

Content-based Information Retrieval from Handwritten Documents*

Sargur Srihari, Chen Huang and Harish Srinivasan

Center of Excellence for Document Analysis and Recognition (CEDAR)

University at Buffalo, State University of New York, Buffalo, USA

Email:{*srihari, chuang5, hs32*}@cedar.buffalo.edu

Abstract

This paper is about retrieving the closest matches from a set of scanned handwritten documents based on a query that is a document image. System indexing and retrieval is based on writer characteristics, textual content as well as document meta data such as writer profile. 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. Retrieval is based on a probabilistic model of information retrieval, where feature difference between documents are modeled as probability distributions. The system has been implemented using Microsoft Visual C++ and a relational database system. This paper reports on the performance of the system for retrieving documents based on writing style.

1 Introduction

Methods for indexing and retrieval of scanned handwritten documents are needed for various applications such as historical manuscripts, scientific notes, 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 database of handwritten documents not only for textual content but also for visual content such as writer characteristics. This paper describes the indexing and retrieval aspects of a system, known as CEDAR-FOX [Srihari *et al.*, 2002; 2003], that attempts to provide the full range of functionalities for a digital library of handwritten documents.

As a document management system for handwritten documents, CEDAR-FOX provides sev-

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

eral functionalities: image indexing for content-based retrieval, keywords searching in a digital library of handwritten documents, and interactive document analysis. For the purpose of indexing based on writer characteristics, features are automatically extracted at various levels: entire document, word/phrase and character level. The user can also use interactive graphical tools to assist in obtaining more accurate features, such as image enhancement or providing a transcript of the document.

Information retrieval can be performed using several query modalities: (i) the entire document image, (ii) a region of interest(ROI) of a document, (iii) a word/phrase image, (iv) textual keyword, and (v) meta-data, e.g., document profile.

Section 2 describes the indexing aspects of CEDAR-FOX which has two parts: image features and meta data. Section 3 describes the retrieval aspects of the system. Section 4 shows the performance results of the system and concluding remarks are presented in Section 5.

2 Indexing

Indexing of handwritten document images does not only involve textual information like in the IR counterpart but also includes the image features and meta-data regarding the document. All the feature-based indexes are computed when the document is opened (or indexed in the batch mode) and are automatically stored into the database along with the meta-data.

2.1 Image Features

Image features are computed at the document, word and character levels. The document level features are called macro features whereas the word and character level features are called the micro features.

2.1.1 Micro Features

The micro features for characters and words used in CEDAR-FOX are known as GSC features [Favata and Srikantan, 1996] corresponding to gradient, structural and concavity. For character, it consists of 512 bits (G: 192 bits, S: 192 bits and C: 128 bits) features. Each of these three sets of features relies on dividing the character image into a 4×4 region. For word, the original GSC were extended to 1024 bits (G: 384 bits, S: 384 bits and C: 256 bits) features resulted from a 4×8 division scheme for a word image [Zhang *et al.*, 2004]. The gradient features capture the stroke flow orientation and its variations using the frequency of the gradient directions, as obtained by convolving the image with a Sobel edge operator, in each of 12 directions. The structural features representing the coarser shape of the character/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 binary features are converted from the original floating number computations.

To match documents based on micro-features, it is necessary to recognize characters, either manually or automatically, so that matching can be performed between the same characters. Character sizes are estimated knowing average height of text lines and word information. The estimates are used to filter connected components, so that those of appropriate size are candidate characters. In the case of cursive writing touching characters are separated using a word recognizer. The word recognizer takes the word image and its text transcription to segment word into characters before sending them to the character recognizer. Word transcription is made in one of two ways: user types the content of each word during document registration time using an interface or using automatic transcript mapping functionality. This allows a pre-typed transcript being automatically read in and the content of each word matched with the corresponding word image automatically.

2.1.2 Macro Features

In the current system, we have 13 macro features including the initial 12 features reported in [Srihari *et al.*, 2002] and average word gaps. All of these macro features are truly computed at the document level, meaning they were computed either directly from the entire document (number of black pixels, threshold, etc.) or on a line-by-line basis and applied to the entire document (number of interior contours per line, horizontal slope, etc.). In addition, recently we also develop a new set of cognitive document-level features that capture: (i) the legibility of a handwritten document at character and word level, (ii) the relative proportions existing between the characters, and (iii) the distribution of writing-styles (lexemes) existing in the handwriting of a writer. Some of these features as well as some of the original macro features are illustrated in Figure 1.

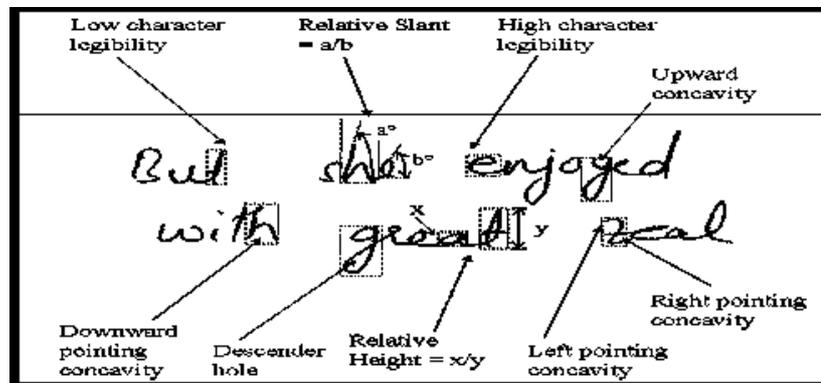


Figure 1: Macro features used in current and future version of CEDAR-FOX

2.2 Meta Data

Storing indexes for handwritten documents also involves use of data related to the origin of the document, document identification number, style of writing (cursive, handprint) etc., which makes the task of retrieval much more flexible for the user. For the use of forensic document examiners, who are potential users of such a retrieval system, we chose the following set of meta-data to describe each document: Case Number, Date of Writing, Date of Input, Writing Instrument, Geographical Area, Handedness, Keywords, Suspect's First and Last Names, Transcript, etc.

3 Retrieval

Related to the issue of retrieving documents from a database that is relevant to the query supplied is the need for associating a quantitative measure of similarity between two samples. For the task of writer identification, the goal is to take the document as a query to compare with some or all of the document data in the database. The matching between the query document and each document in the database is performed using many attributes of the documents. The matching may be done using the micro features, macro features or their combination.

3.1 Micro Feature Similarity

To measure the similarity between two characters (or words) whose shapes are represented using binary feature vectors (the GSC features), several methods have been evaluated recently [Zhang and Srihari, 2003]. This has led to the choice of the Correlation measure defined in [Tubbs, 1989]. 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 the two feature vectors X and Y is given by the formula:

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

The distributions of distances in both the same-writer and different-writer categories follow a univariate Gaussian density. The parameters of these two distributions are estimated by using the training data. The parameters include: the mean μ_{SW} and variance σ_{SW} for same-writer category and the mean μ_{DW} and variance σ_{DW} for different-writer. They define the two distribution densities $p_{SW}(x)$ and $p_{DW}(x)$. The training data set consists of 2000 pairs of sample documents for each of the two categories. The details of the data set will be described in Section 4.1. The parameters for all the 62 characters (0-9, a-z, A-Z) were computed. Table 1 gives the means and variances for numeral characters.

During matching, for each character $c_i, i = 1, \dots, N$, where N is the number of characters considered, we compute the distance d_i^j between the two samples of pair j for that character. We can

Character	Same Writer		Different Writer	
	Means	Variances	Means	Variances
0	0.653	0.087	0.747	0.069
1	0.674	0.143	0.799	0.096
2	0.585	0.077	0.709	0.073
3	0.612	0.069	0.708	0.065
4	0.602	0.085	0.735	0.072
5	0.592	0.076	0.709	0.066
6	0.641	0.084	0.761	0.066
7	0.641	0.096	0.766	0.077
8	0.577	0.082	0.698	0.074
9	0.624	0.082	0.745	0.068

Table 1: Parameters of the distance distribution for micro features of numeral characters

have $j = 1, \dots, M_i$ possible pairs of samples for a given character. For characters $c_i, i = 1, \dots, N$ and M_i pairs of samples for each character we estimate the log-likelihood ratio:

$$LLR(micro) = \ln \left(\frac{\prod_{i,j} p_{SW}(d_i^j)}{\prod_{i,j} p_{DW}(d_i^j)} \right) \quad (2)$$

If $LLR(micro) > 0$, we have a same-writer decision, if $LLR(micro) < 0$ we have a different-writer decision.

3.2 Macro Feature Similarity

The macro features are mapped into a distance vector of differences. Each distance element consist of the absolute difference in value between two documents for that macro feature. The distance distributions for both the same-writer and different-writer categories using the macro features are approximated as univariate Gaussian distribution. Table 2 gives the means and variances of distance distribution for the macro features. The variances indicate that non-zero probabilities will be assigned to negative values of distances. Although a Gamma distribution would be more appropriate we have found reasonably good results using Gaussian.

The verification decision for any pair of documents using macro features is made also using equation 2 above, resulting in $LLR(macro)$ value where the distance are measured using absolute differences for each of the real-valued macro features. The degree of match between two documents is measured using the total LLR score as follows: $LLR(micro) + LLR(macro)$.

This approach to similarity measurement is similar to the probabilistic model of information retrieval [Jones, 1998]. In this model, the LLR gives a relevance/similarity measurement for document retrieval.

Macro Features	Same Writer		Different Writer	
	Means	Variances	Means	Variances
Gray Scale Entropy	0.022	0.022	0.018	0.018
Gray Binary Threshold	0.572	0.623	2.612	1.708
No. of Black Pixels	1915.52	1417.73	4566.41	2973.19
No. of Exterior Contours	1.490	1.270	6.358	5.040
No. of Interior contours	3.009	2.488	7.094	4.368
Horizontal Slope	0.017	0.011	0.064	0.040
Positive Slope	0.014	0.010	0.100	0.061
Vertical Slope	0.018	0.013	0.089	0.055
Negative Slope	0.008	0.006	0.053	0.037
Average Stroke width	0.247	0.374	0.822	0.644
Average Slant	1.604	1.166	11.679	7.171
Average Height	2.221	1.580	7.640	4.892
Average Word Gap	0.038	0.017	0.105	0.052

Table 2: Parameters of the distance distribution for macro features

3.3 Query Methods

3.3.1 Meta Data

With an easy data entry interface, all the meta-data mentioned in section 2.2 can be added to the database along with the image features and will be stored as a record corresponding to the document. Thus a document in terms of a database entity can be considered as a single record or a ‘tuple’.

Our system provides efficient retrieval of such a database. Several query modalities are permitted for retrieval. The database functionality has been implemented using MySQL, a database management system by MySQL ABTM. The interface to the database is through the MySQL libraries. From the user point of view, a graphical user interface has been provided which takes as input the fields on which the user wants to query the database. This query is mainly based on textual information and meta-data. The user can query the database in order to retrieve documents relevant to the meta-data and textual information he uses in the query.

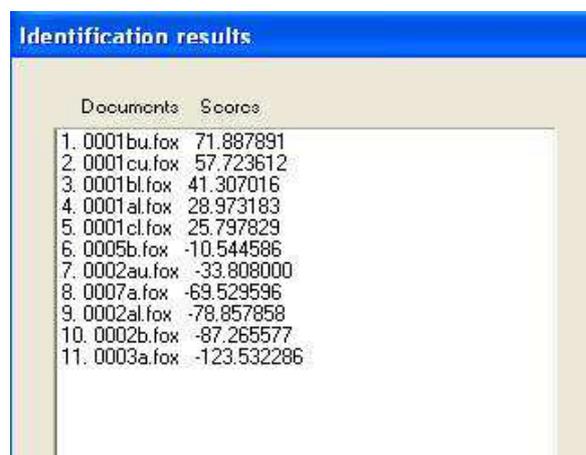
3.3.2 Document Features

Another type of querying is based on the features of the documents and this is where identification comes to play. The process of identification is nothing but querying in which the query consists of a document. Unlike the IR counterpart where the query has to be considered as a pseudo-document for which similar relevant documents are retrieved, here the query is a real document and the task of this information retrieval system is to use the similarity measurements to pull-out documents that it thinks is relevant. In effect, when we give a document as a query, we expect the system to filter out only those documents that it thinks is written by the

same writer who wrote the query document. The relevance or similarity of the retrieved document is measured using the similarity metrics presented above in equation 2. The LLR (Log Likelihood Ratio) for the query document is computed against each document in the database and it is used to decide similarity, just like the probabilistic model of information retrieval does. When the query involves a document rather than textual information entered by the user, there are three possible options the user can use to define what part of the document he wants to use to query the database.

- **Document Level: the entire document image is the query.** When the system loads in a document image, it can be directly used as query. For identifying the document from the database, the automatically extracted features are used for the matching. The query returns a ranked list of documents in the library.
- **Partial Image: a region of interest (ROI) of a document.** A document image may include many text or graphical objects. User often needs to specify a local region of the most interest. Using a system cropping tool, user can easily crop a rectangular region and use it as a query.
- **Word Level Image: any word in the document.** Any word from the query document can be cropped and compared against other words from documents in the database. The result will be a ranking of the words that are most similar with the query. The document containing the retrieved words can then easily be obtained from the database.

Figure 2 shows a sample screen shot of the writer identification using entire document as the query.



The screenshot shows a window titled "Identification results" with a table of document scores. The table has two columns: "Documents" and "Scores". The data is as follows:

Documents	Scores
1. 0001bu.fox	71.887891
2. 0001cu.fox	57.723612
3. 0001bl.fox	41.307016
4. 0001al.fox	28.973183
5. 0001cl.fox	25.797829
6. 0005b.fox	-10.544586
7. 0002au.fox	-33.808000
8. 0007a.fox	-69.529596
9. 0002al.fox	-78.857858
10. 0002b.fox	-87.265577
11. 0003a.fox	-123.532286

Figure 2: Screen shot of writer identification results

From
 Jim Elder
 809 Loop Street 300
 Abertown, New York 14707

To
 Dr. Bob Grant
 600 Summersby Parkway
 Oak, West Virginia 25638

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zuck.

It all started ~~about~~ around six months ago while attending the "Ruben" Jazz Organizing such an event is no picnic, and as President of the Alumni as co-sponsor of the event, Kate was overworked. But she enjoyed her job, and was required of her with great zeal and enthusiasm.

However, the extra hours affected her health; halfway through the show we we noted her to the hospital, and several questions, x-rays and blood tests later it was just exhaustion.

Kate's been in very bad health since. Could you kindly take a look at the and give us your opinion?

Thank you!

JIM

Nov 1987

From
 Jim Elder
 809 Loop Street
 Abertown, New York 14707

To
 Dr. Bob Grant
 600 Summersby Parkway

We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zuck.

It all started around six months ago while attending the "Ruben" Jazz Organizing such an event is no picnic and as President of the Alumni as co-sponsor of the event, Kate was overworked. But she enjoyed her job, and was required of her with great zeal and enthusiasm.

However, the extra hours affected her health; halfway through the show we we noted her to the hospital, and several questions, x-rays and blood tests later, we we noted her to the hospital. Several questions, x-rays and blood tests later, we we noted her to the hospital. Several questions, x-rays and blood tests later, we we noted her to the hospital.

Kate's been in very bad health since. Could you kindly take a look at the results and give us your opinion?

Thank you!

Figure 3: Two samples of the CEDAR letter

4 Results and Analysis

4.1 Experimental Framework

We have a document image collection consisting of 3000 samples written by 1000 writers, where each writer wrote three samples of a pre-formatted letter called CEDAR letter. (This means all the samples have the same content of CEDAR letter.) All the documents were scanned using 300 dpi resolution. Portions of two CEDAR letters from two individuals are shown in Figure 3.

Since in most cases of document retrieval the documents to be compared will have different content, for the different content experiments in this study, for each document, the y-coordinate of an imaginary center line was computed and it was used to cut a whole image into two roughly, logically equal size images, the upper and lower parts. This was done for all the 3000 documents, which led to 6000 half documents.

To evaluate parameters for the distribution density of same-writer category and different-writer category, the training set consists of 2000 pairs of half documents for each of the two categories. And all the pair documents have different content.

In the testing, 412 cases of identification were tested. For each case, one document was selected as the query document and a database of 412 different content documents was formed to be queried. Among these 412 documents, there exists only one relevant document. For example, if the document '0001U' (means the upper part of document 0001) is the query, then the database

will contain document '0001L'(means the lower part of document 0001) and 411 lower part documents randomly selected from the other writers.

The macro features used in this study consist of the following eight features: (i) number of exterior contours per line, (ii) number of interior contours per line, (iii) horizontal slope (0 degree or flat stroke), (iv) positive slope (45 or 225 degree), (v) vertical slope (90 or 270 degree), (vi) negative slope (135 or 315 degree), (vii) average height per line and (viii) average slant per line.

The characters that are used to generate the micro features are recognized by using the word recognizer method mentioned in Section 2.1.1. The word transcription is provided by using the automatic transcript mapping functionality.

4.2 Experimental Results

The document management system CEDAR-FOX on which this experiment is conducted is able to perform document retrieval at the document image, document ROI or word image level. Document image retrieval, for which the experimental results presented here have been obtained, provides the user with an automatic and efficient way to identify a questioned document from a large set of known documents in the database.

Figure 4 shows the results of our identification experiments. The x-axis represents the number of top ranked documents retrieved, and the y-axis represents the accuracy. As seen from the figure, the accuracy of the correct writer being identified (from a pool of 412 writers) in the top 5 ranked documents is 89.3% and if top 10 ranked documents are retrieved the accuracy stands at 91.9%.

5 Conclusion

We have described the information retrieval aspects of a document analysis and management system for handwritten documents. Writer characteristics as well as document content and meta data can be used for retrieval. The performance of the system using a large set of document and character/word level features has been presented.

6 Acknowledgments

CEDAR-FOX is the result of the effort of many students and researchers at CEDAR. We thank the following students who have contributed to the system implementation and its testing: Vinay Shah, Vivek Shah, Bandi Kartik Reddy, Siyuan Chen, Meenakshi Kalera and Devika Kshirsagar.

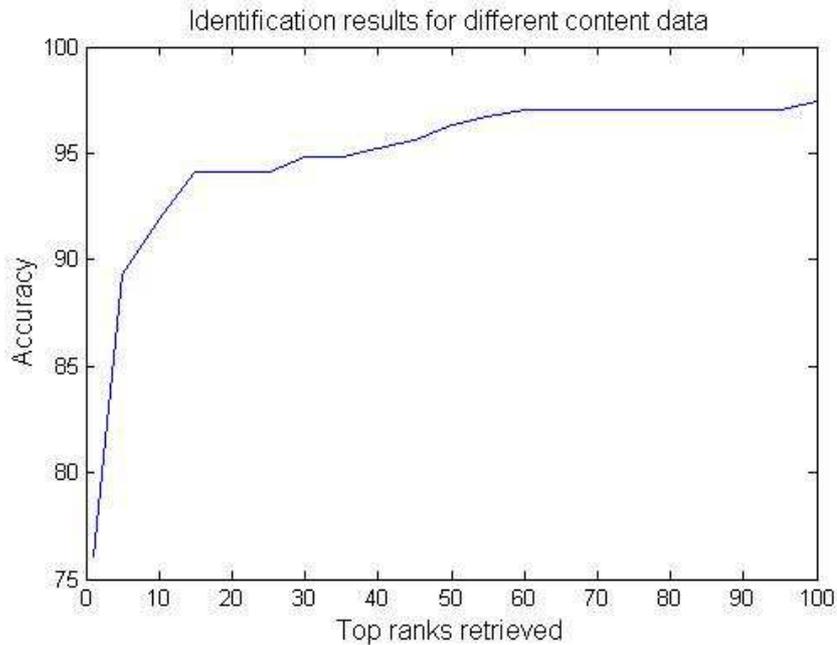


Figure 4: Identification results for different content documents (412 tests conducted)

References

- [Favata and Srikantan, 1996] John T. Favata and Geetha Srikantan. A multiple feature resolution approach for handprinted digit and character recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [Jones, 1998] Karen Sparck Jones. A probabilistic model of information retrieval: Development and status. Technical report, Computer Laboratory, University of Cambridge, UK, 1998.
- [Srihari *et al.*, 2002] Sargur N. Srihari, Sung-Hyuk Cha, Hina Arora, and Sangjik Lee. Individuality of handwriting. In *Journal of Forensic Sciences*, volume 47(4), pages 856–872, 2002.
- [Srihari *et al.*, 2003] Sargur N. Srihari, Bin Zhang, Catalin Tomai, Sangjik Lee, Zhixin Shi, and Yong-Chul Shin. A system for handwriting matching and recognition. In *Proceedings of the Symposium on Document Image Understanding Technology (SDIUT 03)*, Greenbelt, MD, 2003.
- [Tubbs, 1989] J. D. Tubbs. A note on binary template matching. In *Pattern Recognition*, volume 22(4), pages 359–366, 1989.
- [Zhang and Srihari, 2003] Bin Zhang and Sargur N. Srihari. Binary vector dissimilarity measures for handwriting identification. In *Document Recognition and Retrieval X, SPIE, Bellingham, WA*, volume 5010, pages 28–38, 2003.
- [Zhang *et al.*, 2004] Bin Zhang, Sargur N. Srihari, and Chen Huang. Word image retrieval using binary features. In *Document Recognition and Retrieval XI, SPIE, San Jose, CA*, 2004.