

Handwritten Arabic Word Spotting using the CEDARABIC Document Analysis System

Sargur Srihari, Harish Srinivasan, Pavithra Babu and Chetan Bhole

Center of Excellence for Document Analysis and Recognition (CEDAR)

University at Buffalo, State University of New York

Amherst, New York 14228

srihari @cedar.buffalo.edu

Abstract

An algorithm and a system for searching handwritten Arabic documents to locate key words is presented. Three main components of the system are a word segmenter, a shape based matcher for words and a search interface. The user types in a query in English within a search window, the system finds the equivalent Arabic word, e.g., by dictionary lookup, locates word images in an indexed (segmented) set of documents. A two-step approach is employed in performing the search: (1) prototype selection: the query is used to obtain a set of handwritten samples of that word from a known set of writers (these are the prototypes), and (2) word matching: the prototypes are used to spot each occurrence of those words in the indexed document database. A ranking is performed on the entire set of test word images—where the ranking criterion is a similarity score between each prototype word and the candidate words based on global word shape features. A database of 20,000 word images contained in 100 scanned handwritten Arabic documents written by 10 different writers was used to study retrieval performance. Using five writers for providing prototypes and the other five for testing, using manually segmented documents, 55% precision is obtained at 50% recall. Performance increases as more writers are used for training.

1 Introduction

Spotting handwritten words in documents written in the Latin alphabet has received considerable attention [1],[2], [3]. A systematic comparison of seven methods was made in [4] who showed that the highest precision was obtained by their method based on profiles and dynamic time warping (DTW). A method based on word shape was shown to perform better than the method based on DTW, both in terms of efficiency and effectiveness, in [5]. This paper discusses the performance of the word shape method presented in [5] when applied to handwrit-

ten Arabic documents.

The word spotting algorithm is implemented as a tool within an interactive software system known as CEDARABIC which is designed for the analysis of handwritten Arabic documents. Most functionalities of CEDARABIC are the same as that found in a predecessor system known as CEDAR-FOX which was developed for forensic examination by U.S. law enforcement [6]. Being a software branch of the CEDAR-FOX system CEDARABIC also has capabilities for image processing operations such as enhancement, thresholding, underline removal, etc., as well as for writer verification, writer identification and signature verification. The interfaces have been modified to be appropriate for analyzing Arabic language documents by an English-speaking user. The system includes a GUI for various functions including scanning, selecting a region of interest, feature computation, ground truthing (manual segmentation and entering word truth), searching of a word based on its English equivalent, etc. The primary focus of this paper is on the word-spotting tool of CEDARABIC—whose goal is to find each occurrence of a user-specified query word in a given set of documents.

In the CEDARABIC system the task of word spotting has two phases: indexing and search. In the indexing phase each document image in the database is indexed by segmenting into its component word images. The search phase has two parts: in the first part the typed in query is used to generate a set of query images corresponding to handwritten versions of it. In the second part each query image is matched against each indexed image to find the closest match.

The paper consists of the following sections: (i) scanned document collection used to design the system and test its performance, (ii) user interface of CEDARABIC, (iii) word segmentation algorithm, (iv) word-shape matching algorithm, (v) conclusion and acknowledgment.

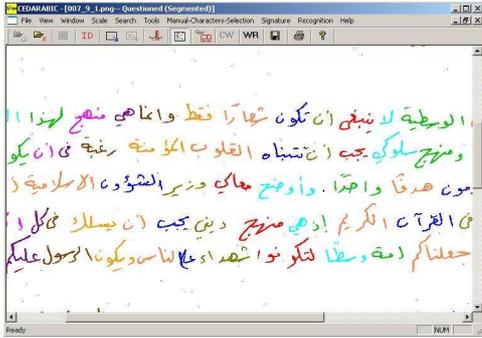


Figure 1: Sample “.arb” file opened with CEDARABIC. Each segmented word is colored differently from adjacent words.

2 Document Image Database

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. The documents were scanned at a resolution of 300 dots per inch. This is the resolution to be used for optimal performance of the system. Figure 1 shows a sample scanned handwritten Arabic document written by writer 1.

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 scanned handwritten document comprising of word images were mapped with the corresponding truth information given, and saved as a separate file with a “.arb” extension. The naming convention for the file is explained with this example. File 007_9.1.arb is the ninth document written by writer number 7. When some writers use more than one full page to transcribe a document, the last digit 1 or 2 identifies the page of the document. These “.arb” files have all the necessary information in them to perform word spotting using the tool that is described in the following section. The preparation of these “.arb” files involved careful manual word segmentation of each scanned document. The software tool described below also provides for a way to manually segment each word from a raw unprocessed document.

3 User Interface

Most of the functionalities of the CEDARABIC system, including word spotting, are the same as that found in a system for Latin script, known as CEDAR-FOX, which was developed for forensic examination by U.S. law enforcement [7] [8]. Being a modification of the CEDAR-FOX system for the analysis of handwritten documents in the Latin alphabet, CEDAR-FOX also has capabilities for image

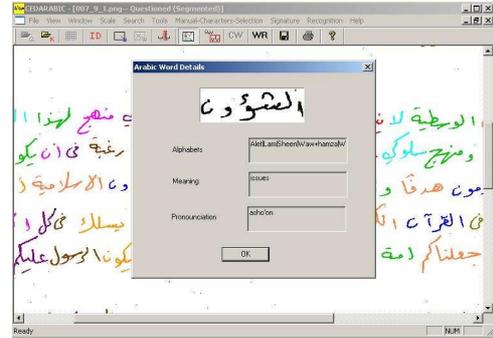


Figure 2: Truth for a particular word. This window pops up when a particular word image is right-clicked. The truth consists of three parts: the individual characters, its meaning in English and its Arabic pronunciation shown in English.

processing operations such as enhancement, thresholding, underline removal, etc., as well as for writer verification, writer identification and signature verification. The interfaces have been modified to be appropriate for analyzing Arabic language documents by an English-speaking user.

CEDARABIC has a GUI for analyzing and using the prepared “.arb” for performing word spotting and other information lookup tasks. A particular “.arb” file can be opened by selecting the “Q” or “K” option on the tool bar and Figure 1 shows a screen shot of CEDARABIC when a particular file is opened. By clicking on the tool bar icon for “Show Words”, the different words get marked with different colors thereby differentiating one Arabic word image from another.

Right-clicking on any word image pops up a dialog- box displaying the complete set of truths associated with that word. Figure 2 shows a screenshot when a particular word image was right-clicked. The Arabic alphabet sequence, the English equivalent meaning, and the pronunciation of each word can be viewed with this tool.

The query for word spotting is typed as English text– one that corresponds to the English equivalent meaning of the word to be searched for.

The first step in word spotting is to identify prototypes (template word images) corresponding to the Arabic word for the queried text. These template word images are obtained from documents which are set aside as training documents. For example, we can say documents pertaining to writers 1 through 5 can be used solely for the purpose of training and they provide for the template word images. The default set consists of writers 1 through 5. These training “.arb” should contain the complete truth for each word image, so as to pick up those that match the query text. The word images for those that matched serve as template word images to per-

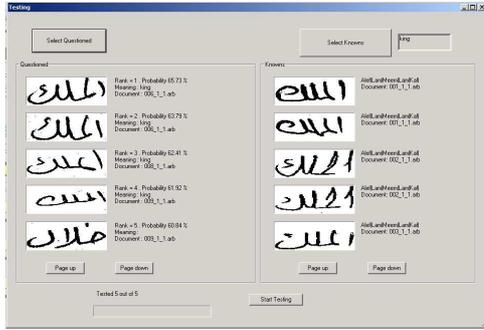


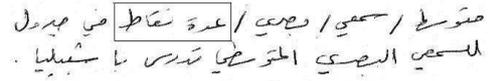
Figure 3: Word spotting results. The template word images on the right side, obtained in the first step, are used to spot the images shown on the left. The results are sorted in rank by probability. The images on the right are the selected prototypes and those on the left are the ranked set of test images.

form word spotting in other documents that do not have the truth. The matching words are automatically selected and the user is also given an option of selecting the ones that need to be used as template word images. Once this is done, upon clicking “Use Selected for Training” these images will be used to spot similar occurrences for words in other documents purely based on the shape of the word. The template word images serve as a knowledge base to learn the writing styles of the particular word and then use this to identify other similar words purely based on matching the image. A set of test documents can now be specified, and the tool will look for the occurrence of the particular word in every document specified and rank the retrieved words in the order of probability of match. Figure 3 shows the result of one such word spotting task performed. In this example the search is for the word “king” whose Arabic equivalent is pronounced as “al malik”. The images on the right are the template word images and those on the left are the word images in order by rank.

4 Word Segmentation

Automatic segmentation into words is based on taking several features on either side of a potential segmentation point and using a neural network for deciding whether or not the segmentation is between two distinct words. Some of the differences between the tasks of segmenting Arabic script and segmenting Latin script are: presence of multiple dots above and below the main body in Arabic and the absence of upper case letters at the beginning of sentences in Arabic.

The process of word segmentation begins with obtaining the set of connected components for each line in the document image. Figure 4 shows a region of a document image the connected components (exterior and interior contours) of a small section. The



(a) Region of a document image

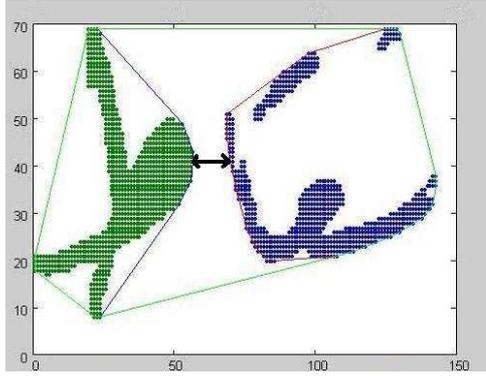


(b) Connected component exterior and interior contours

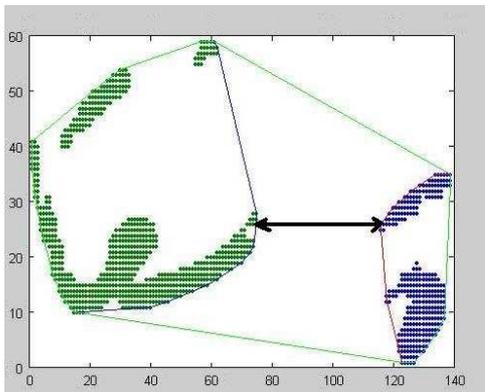
Figure 4: Connected components of a small section of the image.

interior contours or loops in a component are ignored for the purpose of word segmentation as they provide no information for this purpose. The connected components are grouped into clusters, by merging minor components such as dots above and below a major component. Also particular to Arabic, most words starts with the Arabic character “Alef”. The presence of an “Alef” is a strong indicator that there may be a word gap between the pair of clusters. The height and width of the component are two parameters used to check if the component is the character “Alef”. Figure 6 shows samples of the Arabic character Alef. Every pair of adjacent clusters are candidates for word gaps. 9 features are extracted for these pairs of clusters and a neural network is used to determine if the gap between the pair is a word gap. The 9 features are: width of the first cluster, width of second cluster, difference between the bounding box of the two clusters, flag set to 1 or 0 depending on the presence or absence of the Arabic character “Alef” in the first cluster, the same flag for the second cluster, number of components in the first cluster, number of components in the second cluster, minimum distance between the convex hulls enclosing the two clusters and the ratio between, the sum of the areas enclosed by the convex hulls of the individual clusters, to the total area inside the convex hull enclosing the clusters together. Figure 5 shows two pairs of adjacent clusters tested for word gaps.

A neural network was trained using these 9 features with feature vector labeled as to whether it is a word gap or not. This is similar to the neural network approach used for English postal addresses [9].



(a) Not word gap



(b) Word gap

Figure 5: Two pairs of adjacent clusters. Dotted lines indicate convex hulls around the individual clusters. Solid lines indicate convex hull around the two clusters taken together. (a) The gap between the convex hulls is truly not a word gap, and (b) The gap between the convex hulls is truly a word gap.

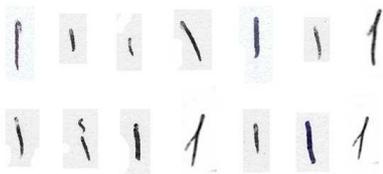


Figure 6: Samples of Arabic character "Alef". The height and width are two parameters that are used to detect the presence of "Alef" in the clusters.



Figure 7: A document segmented into words by the neural network. Figure 8(a) 8(b) 8(c) and 8(d) show the first four clusters of the document for the first line from the left and the classification result.

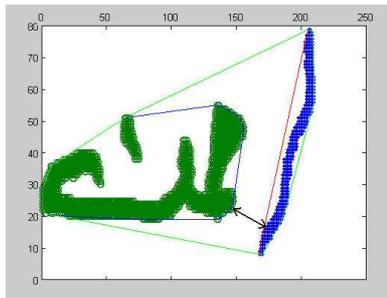
Table 1: NN Features

| Fig | f1 | f2 | f3 | f4 | f5 | f6 | f7 | f8 | f9 | NN |
|------|-----|----|----|----|----|----|----|-------|------|-------|
| 8(a) | 156 | 42 | 12 | 0 | 0 | 2 | 1 | 24.59 | 0.58 | -0.75 |
| 8(b) | 42 | 13 | 7 | 0 | 1 | 1 | 1 | 8.00 | 0.52 | -0.99 |
| 8(c) | 13 | 84 | 14 | 1 | 0 | 1 | 2 | 24.75 | 0.40 | 0.91 |
| 8(d) | 84 | 82 | 51 | 0 | 0 | 2 | 2 | 51.00 | 0.56 | 0.80 |

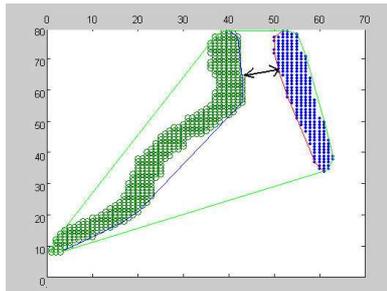
However the features are different. Figure 7 shows a document segmented into words by the neural network. Figure 8 shows the first 4 pairs of clusters whose features are extracted and fed to the neural network. These 4 pairs of clusters correspond to the words starting from the left on the first line of the document in figure 7. Table 1 shows the 9 features for these pairs respectively and the last column shows the output of the neural network and table 2 shows what each feature is. A positive output from the neural network indicates a presence of word gap and a negative output indicates the absence of a word gap. From the table 1, it is evident that though the distance between the clusters for pair 1 and pair 3 are almost the same, the neural network gives the correct classification. This is due to presence of "Alef" that strongly indicates the presence of a word gap.

Table 2: Features

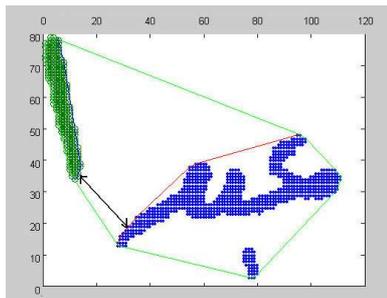
| Number | Feature |
|--------|---|
| 1 | Width of 1 st cluster |
| 2 | Width of 2 nd cluster |
| 3 | Bounding box gap between clusters |
| 4 | Alef flag for 1 st cluster |
| 5 | Alef flag for 2 nd cluster |
| 6 | No of components 1 st cluster |
| 7 | No of components 2 nd cluster |
| 8 | Min. Convex hull distance |
| 9 | Ratio of sum of cluster areas to total area |



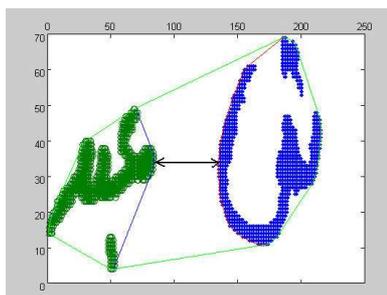
(a) Not word gap



(b) Not word gap



(c) Word gap



(d) Word gap

Figure 8: Four pairs of adjacent clusters. Note that the distance between the clusters in (a) and (c) are the same (the scale in the plot are different). The neural network still classifies (c) as a word gap due to the presence of the “Alef”.



Figure 9: Example of a reasonably well-segmented Arabic document by automatic segmenter designed for English. The truth is shown on the right screen-shot and the result of the automatic word segmenter is shown on the left. In this case 86% of the automatically segmented words are correct.

4.1 Word Segmentation Performance

An example of a reasonably well-segmented Arabic image by the automatic segmenter, together with the perfectly segmented image (truth), are shown in Figure 9. There are 130 words in the truth shown on the right-hand side. The result of the automatic segmenter, shown on the left, has the following error counts: (i) 9 words have been split into two parts (9 errors), (ii) 3 adjacent pairs of words have been merged (6 errors) and (iii) one set of three words have been merged (3 errors). Since there are a total of 18 errors we can say that 86% of the word segmentations are correct.

An example of a poorly segmented image together with its truth are shown in Figure 10. In this case there are 120 words in the truth shown on the right. the error counts of the automatically segmented document on the left are as follows: (i) 17 words have been split (17 errors), (ii) 18 adjacent words have been merged (36 errors), (iii) 3 merges of three adjacent words (9 errors), and (iv) one merge of four adjacent words (4 errors). Since there are 66 errors we can say that 45% of the word segmentations are correct.

Overall performance over ten writers writing 10 documents each, when the lines are segmented correctly, is about 60% using a set of seven segmentation features. A more complex set of features is expected to yield a higher level of performance.



Figure 10: Example of a poorly segmented Arabic document by automatic English segmenter. Truth is on the right and the result of automatic segmentation is on the left. In this case 45% of the segmented words are correct.

5 Word Shape Matching

The word segmentation is an indexing step before using the word shape matching for word retrieval. A two-step approach is employed in performing the search: (1) prototype selection: the query (English text) is used to obtain a set of handwritten samples of that word from a known set of writers (these are the prototypes), and (2) word matching: the prototypes are used to spot each occurrence of those words in the indexed document database. A ranking is performed on the entire set of test word images where the ranking criterion is the mean similarity score between prototype words and the candidate word based on global word shape features.

5.1 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.

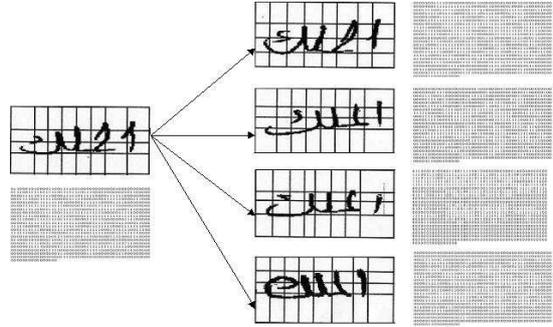


Figure 11: One query word image (left) matched against 4 selected prototype images. The 1024 binary GSC features are shown next to the images.

5.2 Word Matching

The word matching algorithm uses a set of 1024 binary features for the word images. These binary features are compared using the correlation similarity measure 5.2.1 to obtain a similarity value between 0 and 1. This similarity score represents the extent of match between two word images. The smaller the score, the better is the match. For word spotting, every word image in the test set of documents are 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 5.2 shows an example query word image compared with the a set of 4 selected prototypes.

5.2.1 Similarity measure

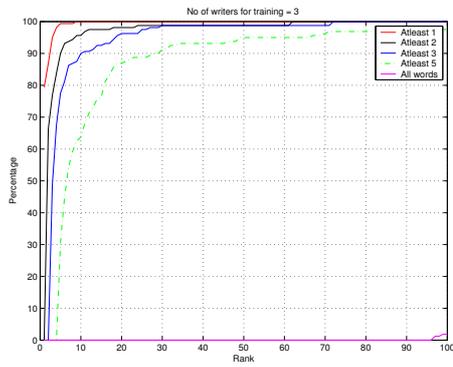
The method of measuring the similarity or distance between two binary vectors is essential. The correlation distance performed best for GSC binary features [10] which is defined for two binary vectors X and Y , as in equation 5.2.1

$$d(X, Y) = \frac{1}{2} \left(1 - \frac{s_{11}s_{00} - s_{10}s_{01}}{[(s_{10} + s_{11})(s_{01} + s_{00})(s_{11} + s_{01})(s_{00} + s_{10})]^{\frac{1}{2}}} \right)$$

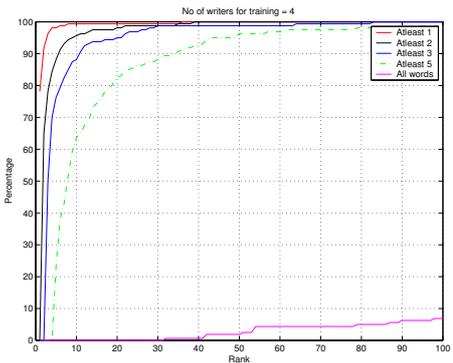
where s_{ij} represent the number of corresponding bits of X and Y that have values i and j .

5.3 Word-Matching Performance

The performance of the word spotter was evaluated using manually segmented Arabic documents. 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 template word images and rest of the documents were used for testing. Figure 12 shows the percentage of times when correct word images were retrieved within the top ranks scaled on the x-axis. More than one correct word image does exist in the test set. For example, if the word “king” occurs in document 1 twice, and we use writers 1 through 6 for training, then



(a) Recall with 3 writers for training.

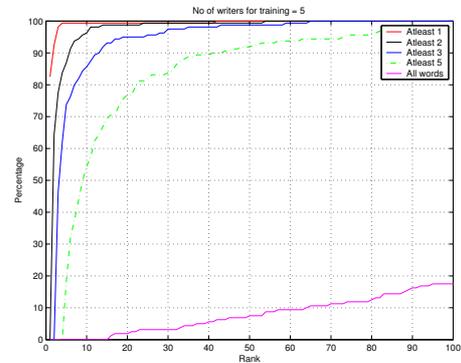


(b) Recall with 4 writers for training.

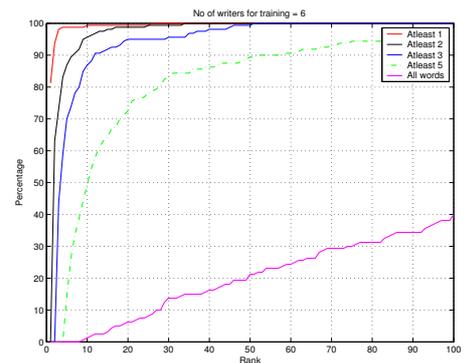
Figure 12: Retrieval performance of the system. All results were averaged over 150 queries. The styles of 3 and 5 writers were used for training (by using their documents) to provide template word images. The rest were used for testing. The different curves report on the percentage of times the word images were retrieved within the ranks scaled on the x-axis.

we have $6 * 2 = 12$ template word images to perform word spotting. The remaining test set contains $(10 - 6) * 2 = 8$ correct matches in the test set. The top curve of the figure shows that at least one correct word is retrieved 90% of the time within the top 12 ranks. The other curves show the same when at least 2 and 3 correct words are retrieved. Similar plots can be obtained when different numbers of writers are used for training. These plots are shown in Figs. 12 13 and 14 where three, five and seven writers were used for training.

Performance of word spotting can be specified in terms of precision and recall. If n_1 is the number of relevant words in the database, n_2 is the number of words retrieved, and n_3 is number of relevant words among those that are retrieved then $Precision = n_3/n_2$ and $Recall = n_3/n_1$. Precision-Recall curves corresponding to using three, five and

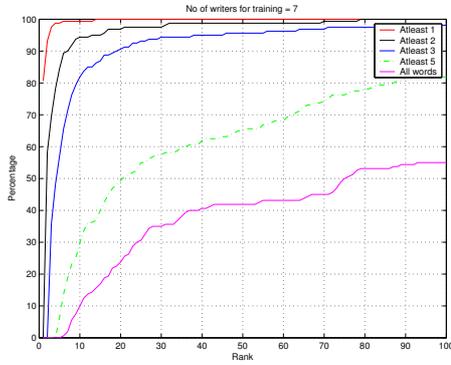


(a) Recall with 5 writers for training.

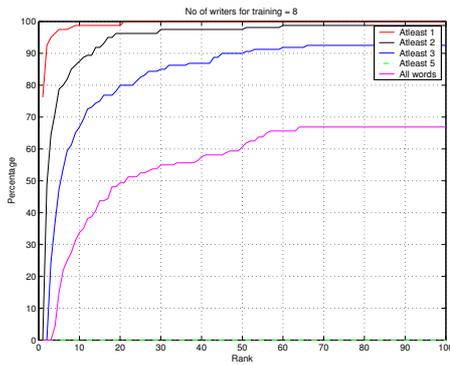


(b) Recall with 6 writers for training.

Figure 13: Retrieval performance of the system. All results were averaged over 150 queries. The styles of 5 and 6 writers were used for training (by using their documents) to provide template word images. The rest were used for testing. The different curves report on the percentage of times the word images were retrieved within the ranks scaled on the x-axis.

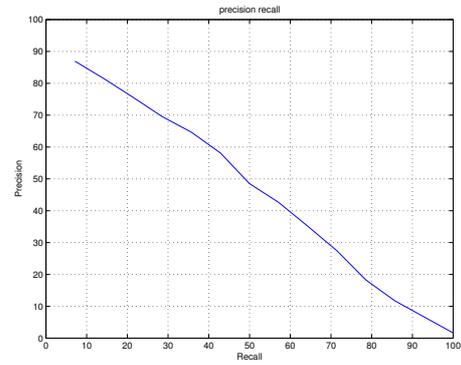


(a) Recall with 7 writers for training.

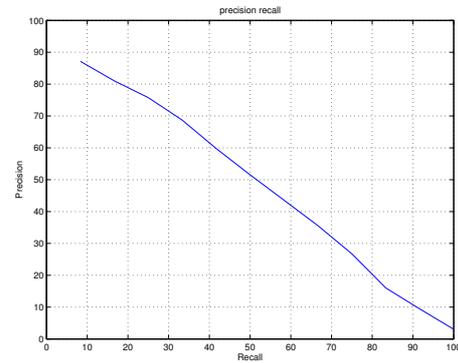


(b) Recall with 8 writers for training.

Figure 14: Retrieval performance of the system. All results were averaged over 150 queries. The styles of 7 and 8 writers were used for training (by using their documents) to provide template word images. The rest were used for testing. The different curves report on the percentage of times the word images were retrieved within the ranks scaled on the x-axis. Results show improved performance over the case when five styles were used as shown in Figure 12 and 13.



(a) Precision Recall with 3 writers for training.



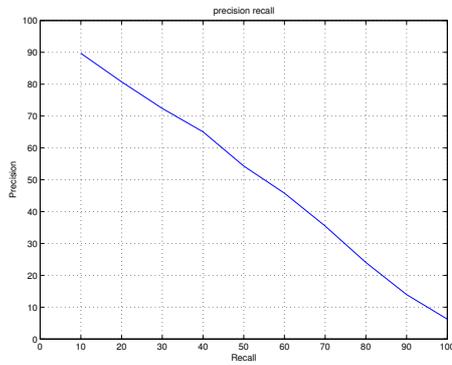
(b) Precision Recall with 4 writers for training.

Figure 15: Precision-Recall curves where the styles of 3 and 4 writers were used for training (by using their documents) to provide template word images.

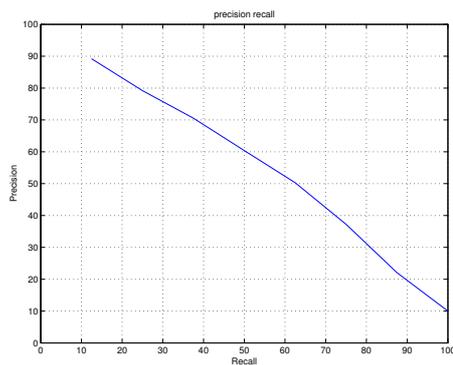
seven writer styles are shown in Figures 15 , 16 and 17 respectively. Finally figure 18 compares the precision at a particular recall with varying number of writers used for training. As the number of writers is increased, the precision at a given recall increases.

6 Conclusion

A technique for Arabic word image spotting in previously segmented images and that of automatic Arabic word segmentation was discussed. A two-step approach involving (i) Prototype selection and (ii) Word matching were described. Performance of the system is very promising with increased recall and better precision as the number of writers in the training set is increased. A precision of 70% is achieved at a recall of 50% when 8 writers were used for training. The CEDARABIC system described as the user interface provides tools for (i) manual word segmentation, (ii) information look-up on truth of every word, and (iii) interface to perform word spotting. A set

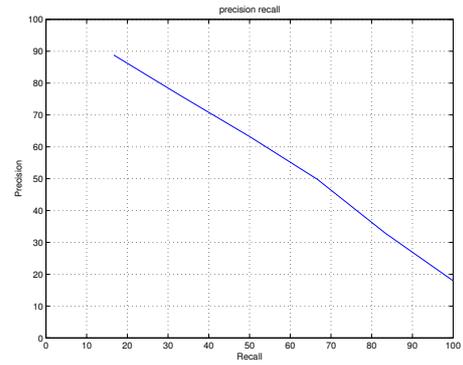


(a) Precision Recall with 5 writers for training.

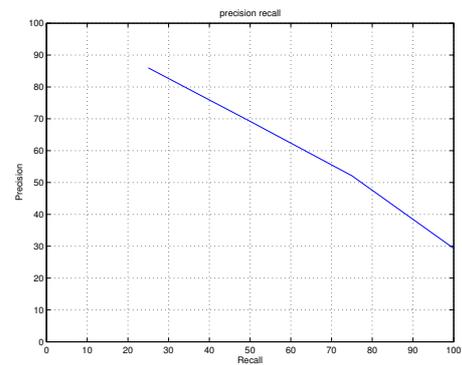


(b) Precision Recall with 6 writers for training.

Figure 16: Precision-Recall curves where the styles of 5 and 6 writers were used for training (by using their documents) to provide template word images.



(a) Precision Recall with 7 writers for training.



(b) Precision Recall with 8 writers for training.

Figure 17: Precision-Recall curves where the styles of 7 and 8 writers were used for training (by using their documents) to provide template word images. Results show an improvement over the case when five styles were used as shown in Figure 15 and 15.

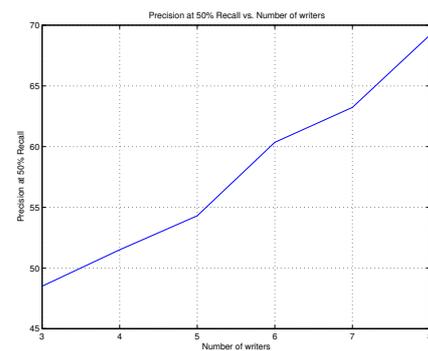


Figure 18: Precision at 50% recall is plotted against different numbers of writers used for training. As more number of writers are used for training, the precision at a given recall increases.

of more complicated features tuned for Arabic will result in better automatic word segmentation performance and it will be a part of future research.

7 Acknowledgment

The authors wish to thank Chen Huang, Vivek Shah, Manivannan Arivazhagan and Thanh Nguyen for their assistance.

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 Vol. 5296**, pp. 45–53, 2004.
- [6] S. N. Srihari, B. Zhang, C. Tomai, S. Lee, and Y.-C. Shin, "A system for handwriting matching and recognition," in *Symposium on Document Image Understanding Technology Greenbelt, MD*, April, 2003.
- [7] 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.
- [8] 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*, pp. 67–75, 2003.
- [9] G. Kim, V. Govindaraju, and S. N. Srihari, "A segmentation and recognition strategy for handwritten phrases," in *International Conference on Pattern Recognition, ICPR-13*, pp. 510–514, 1996.
- [10] B. Zhang and S. N. Srihari, "Binary vector dissimilarity measures for handwriting identification," *Proceedings of the SPIE, Document Recognition and Retrieval*, pp. 155–166, 2003.