a method based on profiles and DTW [4]. A word shape based method was shown to perform better than the DTW method, in terms of efficiency and effectiveness [5] – over 11% in precision and recall and 893 times faster.

This paper discusses the performance of the word shape method presented in [5] when applied to other scripts (Devanagari, Arabic) and languages arising from these scripts such as Sanskrit/Hindi and Arabic/Urdu. The paper also compares the retrieval performances in these scripts and those between printed and handwritten documents. The word spotting technique presented has been found to be effective for all the three languages under consideration (Sanskrit, Arabic, English), and for both printed and scanned handwritten documents. The search capability of our system comes with an interactive graphical user interface that allows searches (i) within the current document (ii) from another opened document or (iii) from a collection of preprocessed indexed documents.

The word spotting method presented here involves the indexing of documents that involves (i) Line and word segmentation (ii) Feature extraction: Global word shape features known as GSC features [6] are extracted for each of the words in the documents. The similarity between word images is measured using a correlation similarity measure. The system that is mentioned in this paper includes functionalities obtained from the CEDARABIC system[7], and the CEDAR-FOX system, which was developed for forensic examination [8], [9].

The organization of the rest of the paper is as follows. Section 2 describes the word spotting technique, including the indexing of the documents and the computation of the similarity between words. Section 3 describes the dataset used and the experiments conducted. Also, the retrieval performance for word spotting in each of the three languages has been discussed and evaluated in this section. Section 4 concludes the paper.

## 2 Word Spotting Technique

The word spotting technique involves the segmentation of each document into its corresponding lines and words. Each document is indexed by the visual image features of its words. The indexing technique used by us is the same for all the three languages: English, Sanskrit and Arabic.

### 2.1 Preprocessing

In the image preprocessing stage, the scanned image of each of the documents is subjected to several processing functions that leads to the segmentation of the document into lines and the words in each of these lines. This information regarding the number of lines and words and their positions in their respective documents is stored along with the document image. Figure 1 shows the line and word segmentation for handwritten English, printed Devanagari and handwritten

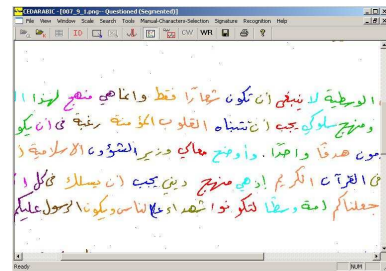
Arabic text. The line segmentation is performed using a clustering method. For word segmentation, the problem is formulated as a classification problem as to whether or not the gap between two adjacent connected components in a line is word gap or not. An artificial neural network with features characterizing the connected components was used to make a decision on this classification problem.



(a) English



(b) Sanskrit



(c) Arabic

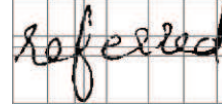
Figure 1: Segmentation examples where segmented words are coloured differently: (a) handwritten English (b) printed Sanskrit and (c) handwritten Arabic

## 2.2 Feature Extraction

The next step in the word spotting technique involves the computation of features for each of the words identified. These features form the key index for image search. The features used here for each word image, are the Gradient, Structural and Concavity (GSC) features which measure the image characteristics at local, intermediate and large scales and hence approximate a heterogeneous multi resolution paradigm to feature extraction. The features are extracted under a  $4 \times 8$  division and contain 384 bits of gradient features, 384 bits of structural features and 256 bits of concavity features, giving us a binary feature vector of length 1024 [6]. Figure 2 shows examples of the region division for words in all three languages. 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. gradient map to capture middle-level geometric characteristics including the presence of corners and lines at several directions. The structural features represent 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 [10]. The concavity features capture the major topological and geometrical features including direction of bays, presence of holes, and large vertical and horizontal strokes.

## 2.3 Measure for Word Spotting

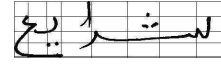
The distance between a word to be spotted and all the other words in the documents being matched against is computed using a normalized correlation similarity measure. This similarity measure is used to measure the similarity between two word images whose shapes are represented using the 1024-bit binary feature vectors described above. The Correlation measure was chosen after evaluating several other distance measures [6]. This dissimilarity measure



(a)



(b)



(c)

Figure 2: Examples of word images with  $4 \times 8$  division for feature extraction: (a)English (b)Sanskrit and (c)Arabic

$D$  between two feature vectors  $X$  and  $Y$  is given by equation 1.

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

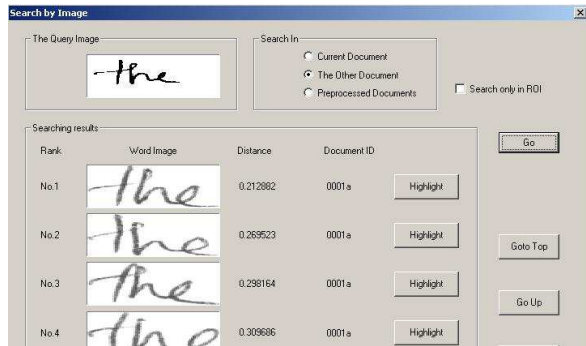
Here,  $S_{ij}, i, j \in \{0, 1\}$ , is defined as the number of occurrences of pattern  $i$  in the first feature vector and pattern  $j$  in the second feature vector at the corresponding positions. The smaller the value of  $D$ , the greater is the similarity between the 2 words. Thus, the results returned are ranked in the order of this similarity measure.

## 2.4 Word Spotting

### 2.4.1 Image as query: Method 1

Here the query is a word image and the goal is to retrieve documents which contain this word image. For this the system supports three kinds of word-spotting searches : 1) search for all relevant words within the same document 2) search for all relevant words in another document which is currently open 3) search for all relevant words from a collection of preprocessed documents. Each

of the retrieved word images is also linked with its corresponding document ID, which allows the user to easily retrieve its location and the document it belongs to. The queried word is selected as an image from the document containing it. The similarity measure for binary feature vectors (calculated as described in Section 2.3) is used to match the queried image against all the words in the selected documents. The results returned are word images and they are ranked according to their distances to the queried image. Figure 3 shows screen shots of the system after doing word spotting for words in English and Sanskrit.



(a) English Search with an image query of the word "the"



(b) Sanskrit Search with an image query of the word "paathu"

Figure 3: Screen shot of the system when performing word spotting search in (a) English and (b) Sanskrit.

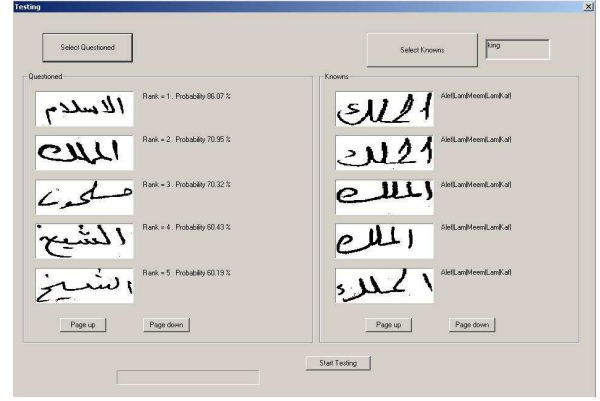


Figure 4: Arabic Search using Text Query. The prototype word images on the right side, obtained in the first step, are used to spot the images shown on the left. The results are sorted in rank by probability.

## 2.4.2 Text as query: Alternate Method

This method of retrieval was used for spotting Arabic words, with the idea of using the English equivalent meaning of that Arabic word as the query. Here the query is text and the goal is to retrieve the documents containing the word. The search phase has two parts: (i) Prototype selection: Here the first part of the typed in query is used to generate a set of query images corresponding to handwritten versions of it. (ii) Word Matching: Here each query image is matched against each indexed image to find the closest match. The purpose of such a two step search is the following (i) The prototype selection step provides for a set of handwritten versions of the query, to account for the different ways in which the word can be written. And hence this improves the word matching capability across different styles of writing. (ii) Also it enables the query to be in text and in a different language such as English. The two step approach is described below.

1. Prototype selection: Prototypes which are handwritten samples of a word are obtained from an indexed (segmented) set of documents. These indexed documents contain the truth (English equivalent) for every word image. Synonymous words if present

in the truth are also used to obtain the prototypes. Hence queries such as "country" will result in selecting prototypes that have been truthed as "country" or "nation" etc... A dynamic programming Edit Distance algorithm is used to match the query text with the indexed word image's truth. Those with distance as zero are automatically selected as prototypes. Others can be selected manually. This step can be thought of as a query expansion step where the text query is mapped into a number of query word images.

2. Word matching: Now that there are prototype images, these can be used as queries in a similar manner as described in section 2.4.1. For word matching, each word image in the test set of documents is compared with every selected prototype and a distribution of similarity values is obtained. The distribution of similarity values is replaced by its arithmetic mean. Now every word is sorted in rank in accordance with this final mean score.

Figure 4 shows the screenshot for word spotting in Arabic. The images on the right hand side are a result of the prototype selection and represent the different styles and the images on the left are the results obtained on matching the prototype word images.

## 3 Experiments and Results

In this section we describe the experiments and the results obtained for Word Spotting in Latin, Devanagari and Arabic.

### 3.1 Dataset

#### 3.1.1 Latin(English)

The data sets used for Word Spotting in English are subsets of a document image collection consisting of 3000 samples written by 1000

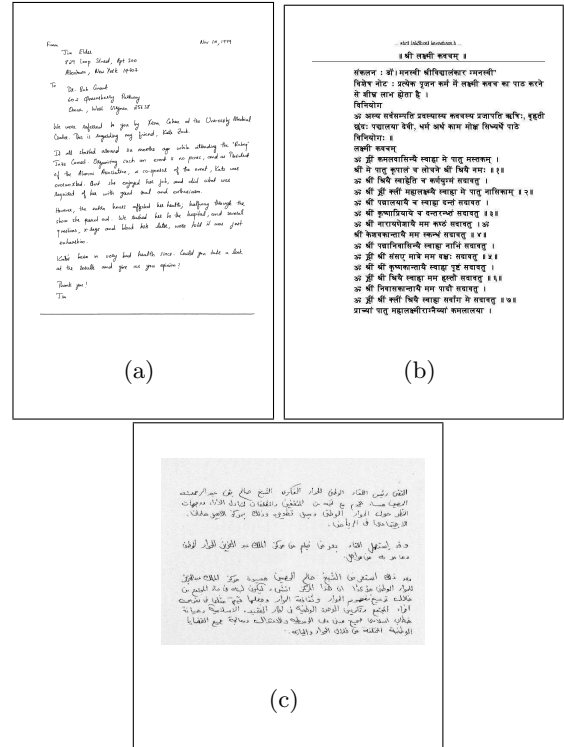


Figure 5: Sample documents: (a) English, (b) Sanskrit and (c) Arabic

writers, where each writer wrote 3 samples of a pre-formatted letter called "CEDAR Letter". The documents used for this experiment were randomly selected from this set. Each document consists of about 150 words.

#### 3.1.2 Devanagari(Sanskrit)

The data set used for Word Spotting in Sanskrit consist of 18 documents with printed text in Devanagari script which were randomly selected from <http://sanskrit.gde.to/> . Each document consists of a different text and the number of words vary from 45 to 250.

#### 3.1.3 Arabic

A document collection was prepared with 10 different writers, each contributing 10 different full page documents in handwritten Arabic. Each

document comprises of approximately 150-200 words each, with a total of 20,000 word images in the entire database. For each of the 10 documents that were handwritten, a complete set of truth comprising of the alphabet sequence, meaning and the pronunciation of every word in that document was also given. The above completes the indexing(preprocessing) of these documents.

All the documents were scanned in 8-bit gray scale (in other words, 256 shades of gray) and using 300 dpi resolution. Figure 5 shows portions of three of the sample documents used. All the document images were first automatically segmented into lines and words during the preprocessing stage as discussed in Section 2. Precision-recall curves have been used to evaluate the word spotting performance in each of the three languages. Precision is the ratio of the number of relevant images retrieved to the total number of images retrieved and recall is ratio of the number of relevant images retrieved to the total number of relevant images present.

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

Table 1: Performance of Word Spotting.

### 3.2 English Word Spotting

The word spotting in English documents was tested for by searching for several different word images. Each word image selected was searched for in another document written by the author of the queried image. A total of 100 queries from different documents were performed. The results are represented as a percentage of the number of times the correct match was returned (i) as the

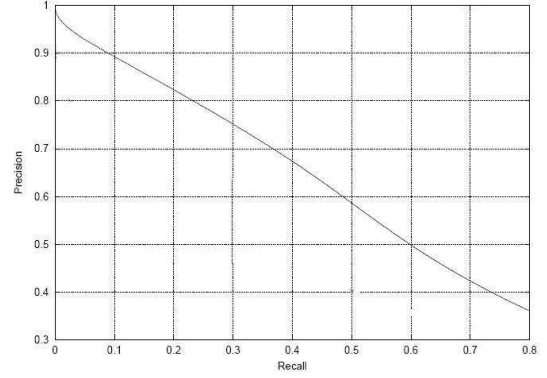


Figure 6: Average Precision Recall Curve for handwritten English

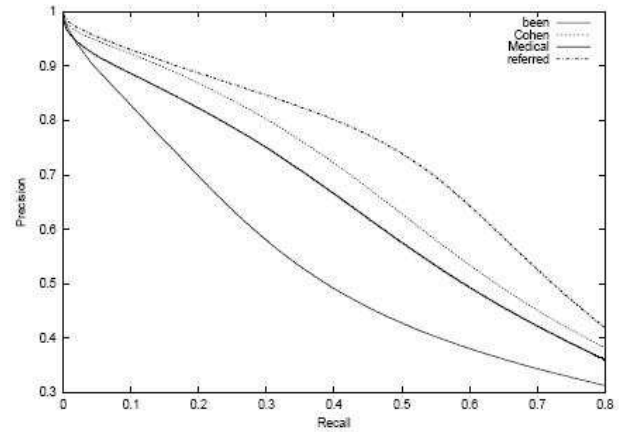


Figure 7: Precision Recall Curve for the four English words.

top rank, (ii) within top 5 ranks, (iii) within top 10, (iv) within top 20, and (v) within top 50. Table 1 shows the evaluation results.

In addition, we also calculated the Precision-Recall values for four randomly selected English words which occurred atleast twice in each document, "referred", "Cohen", "been" and "Medical". Figure 6 shows the average Precision-Recall curve for these four words and Figure 7 shows the individual Precision-Recall curves for these four words.

Words	Rank 1	Within Rank 2	Within Rank 5	Within Rank 10	Within Rank 20
yah	1	2	5	8	9
hoye	1	2	5	9	9
sadavatu	1	2	5	7	7
prabhu	1	2	5	8	8
tum	1	2	5	10	10

Table 2: Results for Word Spotting of Sanskrit words occurring 10 times

### 3.3 Sanskrit Word Spotting

The testing for Word Spotting in Sanskrit documents was done by searching for several randomly selected Sanskrit words in 18 documents. Each of the words tested occurred atleast 5 times in the documents with varying frequencies. Table 2 displays the ranks of the retrieved results for 5 Sanskrit words which occurred 10 times in the selected documents. We have also calculated the precision and recall values for each of the queried words. Figure 8 displays a precision recall curve of the average precision recall values of all the queried words.

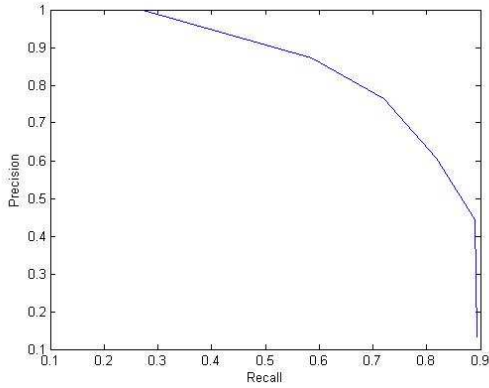


Figure 8: Average Precision Recall Curve for printed Sanskrit words.

### 3.4 Arabic word spotting

The performance of the word spotter was evaluated using manually segmented Arabic docu-

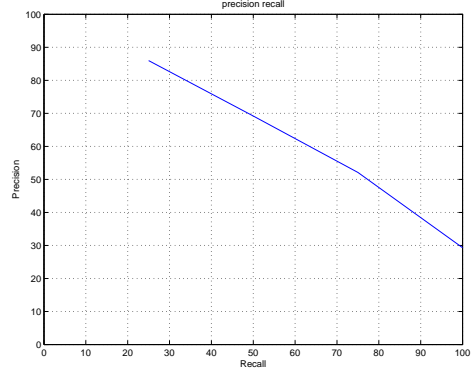


Figure 9: Retrieval of Arabic words. All results were averaged over 150 queries. The styles of 8 writers were used for training (by using their documents) to provide template word images.

ments. All experiments and results were averaged over 150 queries. For each query, a certain set of documents are used as training to provide for the prototype word images and rest of the documents were used for testing. When documents from 8 writers were used to provide prototype styles, the corresponding Precision-Recall curves are shown in Figure 9. When the number of writers used to provide prototype styles are increased, the precision at a given recall can be seen to increase in Figure 10. The reason for this is intuitive as more writers provide for more prototype images that correspond to learning a greater number of styles in which the query word can be written.



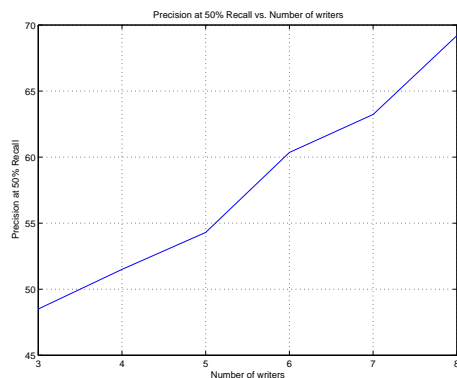


Figure 10: Precision at 50% recall is plotted against different numbers of writers used for training. As more number of writers are used for training, the precision at a given recall increases.

## 4 Conclusion

Content-based search in documents using Latin, Devanagari and Arabic scripts has been described here. Word images were used for the retrieval of document data in each of these languages. The performance of the system in these various types of searches has been presented as Precision-Recall curves and the ranks of the returned results. A very high precision can be obtained for spotting words from printed documents. A two step approach involving a query expansion step yields promising results for spotting words in handwritten documents.

## References

- [1] S. Kuo and O. Agazzi, "Keyword spotting in poorly printed documents using 2-d hidden markov models," in *IEEE Trans. Pattern Analysis and Machine Intelligence*, 16, pp. 842–848, 1994.
- [2] M. Burl and P. Perona, "Using hierarchical shape models to spot keywords in cursive handwriting," in *IEEE-CS Conference on Computer Vision and Pattern Recognition*, June 23–28, pp. 535–540, 1998.
- [3] A. Kolz, J. Alspector, M. Augusteijn, R. Carlson, and G. V. Popescu, "A line-oriented approach to word spotting in handwritten documents," in *Pattern Analysis and Applications*, 2(3), pp. 153–168, 2000.
- [4] R. Manmatha and T. M. Rath, "Indexing of handwritten historical documents-recent progress," in *Symposium on Document Image Understanding Technology (SDIUT)*, pp. 77–85, 2003.
- [5] B. Zhang, S. N. Srihari, and C. Huang, "Word image retrieval using binary features," in *Document Recognition and Retrieval XI, SPIE, San Jose, CA*, 2004.
- [6] 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.
- [7] S. N. Srihari, H. Srinivasan, P. Babu, and C. Bhole, "Handwritten arabic word spotting using the cedarabic document analysis system," in *Proc. Symposium on Document Image Understanding Technology (SDIUT-05)*, pp. 123–132, 2005.
- [8] 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.
- [9] 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.
- [10] J. Favata, G. Srikantan, and S. N. Srihari, "Handprinted character/digit recognition using a multiple feature/resolution philosophy," in *The Fourth International Workshop on Frontiers in Handwriting Recognition (IWFHR 4)*, pp. 57–66, 1994.