

Word Image Retrieval Using Binary Features

Bin Zhang, Sargur N. Srihari and Chen Huang

CEDAR, Computer Science and Engineering Department
State University of New York at Buffalo, Buffalo, NY 14228

ABSTRACT

Indexing of handwritten documents for subsequent retrieval is still a challenge. Existing word image retrieval algorithms suffer from either low retrieval precision or high computation complexity. We present an effective and efficient approach for word image matching by using gradient-based binary features. Experiments over a large database of handwritten word images show that the proposed approach consistently outperforms the existing best handwritten word image retrieval algorithm, Dynamic Time Warping (DTW) with profile-based shape features. Not only does the proposed approach have much higher retrieval accuracy, but also it is 893 times faster than DTW.

Keywords: Word spotting, image retrieval, document retrieval, shape matching, binary features.

1. INTRODUCTION

Indexing of handwritten documents for subsequent retrieval is still a challenge. The key for indexing handwritten documents is to measure the similarity between word images.

Kuo and Agazzi¹ use pseudo 2-D hidden Markov models to represent a key word and all the other extraneous words. This algorithm works only under simulated environment with limited number of words as it requires additional training process to build up the statistical models. In real applications of keyword spotting query example can be any word and may be written in any style, so it is almost impossible to set up a training set before retrieval.

In the work by Burl and Perona,² keywords are represented by hierarchical shape models where a keyword is hierarchically decomposed into a set of segments described by lower-level features like humps and cusps etc.. The same as hidden Markov models in the work,¹ the hierarchical shape models have to be obtained through a number of training samples.

Kolcz et al.³ extract profile-based holistic shape features from a word or line image and utilize dynamic time warping (DTW) for pattern matching. The same word-spotting scheme was recently applied to indexing handwritten historical documents.⁴ In the work,⁴ Manmatha and Rath systematically compared seven word matching algorithms for purpose of word image retrieval and found that DTW using profile-based features presents the highest retrieval precision. However, DTW suffers from its high computation complexity even with speeding-up technique, and retrieving word images from large database of handwritten documents is still unacceptable.

In this study, we present an effective and efficient approach for word image matching by using binary-valued features. The binary features for a word image consist of 1024 bits corresponding to gradient (384 bits), structural (384 bits), and concavity (256 bits) features. The similarity between two binary vectors is computed by a Correlation-based measure. Word retrieval is experimented over 9,312 word images extracted from 2328 handwriting samples by 776 individuals in U.S.. Dynamic Time Warping using profile-based shape features is also implemented for comparison.

The rest of the paper is organized as follows. In Section 2, we describe binary-feature extraction for handwritten word images and present examples. In Section 3 we simply introduce the profile-based features for word images. In Section 4, we compare the proposed approach and DTW one by retrieval experiments. Section 5 concludes the paper.

2. BINARY FEATURES EXTRACTION FOR WORD IMAGES

The proposed feature extractor for word image is a modified version of a well-known gradient-based feature extractor for handwritten character images. We first describe the principle of the gradient-based feature extractor, then present several word image examples, finally define a similarity measure for binary features.

Further author information: send correspondence to Bin Zhang by e-mail binzhang@cedar.buffalo.edu

2.1. Gradient-based Binary Feature Extraction

Gradient-based binary features (known as GSC features) were developed for handwritten character recognition and yielded one of the best numeral/character recognizers.^{5,6} Gradient-based binary feature for a character image consist of 512 bits corresponding to gradient (192 bits), structural (192 bits), and concavity (128 bits) features. Each of these three sets of features rely on dividing the character image into a 4 x 4 region.

GSC features for handwritten alphanumeric characters were directly employed for handwriting identification and verification.^{7,8} GSC features for ten numerals and fifty-two alphabetic characters can correctly identify 97.83% of writers when 875 writers were considered. The GSC feature extractor was further extended to adapt to handwritten word images for purposes of studying handwriting individuality.⁹ After examining several different division schemes for word image, the work in⁹ found that 4 x 8 division was the best one and GSC features at 4 x 8 division from four characteristic words led to an identification accuracy of 83% and a verification rate of 90.94%. For a word image under 4 x 8 division, the corresponding GSC features contain 384-bit gradient features, 384-bit structural features, and 256-bit concavity features, leading to a binary vector of 1024 dimensions.

As handwriting identification is essentially an image retrieval task with an extra effort to determine authorship, the success of word-level GSC features in handwriting identification indicates that they will be the same effective for word image retrieval or spotting. In this study, we examine the performance of GSC features for the task of word image retrieval. GSC features are extracted from each word image under 4 x 8 division.

Figure 1 presents three examples of 4 x 8 division and their corresponding 1024-bit feature vectors.

2.2. Similarity Measure for Binary Feature Vectors

Let Ω be the set of all n -dimensional binary vectors. To measure the similarity between two binary vectors, we use the Correlation measure defined in.¹⁰ This measure has shown very good performance for recognition.¹¹

Let S_{ij} ($i, j \in \{0, 1\}$) be the number of occurrences of matches with i in the first pattern and j in the second pattern at the corresponding positions. Given two binary feature vectors $X \in \Omega$ and $Y \in \Omega$, each similarity measure $S(X, Y)$ above uses all or some of the four possible values, i.e., S_{00}, S_{01}, S_{10} and S_{11} . We define a dissimilarity measure $D^b(X, Y)$ corresponding to the Correlation measure as:

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

3. PROFILE-BASED SHAPE FEATURES

Profile-based shape features simply represent a 2-dimensional word image by four 1-dimensional signals, i.e., projection profile, upper profile, lower profile, and transition-count profile.^{3,4} Transition-count profile records number of ink-background transitions along each column of the image.

The distance between two sets of profile features is calculated by a dynamic time warping (DTW) algorithm.³ As the DTW algorithm has high computation complexity, in our experiments an optimization scheme proposed by Sokae and Chiba¹² was adopted to accelerate computation.

4. PERFORMANCE COMPARISON AND ANALYSIS

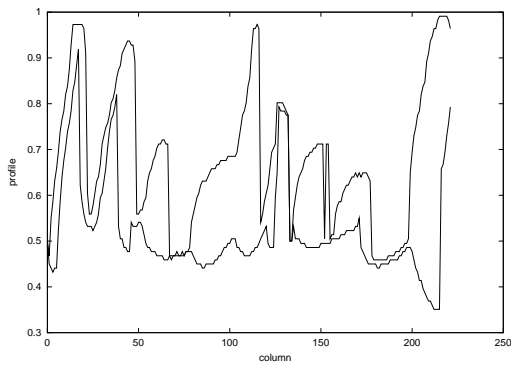
In this section, we describe the experimental settings, present experimental results with GSC and DTW and perform analysis.

Medical

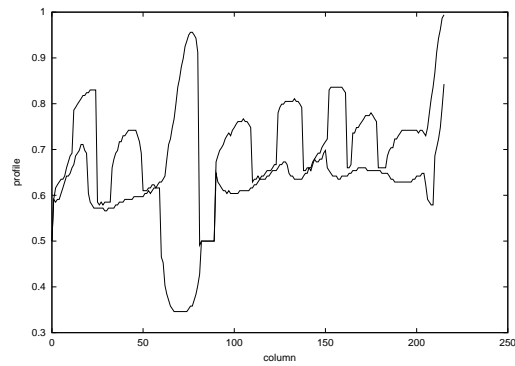
(a)

referred

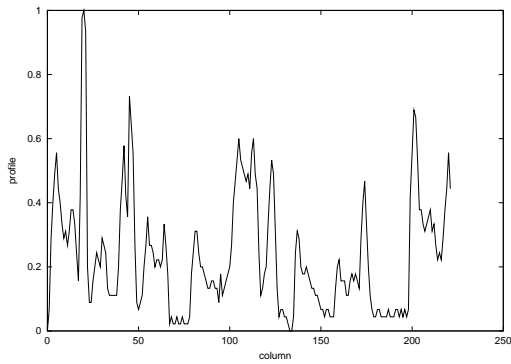
(e)



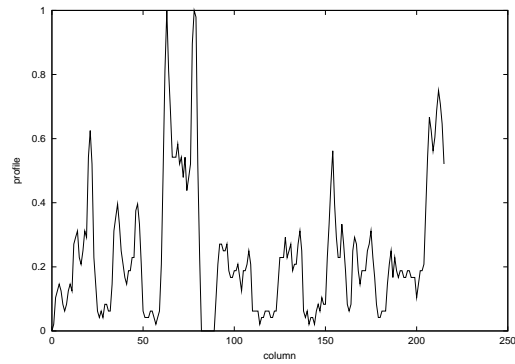
(b)



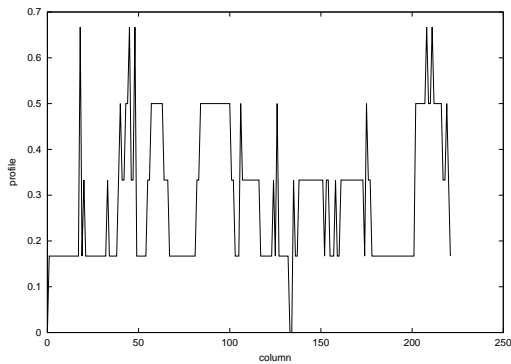
(f)



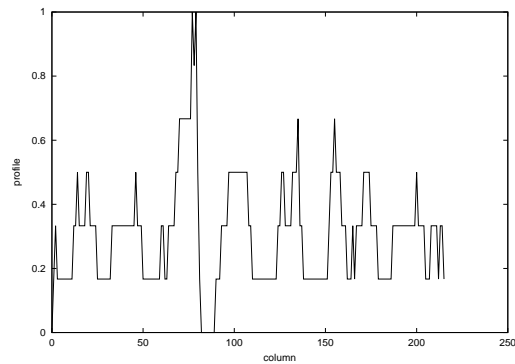
(c)



(g)



(d)



(h)

Figure 2. Profile-based shape features for two word images. Profiles in (b), (c) and (d) are extracted from the word image (a), and Profiles in (f), (g) and (h) are computed from the image (e). (b) and (f) display both upper and lower profiles, (c) and (g) contain projection profile, (d) and (h) give transition-count profile.

been	Cohen	Medical	referred
been	Cohen	Medical	referred
been	Cohen	Medical	referred

(a)

been	Cohen	Medical	referred
been	Cohen	Medical	referred
been	Cohen	Medical	referred

(b)

been	Cohen	Medical	referred
been	Cohen	Medical	referred
been	Cohen	Medical	referred

(c)

Figure 3. Samples of four characteristic words by three writers, each of whom wrote the same word three times.

4.1. Experimental Settings

The study focuses on four characteristic words, “been”, “Cohen”, “Medical”, and “referred”. Word retrieval is performed over 9,312 word images (2304 images for each characteristic word) extracted from 2328 handwriting samples by 776 individuals in U.S.. Each individual copied three times a source document specially designed by CEDAR.⁷ From each handwriting sample, four images, corresponding to the four characteristic words, are manually segmented and ground-truthed. Figure 3 shows the copies of four characteristic words provided by three writers, each of whom wrote the same word three times.

Two word image sets, a query set and an exemplar set, are constructed from the aforementioned 9,312 word images for experiments. For each writer, one of the three samples is randomly selected, then from the selected sample four word images corresponding to four characteristic words are taken as query images, the eight word images for four characteristic words in the other two samples are put into exemplar set. Thus, the query set and the exemplar set are independent. The query set consists of $776 \times 4 = 3104$ word images and the exemplar set contains $776 \times 4 \times 2 = 6208$ images. For each query image, there are $776 \times 2 = 1552$ relevant images in the exemplar set.

4.2. Experiments and Analysis

Precision-recall curve is utilized to evaluate retrieval performance. Precision is the ratio of the number of relevant images retrieved to the number of images retrieved and recall is the ratio of the number of relevant images retrieved to the total number of relevant images.

Experiments were designed to examine two scenarios: (i) performance comparison of GSC and DTW, and (ii) performance comparison of four characteristic words.

Figure 4(a)-(d) compare the retrieval performance of GSC and that of DTW for each of the four individual words, and Figure 4(e) compares the overall retrieval performance (averaged over all four words) of GSC and that of DTW. Table 1 presents the details of overall retrieval performance using GSC and DTW.

The experimental results in Figure 4 show that the proposed GSC consistently outperforms DTW. Specifically, we have following observations according to Table 1:

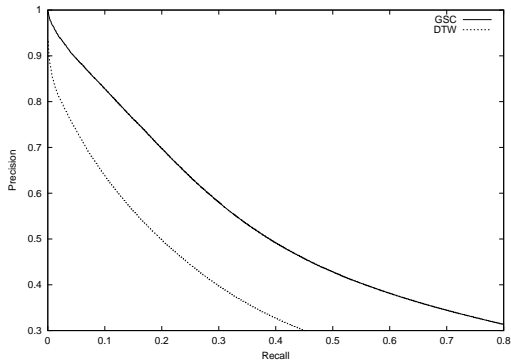
- Considering only top matches, GSC leads to a retrieval precision of 99.19%, 7.6% higher than DTW.
- With top 20 matches, the retrieval precision of GSC is 12% higher than that of DTW.
- For top 500 matches, GSC outperforms DTW by about 20% in precision and by 6% in recall.
- Considering top 1552 matches (for each query, this number is equal to the number of total images relevant in the exemplar set), GSC outperforms DTW by over 11% in both precision and recall.

No. of Top Matches	Precision		Recall	
	GSC	DTW	GSC	DTW
1	0.9919	0.9156	0.00064	0.00059
5	0.9809	0.8845	0.00316	0.00285
10	0.9742	0.8658	0.00628	0.00558
20	0.9647	0.8419	0.01243	0.01085
50	0.9455	0.7997	0.03046	0.02576
100	0.9218	0.7556	0.05939	0.04869
500	0.7858	0.5906	0.25316	0.19025
1000	0.6553	0.4967	0.42221	0.32002
1552	0.54549	0.43450	0.54549	0.43450
2000	0.4817	0.3988	0.62075	0.51394
3000	0.3892	0.3423	0.75225	0.66165
4000	0.3308	0.3024	0.85246	0.77924

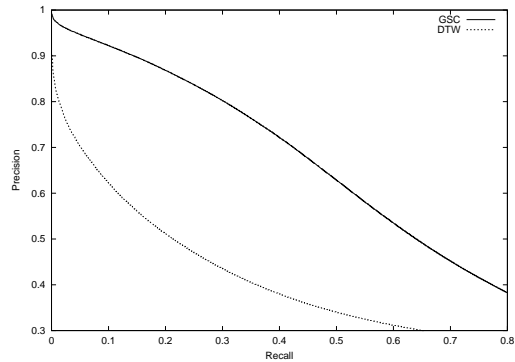
Table 1. Comparison of the average retrieval performance of GSC features and DTW features.

Figure 5 compares four characteristic words under different feature extractions. Tables 2 and 3 present the performance details of four words under GSC feature extraction and DTW one. It is easily found that both GSC and DTW prefer long words. When GSC features are used and the top 1552 matches for each query are considered, the retrieval precision of the long word “referred” is 16.6% higher than that of the short word “been”; When DTW features are used and the top 1552 matches for each query are considered, the precision for retrieving images of “referred” is about 20% higher than that for retrieving images of “been”. From Table 2 and Table 3, we may also find that the length of a word has much bigger impact on DTW than GSC. Consider an extreme case of top 4000 matches, the difference between the precisions of retrieving the word “been” images and the “referred” ones is only 7.6% when GSC is used, but the difference tops 28% when DTW is used.

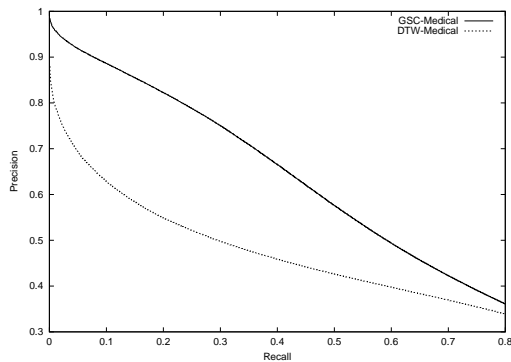
When running on an UltraSPARC workstation with an 800MHz CPU, 512MB RAM, and SunOS 5.8, DTW requires about seven days to finish 3104 queries over an exemplar set of 6208 word images, but GSC takes only 11 minutes, 893 times faster than DTW.



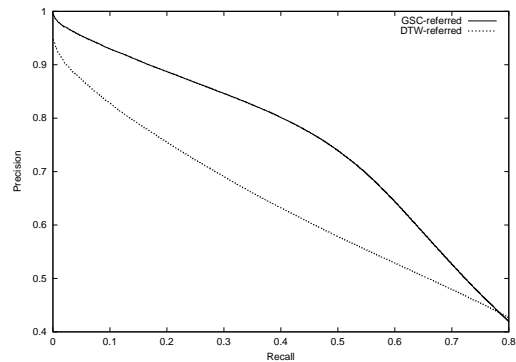
(a) word "been"



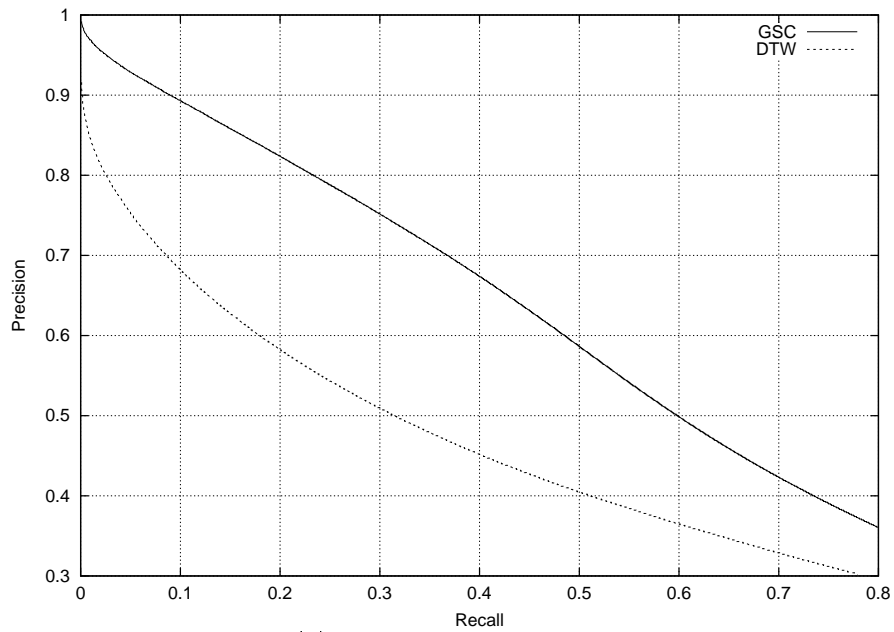
(b) word "Cohen"



(c) word "Medical"



(e) word "referred"



(e) average of four words

Figure 4. Retrieval performance comparison of GSC binary features and DTW features.

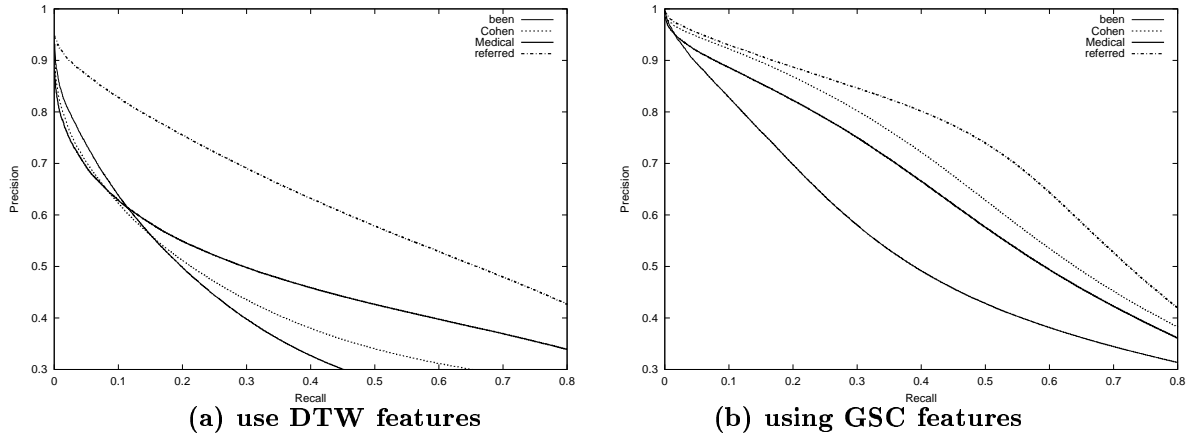


Figure 5. Retrieval performance comparison of four characteristic words.

No. of Top Matches	Precision				Recall			
	been	Cohen	Medical	referred	been	Cohen	Medical	referred
1	0.9961	0.9923	0.9845	0.9948	0.00064	0.00064	0.00063	0.00064
5	0.9814	0.9809	0.9729	0.9884	0.00316	0.00316	0.00313	0.00318
10	0.9720	0.9764	0.9651	0.9835	0.00626	0.00629	0.00622	0.00634
20	0.9569	0.9695	0.9553	0.9771	0.01233	0.01249	0.01231	0.01259
50	0.9251	0.9570	0.9353	0.9648	0.02980	0.03083	0.03013	0.03108
100	0.8847	0.9415	0.9125	0.9485	0.05700	0.06066	0.05879	0.06111
500	0.6752	0.8260	0.7857	0.8563	0.21752	0.26612	0.25312	0.27587
1000	0.5362	0.6846	0.6490	0.7513	0.34550	0.44109	0.41817	0.48406
1552	0.4545	0.5659	0.5411	0.6204	0.45450	0.56592	0.54111	0.62041
2000	0.4121	0.4982	0.4798	0.5366	0.53108	0.64205	0.61833	0.69151
3000	0.3514	0.4001	0.3894	0.4158	0.67925	0.77330	0.75263	0.80380
4000	0.3122	0.3382	0.3308	0.3418	0.80468	0.87154	0.85265	0.88099

Table 2. Comparison of the average retrieval performance of four words using GSC features.

No. of Top Matches	Precision				Recall			
	been	Cohen	Medical	referred	been	Cohen	Medical	referred
1	0.9317	0.9021	0.8802	0.9485	0.00060	0.00058	0.00057	0.00061
5	0.8907	0.8644	0.8420	0.9410	0.00287	0.00278	0.00271	0.00303
10	0.8687	0.8416	0.8225	0.9305	0.00560	0.00542	0.00530	0.00600
20	0.8423	0.8128	0.7932	0.9191	0.01085	0.01047	0.01022	0.01184
50	0.7954	0.7611	0.7477	0.8945	0.02562	0.02452	0.02409	0.02882
100	0.7418	0.7112	0.7025	0.8669	0.04780	0.04582	0.04527	0.05586
500	0.5327	0.5369	0.5612	0.7313	0.17163	0.17298	0.18082	0.23560
1000	0.4225	0.4442	0.4907	0.6292	0.27224	0.28623	0.31618	0.40543
1552	0.3553	0.3865	0.4439	0.5523	0.35530	0.38646	0.44390	0.52333
2000	0.3200	0.3552	0.4158	0.5043	0.41231	0.45770	0.53587	0.64986
3000	0.2712	0.3113	0.3668	0.4199	0.52417	0.60178	0.70894	0.81172
4000	0.2442	0.2854	0.3259	0.3539	0.62931	0.73565	0.83988	0.91213

Table 3. Comparison of the average retrieval performance of four words using DTW features.

The great advantages of GSC over DTW, in terms of both efficiency and effectiveness, arise mainly from two aspects: (i) gradient, structural, and concavity features capture multi-scale characteristics embedded in image while DTW features based on profiles, projection and transition are computed only at low resolutions; (ii) the similarity between binary vectors is calculated using extremely fast bit-based operations while DTW suffers from the very high computation complexity.

5. CONCLUSION

We present an effective and efficient approach for word image matching by using gradient-based binary features. Experiments over a large database of word images show that the proposed approach consistently outperforms the existing best word image retrieval algorithm, Dynamic Time Warping (DTW) with profile-based shape features. Not only does the proposed approach have much higher retrieval accuracy, but also it is 893 times faster than DTW.

Acknowledgments

This work was supported in part by the Department of Justice, National Institute of Justice grant 2002-LT-BX-K007.

REFERENCES

1. S. Kuo and O. E. Agazzi, "Keyword spotting in poorly printed documents using pseudo 2-d hidden markov models," *IEEE Trans. Pattern Analysis and Machine Intelligence* **16**, pp. 842–848, August 1994.
2. M. C. Burl and P. Perona, "Using hierarchical shape models to spot keywords in cursive handwriting," in *IEEE Comp. Society Conf. on Computer Vision and Pattern Recognition*, pp. 535–540, (Santa Barbara, CA), June 23-25 1998.
3. A. Kolcz, J. Alspector, M. Augusteijn, R. Carlson, and G. V. Popescu, "A line-oriented approach to word spotting in handwritten documents," *Pattern Analysis & Applications* **2**(3), pp. 153–168, 2000.
4. R. Manmatha and T. M. Rath, "Indexing of handwritten historical documents-recent progress," in *2003 Symposium on Document Image Understanding Technology (SDIUT03)*, pp. 77–85, (Greenbelt, Maryland), April 9-11 2003.
5. J. T. Favata, G. Srikantan, and S. N. Srihari, "Handprinted character/digit recognition using a multiple feature/resolution philosophy," in *Fourth International Workshop on the Frontiers of Handwriting Recognition (IWFHR IV)*, pp. 57–66, (Taipei, Taiwan), Dec 1994.
6. J. T. Favata and G. Srikantan, "A multiple feature/resolution approach to handprinted character/digit recognition," *International Journal of Imaging Systems and Technology* **7**, pp. 304–311, 1996.
7. S. N. Srihari, S.-H. Cha, H. Arora, and S. Lee, "Individuality of handwriting," *Journal of Forensic Sciences* **47**, pp. 1–17, July 2002.
8. B. Zhang, S. N. Srihari, and S. Lee, "Individuality of handwritten characters," in *Seventh International Conference on Document Analysis and Recognition*, (Edinburgh, Scotland), August 3-6 2003.
9. B. Zhang and S. N. Srihari, "Analysis of handwritten individuality using handwritten words," in *Seventh International Conference on Document Analysis and Recognition*, (Edinburgh, Scotland), August 3-6 2003.
10. J. D. Tubbs, "A note on binary template matching," *Pattern Recognition* **22**(4), pp. 359–365, 1989.
11. B. Zhang and S. N. Srihari, "Binary vector dissimilarities for handwriting identification," in *Proceedings SPIE, Document Recognition and Retrieval X*, pp. 155–166, (Santa Clara, CA), January 2003.
12. H. Sakoe and S. Chiba, "Dynamic programming optimization for spoken word recognition," *IEEE Transaction on Acoustics, Speech and Signal Processing* **26**, pp. 43–49, February 1978.