

# Assessing the Authorship Confidence of Handwritten Items

Sung-Hyuk Cha and Sargur N. Srihari  
Center of Excellence for Document Analysis and Recognition  
State University of New York at Buffalo  
{scha,srihari}@cedar.buffalo.edu

## Abstract

*The Writer Verification is a process to compare questioned handwriting with samples of handwriting obtained from known sources for the purposes of determining authorship or non-authorship. It plays an important investigative and forensic role in many types of crime. Particularly problematic in this regard is the fact that there exists neither true standard nor universal definition of comparison. For this reason, we propose an algorithmic objective approach. Using visual information of two or more digitally scanned handwritten items, we show a method to access the authorship confidence that is the probability of errors. Instead of building a costly handwritten item database to support the confidence, questioned words are simulated from the CEDAR letter image database in order to handle any handwritten items. An Artificial Neural Network is trained to verify the authorship using the synthesized words.*

**Key Words:** *Authorship, CEDAR letter, Dichotomizer, Questioned Document, Writer Verification*

## 1. Introduction

The *Writer Verification* is a process to compare questioned handwriting with samples of handwriting obtained from known sources for the purposes of determining authorship or non-authorship. It is the examination of the design, shape and structure of handwriting to determine authorship of given handwriting samples. Document examiners or handwriting analysis practitioners find important features to characterize individual handwriting as features are consistent with writers in normal undisguised handwriting [1]. Authorship may be determined due to the following hypothesis that people's handwritings are as distinctly different from one another as their individual natures, as their own finger prints. It is believed that no two people write the exact same thing the exact same way.

Not surprisingly, the *Writer Verification problem* has re-

ceived a great deal of attention since it plays an important investigative and forensic role in many types of crime. Hence, there is a wealth of literature regarding *document examination* dating as far back as 1929. Osborn suggested a statistical basis to handwriting examination by the application of the *Newcomb rule of probability* [8] and Bertillon was the first who used the Bayesian theorem to handwriting examination [6]. Hilton calculated the *odds* by taking the likelihood ratio statistic that is the ratio of the probability calculated on the basis of the similarities, under the assumption of identity, to the probability calculated on the basis of dissimilarities, under the assumption of non-identity [6, 5]. Various automatic writer identification by computer techniques, feature extraction, comparison and performance evaluation methods have been studied (see [9] for the extensive survey). However, relatively little study has been carried out to demonstrate its scientific and statistical validity and reliability as forensic evidence. To verify the authorship, it is necessary to determine the statistical validity of individuality in handwriting based on measurement of features, quantification, and statistical analysis.

Particularly problematic in this regard is the fact that there exists neither true standard nor universal definition of similarity. Interaction with questioned document is generally personal, subjective and limited to some degree by geography. Given the nature of the assessment, the improvement in it might be realized if the process were less subjective. The FBI formed the Technical Working Group on Forensic Document Examination (TWGDOC) in May 1997 and the importance of standardizing procedures for handwriting comparison was recognized as a primary task [11]. Such procedures must be based on more than community-based agreement. Procedures must be tested statistically in order to demonstrate that following the stated procedures allows analysts to produce correct results with acceptable error rates. This has not yet been done. For this reason, we propose an algorithmic objective approach for the analysis part of the procedure.

In our earlier work on the validation of individuality of handwriting [3], we asked 1,000 individuals to write the

“CEDAR letter” three times to determine the statistical validity of individuality in handwriting based on measurement of features, quantification and statistical analysis. To establish the inherent distinctness of the classes, i.e., validate individuality, we transform the many class problem into a dichotomy by using a “distance” between two samples of the same class and those of two different classes. A measure of confidence is associated with individuality. Using ten feature distance values, we trained an *artificial neural network* and obtained 97% overall correctness. Various features and their distance measures are given in [2]. Unfortunately, the features used in this experiment do not always occur in the real world writer identification problem. Features are subject to change in case by case.

Consider a comparison between a handwritten item from a ransom note and suspects’ exemplars of the same handwritten item. In an attempt to find the authorship confidence of the questioned documents, one can ask 1,000 or a substantially large number of people that is representative to US population or criminals and do the same analysis as in [3]. Instead of *CEDAR letter database*, the given handwritten item database is constructed to support the authorship confidence. This experiment gives *type I* and *type II* errors helping the decision of the authorship. This is a very costly and time consuming process.

To this end, we propose a simulated method to access the authorship confidence of handwritten items. Instead of constructing the costly handwritten item database, simulated handwritten item database is built from the *CEDAR letter database*. We provide a procedure for determining whether two handwritten items were written by the same person. The suggested procedure is for digitally scanned handwritten items using visual information only.

The subsequent sections are organized as follows. The section 2 discusses the completeness of the *CEDAR letter* for the purpose of synthesization followed by the database specification consisting of writer data and features obtained from a handwriting sample, statistically representative of the US population. In section 3, we outline the procedure and then explain each step with the example of the real ransom note from *JonBenet Ramsey* murder case. Finally, section 4 concludes the work and addresses open problems.

## 2. CEDAR Letter Database

In this section, we introduce the *CEDAR letter*, discuss its completeness for the purpose of synthesization and then describe the specification of the *CEDAR letter image database* consisting of writer data and features obtained from a handwriting sample, statistically representative of the US population.

## 2.1. CEDAR Letter and its Completeness

<p>From Jim Elder 829 Loop Street, Apt 300 Allentown, New York 14707 To Dr. Bob Grant 602 Queensberry Parkway Omar, West Virginia 25638 We were referred to you by Xena Cohen at the University Medical Center. This is regarding my friend, Kate Zack. It all started around six months ago while attending the “Rubeq” Jazz Concert. Organizing such an event is no picnic, and as President of the Alumni Association, a co-sponsor of the event, Kate was overworked. But she enjoyed her job, and did what was required of her with great zeal and enthusiasm. However, the extra hours affected her health; halfway through the show she passed out. We rushed her to the hospital, and several questions, x-rays and blood tests later, were told it was just exhaustion. Kate’s been in very bad health since. Could you kindly take a look at the results and give us your opinion? Thank you! Jim</p>	<p>Nov 10, 1999</p>
---	---------------------

Figure 1. CEDAR Letter

The *CEDAR Letter*, as shown in Figure 1, is concise (it has just 156 words), easy to understand and *complete*. It’s *complete* in that, each alphabet occurs in the beginning of a word as a capital and a small letter, and as a small letter in the middle and end of a word. In addition, it also contains punctuation, numerals, interesting letter and numeral combinations (ff, tt, oo, 00) and a general document structure that would allow us to extract document level features such as word and line spacing, line skew etc. Forensic literature refers to many such documents - the “*London Letter*” [8], the “*Dear Sam Letter*” to name a few. But none of them are *complete* in the sense of the CEDAR Letter as follows. All capitals must appear in the letter and it is desirable to have all small letters in the beginning, middle and terminal positions of the word. We score the letter according to these constraints:

$$\text{score}(\text{letter } x) = \frac{104 - \text{Number of missings}}{104} \quad (1)$$

The CEDAR letter scores 99% whereas the London letter scores 76%. the cedar letter has only 1 missing entry that is a word that ends with a letter “j”. Since there is no common English word that ends with the letter “j”, the cedar letter excludes this entry.

The table in Figure 2 shows counts of each letter in each position. The second row is the appearance count of the capital letter A through Z in the beginning of the word. The fourth to sixth rows are the appearance counts of the small

letter *a* through *z* in the beginning, middle and terminal positions, respectively. Figure 2 shows the comparison be-

	A	B	C	D	E	F	G	H	I	J	K	L	M
Init	1	2	2	3	3	1	1	1	1	1	1	4	2
	a	b	c	d	e	f	g	h	i	j	k	l	m
Init	13	4	0	0	1	2	4	2	1	1	0	1	0
Mid	7	1	2	5	26	1	1	7	13	0	0	10	3
Ter	2	0	2	15	13	0	1	0	0	0	1	2	0

(a) London Letter alphabet frequency counts: score = 76%.

	A	B	C	D	E	F	G	H	I	J	K	L	M
Init	4	2	4	1	1	1	1	1	1	2	3	1	1
	a	b	c	d	e	f	g	h	i	j	k	l	m
Init	17	4	1	1	6	1	2	9	4	2	1	2	2
Mid	33	2	8	6	59	4	5	20	32	1	3	14	3
Ter	5	2	1	21	20	3	3	5	1	0	3	5	2

(b) CEDAR Letter alphabet frequency counts: score = 99%.

	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
Init	1	1	1	1	2	2	2	1	1	1	1	1	1
	n	o	p	q	r	s	t	u	v	w	x	y	z
Init	1	6	2	1	5	8	14	1	1	8	1	3	1
Mid	35	36	4	1	30	19	25	18	7	5	2	2	2
Ter	7	5	1	1	12	15	17	2	1	2	1	8	1

Figure 2. The CEDAR and London Letter: A Comparison.

tween the CEDAR and London letters. Each table has 104 entries and we would like each entries to be non-zero. The London letter has 25 zero entries whereas the cedar letter has only 1 zero entry that is a word that ends with a letter “*j*”. The completeness of the *CEDAR letter* allows to synthesize any English word.

## 2.2. Specification of Database

Each subject provides three exemplars of the CEDAR Letter and they are digitally scanned and stored in the image database. Thus, the *CEDAR letter* database consists of two entities: writer data and document feature data. First, there are seven attributes of writers collected through the questionnaire data-sheet. They are:

- **Gender:** Male/Female
- **Age:**
  - under 15 years
  - 15 through 24 years
  - 25 through 44 years
  - 45 through 64 years
  - 65 through 84 years
  - 85 years and older
- **Handedness:** Left/Right
- **Highest level of education:** High school graduate/ higher
- **Country of Primary Education:** If USA, which state
- **Ethnicity:**

- Hispanic
- White (not Hispanic)
- Black (not Hispanic)
- Asian and Pacific Islander (not Hispanic)
- American Indian, Eskimo, Aleut

- **Country of Birth:** USA/Foreign

We built a database that is “representative” of the US Population. This has been achieved by basing our sample distribution on the US census data (1996 Projections) [7]. There are 510 female and 490 male population distributions and 36% of white ethnicity group, and so forth. The database contains handwriting samples of 1,000 distinct writers.

We asked each participant to copy out the *CEDAR letter* three times in his or her most natural handwriting. Thus

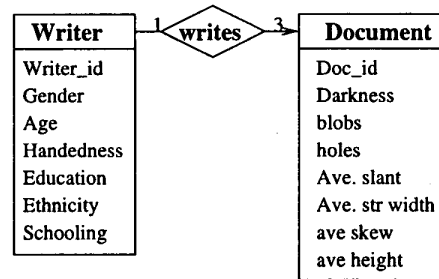


Figure 3. Entity and Relationship Diagram for CEDAR letter Database

the relationship between the writer entity and the document entity is one to many ( $m = 3$ ) as shown in Figure 3. In this data collection, provided uniform writing materials are used: plain unruled sheets and a medium black *Bic* round stic pen.

The experiment was designed to show the individuality of handwriting using the *CEDAR letter* image database [3]. Given two randomly selected handwritten documents, we can determine whether the two documents were written by the same person or not with 97% correctness performance. However, the handwritten items in practice are not *CEDAR letter* but can be any words or documents. The following section 3 shows the procedure to handle any questioned handwritten item utilizing the *CEDAR letter* database.

## 3. Procedure

In this section, we first outline the procedure and then explain each step with the example of the real ransom note

from JonBenet Ramsey murder case as shown in Figure 4.<sup>1</sup> To access the authorship confidence, we use the *dichotomy*

money and ~~have~~ <sup>and</sup> ~~order~~ <sup>order</sup>  
~~pick-up~~ <sup>pick-up</sup> at your daughter.  
 any deviation of my instructions  
 will result in the immediate  
 execution of your daughter. You  
 will also be denied her remains  
 for proper burial. The two  
 gentlemen watching over your daughter  
 do particularly like you so I  
 advise you not to provoke them.  
 Speaking to anyone about your  
 situation, such as police, F.B.I., etc.,  
 will result in your daughter being  
 beheaded. If we catch you talking  
 to a stray dog, she dies. If you  
 alert bank authorities, she dies.  
 If the money is in any way  
 marked or tampered with, she...  
 dies. You will be scanned for  
 electronic devices and if any are  
 found, she dies. You can try to  
 deceive us but be warned that  
 we are familiar with law enforcement  
 countermeasures and tactics. You  
 stand a 99% chance of killing  
 your daughter if you try to out-  
 smart us. Follow our instructions

(a)

beheaded

(b)

Figure 4. (a) scanned QD, a ransom note (b) Extracted TOI, "beheaded".

model [3] as a framework. Only difference, the most important issue of this paper, is how the support image database is simulated.

### 3.1. Outline of Procedure

#### 1. Questioned document item collection

- (a) Extract texts of interest (TOI) from QD.
- (b) Suspects write the TOI as well as CEDAR letter.

<sup>1</sup>a copy of the ransom note recovered in the JonBenet Ramsey murder case was released by District Attorney Alex Hunter. The image was obtained with permission from the URL "<http://www.thesmokinggun.com/archive/ransom2.html>".

- (c) Scan the QD, suspects' TOI and CEDAR letter.

#### 2. Simulated TOI image DB construction

- (a) simulate words in QD using common substrings in CEDAR letter.
- (b) (Word spotting) find selected words in the scanned CEDAR letter images for every writer.
- (c) (TOI segmentation) segment the word out of the CEDAR letter images and then characters of interest out of the word images.

#### 3. Finding the confidence of authorship

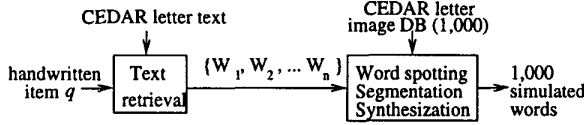
- (a) Generate the character level features.
- (b) Compute the feature distances.
- (c) Train an ANN using the simulated feature distances as input.
- (d) Build the same and different writer histograms from simulated data.
- (e) Find the ANN output value ( $T$ ) for the distances between the TOI in QD and suspects' TOI and find error probabilities.
- (f) Validate the result.

### 3.2. Description of Procedure

The first step to compare QD with suspects' writing is to extract texts of interest, TOI in short, from the questioned document. Possible TOI's are i) words that appear in the CEDAR letter ii) words that can be obtained from suspects' known writing exemplars, or an unusual word such as "beheaded" as shown in Figure 4 (b). All suspects are asked to copy the TOI as well as the CEDAR letter. The reason that suspects copy the CEDAR letter is for the purpose of the validation in the step 3 (f).

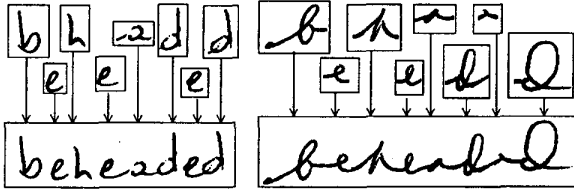
The next step is the simulated TOI image DB construction. synthesized TOI are generated from the CEDAR letter text. To do so requires the text retrieval using the string matching technique [10]. For example, to simulate the TOI, "beheaded", we first find all words that have the longest prefix-sub-string in CEDAR letter text. There is only one word {been} that have the same first two characters "be?????". Next, we find all words that have the longest suffix-sub-string in the CEDAR letter text and they are {referred, started, overworked, enjoyed, required, affected, passed, rushed}. They all have the same last two characters "?????ed". Amongst, we choose "referred". For the remaining parts, we repeat searching the longest sub-string that are neither prefix nor suffix. We choose {Co'he'n, Me'd'ic'a'l}.

After words are selected to generate the TOI, "beheaded", those words need to be spotted and segmented



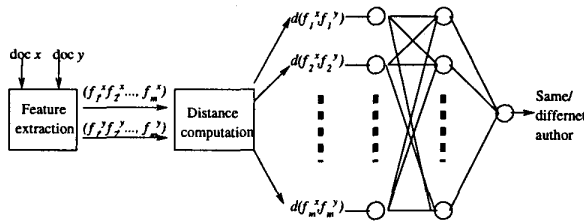
**Figure 5. Simulated word TOI database construction**

from the *CEDAR letter* document images. All characters to be synthesized to generate the *TOI* are also segmented from the words. Although there are automatic systems for spotting and segmentation, we use an image manipulating tool to manually extract characters. This is because the performances of current automatic systems are not yet as perfect as humans. This step is depicted in Figure 5 and sample synthesized *TOI*'s are shown in Figure 6.



**Figure 6. Synthesized TOI, beheaded**

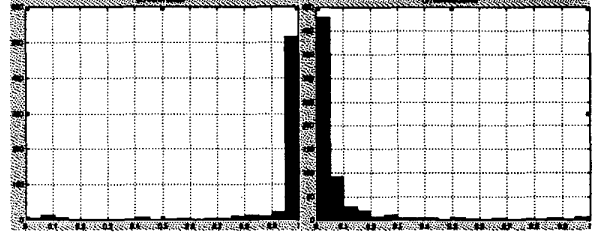
Finally, we find the authorship confidence of the *TOI*. Suppose the *TOI* from the handwritten item  $x$  has  $m$  character level features, then the *TOI* is represented in the vector of  $(f_1^x, f_2^x, \dots, f_m^x)$ . We use only the character level features here as word level features cannot be synthesized from the *CEDAR letter* database. If the *TOI* happen to be the exact word that appear in the *CEDAR letter*, we use the word level features. Those words that frequently appear in both common *English* and *CEDAR letter* are {*From, To, We, were, to, you, at, the, This, is, It, all, an, no, and* }, etc.



**Figure 7. Artificial Neural Network**

We generate all  $m$  features from the synthesized *TOI* and then take the distance values between two documents  $x$  and  $y$  for each feature  $\{d(f_1^x, f_1^y), d(f_2^x, f_2^y), \dots, d(f_m^x, f_m^y)\}$ . These distance values are fed into the *Artificial Neural Net-*

*work*. Figure 7 provides an overview of the entire process of the *artificial neural network*. For training, the output value of 1 is given if the two simulated handwritten items  $x$  and  $y$  were written by a same writer and 0, otherwise. We use the *Artificial Neural Network (ANN)* because it is equivalent to a sound multivariate statistical analysis [4]. When the output value of the *ANN* is above 0.5, we classify it as the “*identity*” that two documents were written by the same writer and otherwise, they were written by two different writers. The distance vectors are divided into training, validating and testing sets. Figure 8 shows the decision



**Figure 8. Decision Histogram on the testing set: (a) Within author distribution (Identity) (b) Between author distribution (Non-Identity)**

histogram on the testing set. The *within author distance* data, also know as *identity* tends to be voluminous toward the value 1 while the *between author distance* data or *non-identity* tends to be voluminous toward the value 0 as plotted in Figure 8 (a) and (b), respectively.

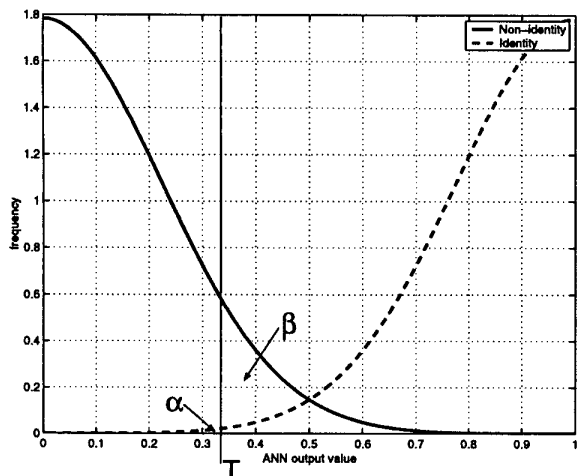
Now compute the distance vector between the *TOI* from the QD,  $q$  and that from the suspect,  $s$ :  $\{d(f_1^q, f_1^s), d(f_2^q, f_2^s), \dots, d(f_m^q, f_m^s)\}$ . When these values are fed into the *ANN* as an input vector, the *ANN* returns an output value. We call the output value as the query bar,  $T$  as depicted in Figure 9. The query bar gives *Type I*,  $\alpha$  and *Type II*,  $\beta$  errors.

$$\alpha = Pr(ANN(d^x, d^y) < T | \text{author}(x) = \text{author}(y))$$

$$\beta = Pr(ANN(d^x, d^y) \geq T | \text{author}(x) \neq \text{author}(y))$$

The *ANN* denotes the whole process in Figure 7 that takes two document images as input and authorship as output. The *type I error*,  $\alpha$  is the probability of others having similar handwriting attributes as the writer. Conversely, the *type II error*,  $\beta$  is the probability of how differently one writer may write.

Finally, we determine whether the query bar and error probabilities are valid by taking the distance vector between the suspects’ *TOI*'s and the synthesized *TOI* from suspects’ *CEDAR letter*. If the similarity between them is low, the confidence is low as well and thus reject the result. If the



**Figure 9. Simulated plot to illustrate  $T$ ,  $\alpha$  &  $\beta$  errors**

output of ANN with this input values is identified as non-authorship, then reject the test, otherwise, accept the result. This is necessary to prevent the unnatural simulation of the *TOI*. Hence, this final step measures the acceptability of the synthesized *TOI*'s.

#### 4. Conclusions

The work reported in this paper is applicable to the area of *Forensic Document Examination*. We have shown a method to access the authorship confidence of handwritten items utilizing the *CEDAR letter* database. It is a procedure for determining whether or not two or more digitally scanned handwritten items were written by the same person. Thanks to the completeness of the *CEDAR letter* database, it is a panacea for the analysis for any handwritten item by synthesization.

Albeit the simulating *TOI* saves enormous time to build the database to support the authorship confidence, it still takes a while to spot and segment the word and character out of the document unless the questioned *TOI* is already segmented and stored in the image database. The automatic line, word and character segmentation problem is one of open problems.

Standardizing procedures require a community-based effort. In order for the proposed procedure to be recognized as a part of the standard forensic document examination protocol, it is necessary to provide the validation that it uniquely identifies the writer using further statistical experimental design and protocol evaluation techniques.

#### Acknowledgments

This research has been possible funded by National Institute of Justice (NIJ) in response to the solicitation entitled *Forensic Document Examination Validation Studies* [11]: Award Number 1999-IJ-CX-K010.

#### References

- [1] Russell R. Bradford and Ralph B. Bradford. *Introduction to Handwriting Examination and Identification*. Nelson-Hall Publishers: Chicago, 1992.
- [2] Sung-Hyuk Cha and Sargur N. Srihari. Multiple feature integration for writer verification. In *Proceedings of 7th IWFHR 2000*, pages 333–342, September 2000.
- [3] Sung-Hyuk Cha and Sargur N. Srihari. Writer identification: Statistical analysis and dichotomizer. In *Proceedings of SS&SPR 2000 LNCS - Advances in Pattern Recognition*, volume 1876, pages –. Springer-Verlag, September 2000.
- [4] Richard O. Duda, David G. Stork, and Peter E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, Inc., 2nd edition, 2000.
- [5] Ordway Hilton. The relationship of mathematical probability to the handwriting identification problem. In *Proceedings of Seminar No. 5*, pages 121–130, 1958.
- [6] Roy A. Huber and A. M. Headrick. *Handwriting Identification: Facts and Fundamentals*. CRC Press LLC, 1999.
- [7] U.S. Department of Commerce. Population profile of the united states. Current Population Reports Special Studies P23-194, Semtember 1998.
- [8] Albert S. Osborn. *Questioned Document*. Albany, N.Y. : Boyd Print. Co., 2nd edition, 1929.
- [9] Rejean Plamondon and Guy Lorette. Automatic signature verification and writer identification - the state of the art. *Pattern Recognition*, 22(2):107–131, 1989.
- [10] Graham A. Stephen. *String Searching Algorithms*. Singapore:World Scientific, 1994.
- [11] Jeremy Travis. Forensic document examination validation studies. Solicitation: <http://ncjrs.org/pdffiles/sl297.pdf>, October 1998.