

Automatic Handwriting recognition and Writer Matching on Anthrax-related Handwritten Mail

Sargur N. Srihari and Sangjik Lee
Center of Excellence for Document Analysis and Recognition (CEDAR)
University at Buffalo, State University of New York (SUNY)
520 Lee Entrance, Suite 202
Amherst, NY 14228, USA
{srihari,sj11}@cedar.buffalo.edu

Abstract

A handwriting recognition technique and a handwriting identification technique were combined and applied on anthrax-related handwritten mail. HWAI (Handwritten Address Interpretation) system interprets the address on the anthrax-related letters. The HWAI process can be modified to provide specific alphabet images which can then be used for writer identification. Micro-level feature values from segmented characters were extracted and an identification test was conducted. Preliminary results based on automatic handwriting recognition and identification are shown.

1. Introduction

This paper reports on an analysis of the writings in recent anthrax-related mailings using combination of HWAI (Handwritten Address Interpretation) system and Forensic Document Examination system.

Scientists at CEDAR have been conducting research on computer interpretation of postal images for nearly two decades. Also, research on software for forensic analysis of handwriting has been conducted for over two years [1]. The recent threat of anthrax mailings have called for combining techniques from the two areas for forensic purposes. The volume of handwritten mail processed by the U. S. Postal Service is very large (over 50 million per day).

Since combination of handwriting recognition and writer identification is highly time-consuming this approach will form a good starting point on automatic forensic document examination system for postal envelopes.

This paper consists of five phases: introduction of anthrax-related images (section 2), experimental setup that used in this application (section 3), handwriting recognition of anthrax-related images using HWAI system (section 4),



Figure 1. Four writings which illustrate the study: the extent to which they were written by the same individual.

handwriting identification on anthrax-related images (section 5), and finally summary and conclusion (section 6)

2. Anthrax-related Images

Images of four envelopes associated with anthrax mailings are shown in Fig. 1. The handwritten letters contained in three of those envelopes are also shown. Non-anthrax control letters are shown in Fig. 2

3. Experimental Setup

A flowchart of handwriting recognition and identification through envelope image to character-level feature values is shown in Fig. 3. Recognition part provides images of alphabets and numeric for the identification part.

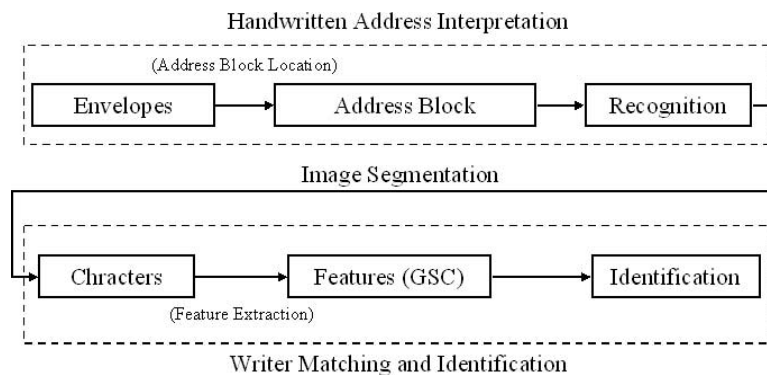


Figure 3. the whole system used in handwriting recognition and identification on this experiment

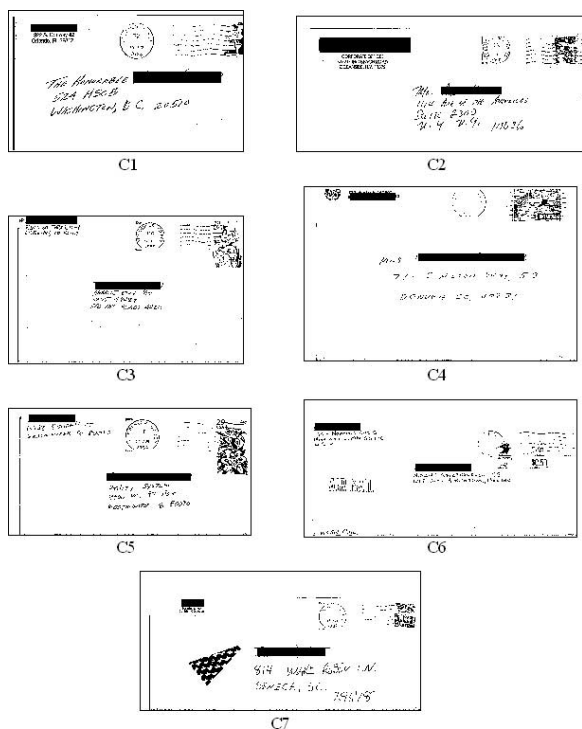


Figure 2. Control Letters C1-C7 (The names of recipients are blanked out for privacy)

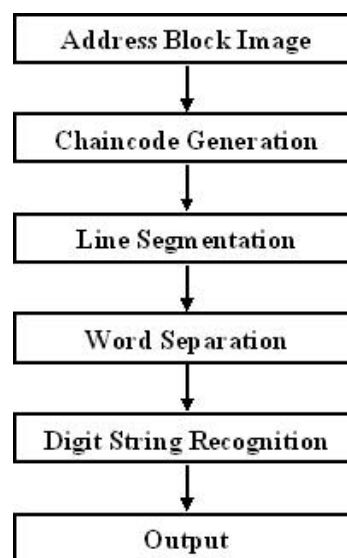


Figure 4. Main flowchart of HWAI system

3.1. HWAI (Handwritten Address Interpretation) system

The Handwritten Address Interpretation [2], or HWAI, project that has been conducted at CEDAR since the early eighties explores the use of knowledge of the postal domain in the recognition of handwritten addresses. The task is considered to be one of interpretation rather than recognition since the goal is to assign the address to its correct destination irrespective of incomplete or contradictory information present in the writing. This work at CEDAR has led to a system that recognizes handwritten addresses that is currently in use by the United States Postal Service. Main flowchart of HWAI system is shown in Fig. 4.

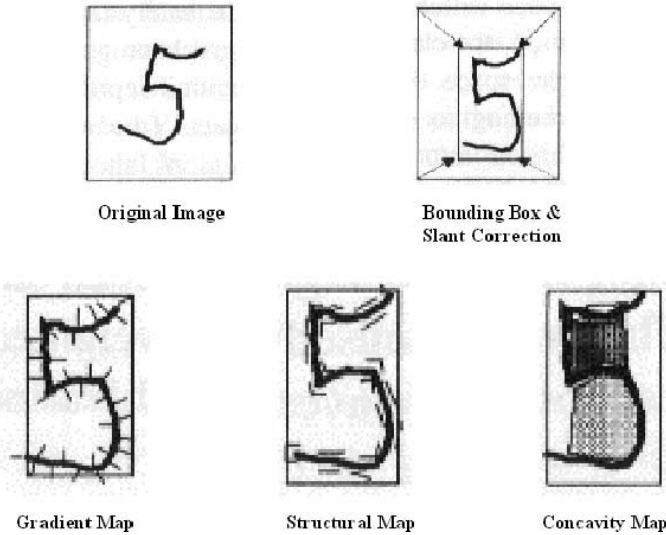


Figure 5. General flow of GSC feature generation

3.2. GSC Feature Extraction

For micro-level features for characters, we used GSC feature extractor. The GSC algorithm [3, 4] takes a quasi-multiresolution approach to feature generation; that is, several distinct feature types are applied at different scales in the image. These computed features measure the image characteristics at local, intermediate, and large scales.

The gradient features detect local features of the image and provide a great deal of information about stroke shape on a small scale. The structural features extend the gradient features to longer distances and give useful information about stroke trajectories. The concavity features are used to detect stroke relationships at long distances which can span across the image. General flow of GSC feature generation is shown in Fig. 5.

3.3. Handwriting Identification

Writer identification is the task of determining the writer when there are n candidates. This classification task has to be learnt from a set of handwriting samples provided by each of the candidates. Given a test sample of an unknown writer, the task is to determine whether it was written by any of the n writers and, if so, to identify the writer. The writer identification procedure uses the features extracted from the test image and from the labeled prototype images to determine the writer of the test document.

Table 1. Results of CEDAR-developed HWAI on four envelope images

Envelope	Zip Code	Read as	Confidence
Brokaw (10112)	1	2	0.83
	0	0	1.00
	1	2	0.82
	1	1	0.50
NY Post (10036)	2	2	1.00
	1	2	0.95
	0	0	1.00
	0	0	1.00
Daschle (20510)	3	3	1.00
	6	6	0.95
	2	2	1.00
	0	0	1.00
	5	5	0.95
	1	(missing)	
	0	(missing)	
	2	2	1.00
Leahy (20510-4502)	0	0	1.00
	5	5	1.00
	1	2	0.94
	0	0	1.00
	4	4	0.94
	5	5	1.00
	0	0	1.00
	2	2	1.00

4. Handwriting Recognition of Anthrax-related Envelopes using HWAI system

The results of applying the CEDAR-developed Handwritten Address Interpretation system to the four anthrax envelopes are shown in Table 1.

The results of digit recognition confidence are in the table. As can be seen, the 1's are recognized as 2's due to the presence of the serif and initial spur. The result of HWAI on the Leahy letter is similar to what has been reported in the general press: that it was incorrectly read as going to the US State Department instead of the Senate Russell Building.

As a historical note, research at CEDAR led to the first deployment of a handwritten address interpretation (HWAI) system by Lockheed Martin at Christmas 1997. Today, the CEDAR HWAI system alone can correctly interpret over 50% of handwritten mail with a less than 2% error rate. In combination with other technology, Lockheed is realizing even higher rates of performance in mail sorting systems.

5. Handwriting Identification on Anthrax-related Envelopes

We also have developed software (known as CEDAR-FDE, for Forensic Document Examination) for matching handwriting for the purpose of determining writer-ship. This software was developed as part a study to quantitatively establish the individuality of handwriting, a project supported by the National Institute of Justice (NIJ). We report here on a preliminary matching of the handwriting styles on the five envelopes and letters using CEDAR-FDE software for determining handwriting similarity.

We have developed a set of 11 macro-features that can characterize a handwritten document at any level: entire document, paragraph, word, or character. We also have developed 512 micro-features that are computed from a connected component (such as a word or character extracted from the writing). As a preliminary test, we applied the matching software using only micro-features at the character level to a few characters found in the images.

5.1. Micro-features

We applied the matching software using only micro-features at the character level to a few characters found in the images. The frequency of alphabets in the four envelopes is shown in Table 2.

The micro-level features consist of 512 binary (0 or 1 value) features corresponding to gradient (192 bits), structural (192 bits), and concavity (128 bits) features. Micro-features of a character "A" from the Brokaw envelope are shown in Fig. 6.

5.2. Similarity of Writing Based on Micro-features

From each of the five envelope images, eight character images were extracted (A, C, E, N, O, R, T, W) for further analysis. For each image, micro-level features were computed, resulting in 4096 bits (512 * 8 characters). The similarity was between a pair of images is calculated as follows: the weighted sum of identical bits where 1.0 is added if both bits are 1, and 0.5 is added if both bits are 0. Then the weighted sum is divided by (number of 0's*0.5) + (number of 1's*1.0) of the query document. The similarity between A1 and A3 is 87.8%, between A1 and A4 is 85.3%, and between A1 and A2 is 88.5%. The lowest match of A1 was to C2 (71%).

Considering only eight alphabets, we can see that each of the four anthrax envelope writings are more similar to each other than to the control letter. The use of macro-features can provide a better match than those shown in Table 3.

When there is a large volume of mail images, such a system can be used to narrow down the search to a smaller set

of similar cases to be presented to a forensic document examiner. The individual characters to be compared can be automatically selected, based on recognition results.

6. Summary and Conclusion

In this paper we used a combination of handwriting recognition techniques and handwriting identification techniques sequentially. Application of the micro-features to the recent anthrax-related mailings shows their usefulness in handwriting identification. This approach will be especially useful where a large number of handwriting must be identified automatically.

Since this approach can save a lot of time and effort, forensic document examiners can have a benefit where they can concentrate what they really need instead of time-consuming repeated work. We will work further on more discriminatory features and look forward to more accuracy.

Acknowledgements

The authors wish to thank Vemulapati Ramanaprasad who assisted in obtaining results of the CEDAR HWAI system on the anthrax images.

References

- [1] Sargur N. Srihari, Sung-Hyuk Cha, Hina Arora, and Sangjik Lee. Individuality of handwriting. *Journal of Forensic Sciences*, in press.
- [2] Sargur N. Srihari, E. Cohen, V.Govindaraju, and A. Shekhawat. Determining delivery point codes on handwritten addresses. In *Fifth United States Postal Service Advanced Technology Conference*, pages 321–335, November-December 1992.
- [3] John T. Favata and Geetha Srikantan. A multiple feature/resolution approach to handprinted digit and character recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [4] G. Srikantan, S. W. Lam, and S. N. Srihari. Gradient-based contour encoding for character recognition. *Pattern Recognition*, 29(7):1147–1160, 1996.

Table 2. Alphabets and numerics available from four anthrax-related envelopes.

Recipient	Document	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	0	1	2	3	4	5	6	7	8	9	Sum
Brokaw	envelope	3	2	2		4	1					3	3	1	3	4	1		4		2		1	2		2	1	2	3	1	1							46
NY Post	envelope	3		1	1	6	1		1	2		2		1	3	5	1		4	2	3		1	2		3		2	4	1	1			1				51
Daschle	envelope	9	1	4	5	9	3	4	5	5	1	2	5		9	5	1		6	5	5	1		1			3	2		1	3				2	1	98	
Leahy	envelope	8	1	3	4	10	3	4	4	5	1	2	7		9	5	1		6	6	4	2		1	1		4	1	3	2	3	3				2	105	
	Sum	23	4	10	10	29	8	8	10	12	2	9	15	2	24	19	4	0	20	13	14	3	2	6	0	6	1	11	8	7	4	4	6	1	0	4	1	300



```

1010000000000010000000000010000000000010001000001000000000000000
0000000000000000000000000000000000000000000000000000000000000000
0000000000010000001000000010000000000000000000000000000000000000
Gradient (192 bits)
0000000001000000010000000000100000000000000000000000000000000000
0000000000000000000000000000000000000000000000000000000000000000
0000000000000000000000000000000000000000000000000000000000000000
Structural (192 bits)
1111111111111111100000000000000000000000000000000000000000000000
0000000000000000000000000000000000000000000000000000000000000000
Concavity (128 bits)

```

Figure 6. Micro-features of the character A from envelope A1 in Fig. 1

Table 3. Similarity between writings A1-A4 and C1-C7 based on eight characters.

Query	A1	A2	A3	A4	C1	C2	C3	C4	C5	C6	C7
A1		0.89	0.88	0.85	0.75	0.71	0.73	0.75	0.75	0.75	0.75
A2	0.89		0.88	0.87	0.76	0.73	0.75	0.75	0.74	0.76	0.75
A3	0.88	0.87		0.89	0.75	0.72	0.72	0.75	0.74	0.75	0.74
A4	0.85	0.86	0.88		0.73	0.70	0.70	0.73	0.73	0.75	0.73
C1	0.76	0.76	0.76	0.74		0.78	0.78	0.77	0.78	0.71	0.77
C2	0.70	0.71	0.71	0.69	0.76		0.77	0.73	0.74	0.71	0.72
C3	0.72	0.73	0.70	0.69	0.76	0.77		0.76	0.73	0.70	0.74
C4	0.74	0.74	0.75	0.72	0.76	0.74	0.78		0.76	0.72	0.75
C5	0.75	0.74	0.74	0.73	0.77	0.75	0.75	0.77		0.73	0.76
C6	0.74	0.75	0.75	0.74	0.70	0.72	0.71	0.72	0.72		0.71
C7	0.73	0.72	0.72	0.70	0.74	0.71	0.73	0.73	0.74	0.69	