# Recognition of Handwritten Phrases as Applied to Street Name Images

Gyeonghwan Kim and Venu Govindaraju

Center of Excellence for Document Analysis and Recognition (CEDAR)
State University of New York at Buffalo
Amherst, New York 14228–2567

#### Abstract

A method for recognition of street name phrases collected from mail pieces is presented in this paper. Some of the challenges posed by the problem are: (i) patron errors, (ii) non-standardized way of abbreviating names, and (iii) variable number of words in a street name image. A neural network has been designed to segment words in a phrase, a street name in this case, using distances between components and style of writing. The network learns the type of spacing (including size) that one should expect between different pairs of characters in handwritten text. Experiments show perfect word segmentation performance at about 85% of cases. Unlike conventional methods, where lexicon entries are expanded to take care of all variations of prefixes and suffixes, substring matching is attempted only between the main body of a lexicon entry and the word segments of an image. Efforts to reduce computational complexity are successfully made by the sharing of character segmentation results between the segmentation and recognition phases. 83% phrase recognition accuracy is achieved on a test set.

#### 1 Introduction

Machine recognition of unconstrained handwriting is a challenging task. Recent research in recognition of general handwritten text has revealed a number of complex problems [1]. While recognition of isolated units of writing, such as a character or a word, has been extensively studied in literature [2, 3], emphasis on the next logical step of recognizing words and phrases has been lacking. One realizes the importance of recognizing phrases in applications where a computer program has to read fields containing handwritten responses. Such applications warrant segmentation of the field into its constituent words as the first processing step.

The stages in the architecture of a phrase recognition system are (i) segmentation, concerns extracting words from a phrase, (ii) word recognition, concerns handwritten word recognition algorithms with and without lexicon and (iii) postprocessing, concerns use of linguistic constraints to improve recognition performance. Our focus in this paper is on the first two stages. The importance of correct word segmentation, should be evident from the sequential architecture of the phrase recognition system. Correct word segmentation is of crucial importance since errors made in the segmentation stage get compounded by the time phrase recognition is performed.

Figure 1 illustrates an overview of the system described in this paper. Inputs to the system are a binary image of a phrase and a lexicon. The phrase in the example has three words: a prefix ("S"), followed by the 'main' body or name ("Michigan") and a suffix ("Ave"). The pre-processing step involves noise removal, slant correction, and smoothing. Objective of word segmentation is to return several sub-images, where each sub-image contains a single word. As shall be seen in the course of this paper, because of the design of our method, it is acceptable for the word segmentation module to return sub-images which contain parts of a word. In the example of Figure 1, the desired output is 3 sub-images containing the words: "S", "Michigan", and "Ave", respectively.

In this paper, we describe a recognition driven strategy for segmenting a handwritten phrase into its constituent words. Algorithms dealing with word segmentation that have been reported in the literature thus far are based primarily on analysis of geometric relationship of adjacent components in an image. Consequently, the leading and trailing ligatures of a word, often exaggerated in writings with flourish, hinder accurate segmentation. We overcome the problem of interference between ligatures of adjacent components

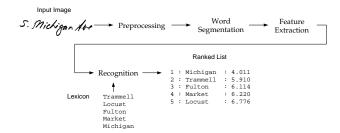


Figure 1: System overview

by making use of the knowledge of how often different characters get split and in what manner. This knowledge is available in the form of a table of statistics generated by our character segmentation algorithm applied to thousands of handwritten words. Computational complexity, a serious concern while matching lexicon entries and all possible combinations of segments, is also addressed by the use of the statistics [4].

The test bed used in this work is a dataset of street name images extracted from live mail pieces. Recognizing the street name is an integral part of an address interpretation system [5]. Street name images usually contain more than a single word and the number of words can vary. Different ways of abbreviating prefixes and suffixes also make the segmentation task difficult. Our approach begins with locating the bounding boxes of character segments. The center-line of each box is determined using the distribution of pixels within the box. Intervals and heights of the center lines are treated as key parameters and input to a neural network designed for locating word breaks. A lexicon is generated so that it contains only the keywords ('main' part of the name). Matching between the lexicon and the input image (street name) is performed by a dynamic programming based word recognizer [4].

Organization of this paper is as follows. Section 2 describes previous work on general word segmentation. Section 3 outlines general problems in the task of street name recognition and the requirements incumbent on any word recognizer that needs to be extended to handle phrases. Section 4 describes a strategy for locating word breaks in a street name image and the methodology of matching the main part of lexicon entries and a group of segments. Section 5 is about experiments and results. Summary of the work and some concluding remarks are presented in Section 6.

#### 2 Previous work

Most previous work in recognizing text assumes that words are already isolated or that isolation is trivial (e.g., words are separated by large amounts of white space) [6], or that words are written in boxes whose

location is known [7]. Few published reports describe methods of separating words in unconstrained handwriting [8, 9].

Deriving a computational methodology for word segmentation is quite complex. Humans use several cues to perform word separation: (i) spatial separation between components, (ii) presence of punctuation marks, (iii) presence of upper case characters following a string of lower case characters, (iv) transition between numerals and alpha characters, and (v) actual recognition of words prior to word separation. Among these, the cue of spatial separation is most commonly used.

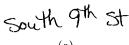
Computing the distance between adjacent components in itself is an open issue which warrants further research. The objective is to obtain an estimate of the inter-component gaps as perceived by humans. The most straightforward of methods computes the horizontal distance between the bounding boxes of adjacent components, where the bounding box of a component is defined as the smallest rectangle enclosing the component. The second method [9], uses run-lengths and  $Euclidean\ distances$  between connected components and heuristics. The heuristics handle cases where adjacent components do not have sufficient gap in the x and y direction. A third technique approximates the gaps between components by the distance between their convex hulls [10].

Methods based on spatial cues alone have limited success for two reasons. First, handwriting does not necessarily adhere to the rule of placing larger gaps between words than between characters. Second, the perceived space between components cannot be easily estimated by a 1-dimensional scalar metric.

# 3 Word recognition overview

In order to recognize phrases, such as street names, we must first preview the architecture of word recognizers and look for ways of extending the architecture to handle phrases. Word recognizers, traditionally require a lexicon of possible words to recognize the word in question. The lexicon is a list of possible words (phrases) that could possibly occur in that image. This lexicon is usually determined by the application.

Generating lexicon for the street name recognition task is somewhat involved. As can be seen in Figure 3, there are several variations of prefixes and suffixes (e.g., "Fwy" instead of "Freeway") that can appear in an image. Moreover, incorrect suffixes (and their omission altogether) introduced by patrons are common (e.g., "Street" instead of "Avenue"). One possible solution would be to expand each lexicon en-



(a)							
S9thSt	South9thSt	9thSt					
SNinthSt	SouthNinthSt	NinthSt					
S9thStreet	South9thStreet	9thStreet					
SNinthStreet	SouthNinthStreet	NinthStreet					
S9th	South9th	9th					
SNinth	SouthNinth	Ninth					
(b)							

Figure 2: An example of lexicon expansion: (a) a numeric street name (b) expanded lexicon entries for the image

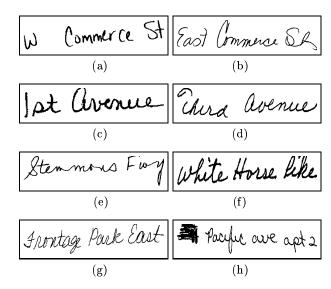
try in all conceivable ways (e.g., "Main Street" to "Main", "Main Street", "Main St"). However, this approach increases the lexicon size, consequently, the matching process becomes inefficient. A particularly difficult instance to be noted in this regard is that of numeric street names (Figure 2) wherein to account for a single street name the lexicon must contain 18 different forms of the name in which it can occur in a real image.

Additional complexities can occur due to the presence of extraneous words such as 'apartment' and 'suite' appearing at the end of a street name (Figure 3(h)). Certainly, the lexicon entries do not anticipate their presence. Furthermore crossing out of writing (Figure 3(h)) is an additional source of concern

The first task of a traditional word recognizer is to locate the character segmentation points. A segment can be a complete character or part of a character. Our recognition method (described in detail in [4]) proceeds as follows. A segment or a combination of consecutive segments is provided to a feature extraction module. The feature extraction module turns the segments into a multi-dimensional vector based on global and local features of the input segment(s). In the recognition module, comparison between the feature vector and code words of a character cluster (built during the training phase) is performed. The procedure is repeated for all segments and lexicon entries and the best matching entry is obtained using dynamic programming.

# 4 Methodology

Sample images of phrases extracted from street lines of addresses are shown in Figure 3. The number of words within the phrase can vary. Moreover, even the number of words in the 'main' part of a street name can be more than one, as in "White Horse" (Figure 3(f)). Typically, street name recognition involves segmenting the phrase into the constituent words and



**Figure 3:** Examples of street name images: (a - b) the same street name with different pre-directionals, (c - d) numeric street name, (e - f) post-directionals, (g) directional suffix, and (h) miscellaneous addons

recognizing the individual words. Several approaches have been described in the literature [1, 5, 9, 10, 11] to address this issue.

In the subsequent sections we present an elegant method of extending an existing word recognizer that operates on single words to handle phrases of several words. There are two major parts to the methodology: (i) segmentation and (ii) recognition.

#### 4.1 Segmentation

The character segmentation algorithm used in our system locates segmentation points between characters in a word and is based on the following two assumptions: (i) the number of segments per character must be at most 4, and (ii) all touching characters should be separated. Segmentation points are determined using features like ligatures and concavities [4].

We propose a word segmentation method (that divides a phrase into words) based on the character segmentation method described above. Our method uses additional information such as author's writing style. The style is captured by characterizing the variation of spacing between adjacent characters as a function of the corresponding characters themselves. The notion of expecting greater space between characters with leading or trailing ligatures is encoded into the segmentation scheme. Furthermore, we make the additional assumption that character segments must have at least one significant vertical stroke – a fair assumption given the structure of the English alphabet.

Consider the input street name image of Fig-

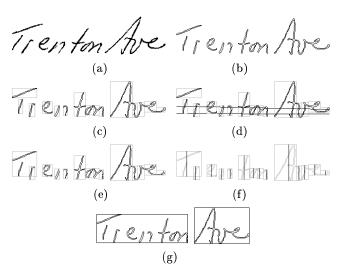


Figure 4: Word segmentation: (a) an input street name image, (b) slant normalized, (c) character based segmentation, (d) base line and middle line, (e) merged segments, (f) center line of each segments, and (g) segmentation result

ure 4(a). It is to be noted that the gap between n and A, is significantly less than gaps between r and e, and n and t. The input image is converted into a contour representation and the slant angle is estimated and corrected (Figure 4(b)). The original character segmentation algorithm segments the image into characters and sub-characters (Figure 4(c)).

Four reference lines have been traditionally used to extract spatial features: upper bound, upper half, lower half, and lower bound [2, 9]. However, all these lines cannot be accurately located, specially when dealing with phrases (as opposed to single words). Therefore, we choose to define two reference lines: base line and middle line (Figure 4(d)). The base line is the best-fit line of concavities in the lower contours. The middle line is the best-fit line of y-coordinates of the center of mass of each character segment. Based on the two reference lines, segments are examined. Rules are developed to determine which segments can be merged, if at all. For example, in Figure 4(d), a horizontal stroke over the vertical stroke of letters T and r is merged into the vertical stroke of T. Figure 4(e)shows bounding boxes after the merge. A vertical line is located based on the distribution of pixels inside the box as shown in Figure 4(f). The heights of the bounding boxes and the intervals between two adjacent boxes characterize the author's writing style in terms of inter-character and inter-word spacing.

We adopt a simple neural network to determine the segmentation points. Input parameters to the neural network are properties of bounding boxes described above. Let  $B_i$  and  $B_{i+1}$  denote a pair of adjacent

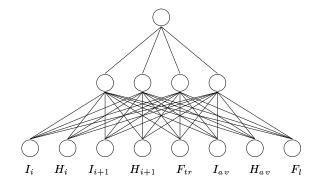


Figure 5: Structure of the neural network for word segmentation

bounding boxes after the merging, and let  $C_i$  be the center line of the bounding box  $B_i$ . Let  $I_i$  be the interval between  $C_i$  and  $C_{i+1}$  and  $H_i$  be the height of  $B_i$  and defined as follows:

$$I_{i} = X(C_{i+1}) - X(C_{i})$$

$$H_{i} = \begin{cases} Y_{u}(C_{i}) - BL(X(C_{i})) & \text{if } BL(X(C_{i})) \\ & \geq Y_{b}(C_{i}) \\ Y_{u}(C_{i}) - Y_{b}(C_{i}) & \text{otherwise} \end{cases}$$

$$(1)$$

where  $X(C_i)$ ,  $Y_u(C_i)$ , and  $Y_b(C_i)$  represent the x-coordinate, upper and bottom y-coordinates of a center line  $C_i$ .  $BL(\cdot)$  returns the y-coordinate of  $C_i$  at the baseline.

 $H_i$ s and  $I_i$ s are normalized and averages and standard deviations are computed. The neural network has 8 inputs units, 4 hidden units and 1 output unit as shown in Figure 5. Following 8 parameters are selected for inputs of the neural network: two consecutive intervals ( $I_i$  and  $I_{i+1}$ ) and heights ( $H_i$  and  $H_{i+1}$ ), average of intervals ( $I_{av}$ ) and heights ( $H_{av}$ ), flag for inter-component transition ( $F_{tr}$ ), and flag for the last segment ( $F_l$ ). It should be noted that the height information as well as interval information of segments are included.

Figure 4(g) shows the word segmentation result of the image of Figure 4(a) using the described neural network.

## 4.2 Recognition

The final goal is to find the best match between a word in the lexicon and the 'main' word in the street name image. We have described in detail a lexicon driven handwritten word recognition algorithm based on dynamic matching elsewhere to handle isolated words [4]. In the matching procedure, comparisons between feature vectors of several possible combinations of segments and reference feature vectors of codewords are made to find the best match. Since codewords are trained at the character level and a character can be

composed of up to four segments, a segment or a combination of segments is compared to the codewords of reference characters within a permissible window in the first phase of matching. For each match, the minimum distance value of each comparison is retained.

In the training phase of the word recognizer, segmentation statistics of each character are obtained. The segmentation statistics represent the possible ways in which a training character image can be split in to segments [12]. However, unlike [12], where segmentation statistics are used to determine the transition probability between segments of a character, we use the statistics to define the size of the matching window. It has been proven that the statistical information can be used to improve recognition accuracy as well as speed [4].

Ways to minimize the computational complexity are sought by maximizing the advantage of the lexicon driven approach. The variable duration is used for the early rejection of lexicon entries (Eq. 2).

$$max\_seg = \sum_{j=0}^{N_c(i)-1} dur(lex\_entry[j])$$
(2)

if ( num\_seg\_of\_word > max\_seg) reject;
else continue;

where  $N_c(i)$  represents number of characters in *i*-th entry in the lexicon,  $lex\_entry[\cdot]$  represents character string of the lexicon entry,  $dur(\cdot)$  returns variable duration of a character defined in [4], and  $num\_seg\_of\_word$  is the number of segments in a word obtained by the character segmentation algorithm.

# 5 Experiments and results

Experiments of extension of the recognition algorithm to street names follows. Street name images collected from live mail pieces were used for training and testing the neural network. 518 images were used for the training. 7802 patterns were collected by visual examination with desirable word segmentation points marked manually. Figure 6 shows a sample image and the corresponding test patterns used in training. The network was tested on two sets of data: the training set, and a new set of 107 images comprising of failure cases reported by the present handwritten address interpretation system at CEDAR [5]. Table 1 shows the word segmentation results on the two sets of images. 82–86% test images are perfectly segmented. Oversegmentation cases (5.0–7.5%) can be actually classified as correct segmentation since the recognition procedure will combine the broken parts anyway. This gives an effective segmentation rate of 91%.

Table 1: Test results of word segmentation

	# of images	perfect	overseg- mentation	error
set1	518	446	26	46
		(86.1%)	(5.0%)	(8.9%)
set 2	107	88	8	11
		(82.2%)	(7.5%)	(10.3%)

Table 2: Recognition performance - numbers in () represents the number of cases resulted in better recognition confidence comparing to that of Method1 (S: Success; F: Failure)

		Method1		Method2		Method3	
accuracy		S	F	S	F	S	F
		79	17	83	13	81	15
				(60)		(56)	
	$_{ m image}$	115.0		89.	0	122	.8
$_{ m time}$	recog	523.8		1087.7		581.6	
(msec)	total	638.8		1176.7		704.4	

Among the 107 street name images, 11 images which were classified as word segmentation error, were excluded during the recognition phase. Table 2 shows the test results in terms of recognition accuracy and speed using three different methods. Method1 uses the word recognizer and an expanded lexicon by regarding the entire street name block as a unit. Method2 uses the word segmentation method described in this paper to locate word breaks, and word segments are submitted to the word recognizer with the lexicon containing only the 'main' part of the names. It should be noted that the ratio of lexicon sizes used in Method1 against Method2 and Method3 is 4.6:1. Method3 uses the word segmentation method described in this paper as well. However, the difference between Method2 and Method3 is in the way in which they use character segmentation results. In Method2, character segmentation is performed twice, once in the beginning to collect information to feed to the neural network and once during recognition of word segments. In Method3, character segmentation results obtained during word segmentation are used during recognition as well. Average time (on Sparc 10) consumed by each method is shown in the lower half of Table 2. It needs to be noted that the number of cases resulted in better recognition confidence comparing to that of Method1 (the numbers in () in Table 2) is significantly increased, as well as increase of total number of correctly recognized cases.

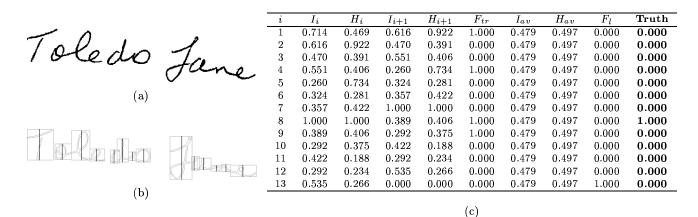


Figure 6: A sample image and training patterns: (a) original image (b) result of character segmentation and merging (c) training patterns

# 6 Summary

In this paper, we have introduced a strategy for phrase recognition as applied to street name images. A word segmentation algorithm was designed to overcome the drawbacks of conventional methods which use distance metrics of components in an image and are easily affected by ligatures of adjacent components. Height information as well as interval information of center-lines of character segments are included as input parameters to a neural network to determine the possible word breaks. This allows for the capture of space information between components as a function of adjacent characters. We have experimentally shown the effectiveness of our approach both computationally and in terms of accuracy. 83% phrase recognition accuracy is achieved on a test of difficult cases comprising of street name images extracted from mail pieces.

#### Acknowledgments

The authors would like to thank Uma Mahadevan, CEDAR, for the discussion on different distance metrics presented in Section 2.

# References

- V. Govindaraju, R. K. Srihari, and S. N. Srihari, "Handwritten Text Recognition," in Proceedings of the conference on Document Analysis System, Kaiserslautern, Germany, pp. 157-171, 1994.
- [2] R. M. Bozinovic and S. N. Srihari, "Off-line Cursive Script Word Recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 11, pp. 68-83, January 1989.
- [3] J. C. Simon, "Off-line Cursive Word Recognition," Proceedings of the IEEE, vol. 80, no. 7, pp. 1150-1161, 1992.
- [4] G. Kim and V. Govindaraju, "Handwritten Word Recognition for Real-Time Applications," in *Third*

- International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada, pp. 24– 27, 1995.
- [5] V. Govindaraju, A. Shekhawat, and S. N. Srihari, "Interpretation of Handwritten Addresses in US Mail Stream," in *The third International Workshop on Frontiers in Handwriting Recognition, Buffalo, New York*, pp. 197–206, 1993.
- [6] H. Baird and K. Thompson, "Reading chess," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 12, pp. 552–559, 1990.
- [7] R. O. Duda and P. E. Hart, "Experiments in Recognition of Hand-printed Text," in AFIPS Conference Proceedings, pp. 1139-1149, 1968.
- [8] E. Cohen, J. Hull, and S. N. Srihari, "Understanding Handwritten Text in a Structured Environment: Determining ZIP codes from addresses," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 5, no. 1&2, pp. 221-264, 1991.
- [9] G. Seni and E. Cohen, "External Word Segmentation of Off-Line Handwritten Text Lines," Pattern Recognition, vol. 27, no. 1, pp. 41-52, 1994.
- [10] U. Mahadevan and R. C. Nagabushnam, "Gap Metrics for Word Separation in Handwritten Lines," in Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada, pp. 124-127, 1995.
- [11] J. C. Simon, O. Baret, and N. Gorski, "A System for the Recognition of Handwritten Literal Amounts of Checks," in Proceedings of the conference on Document Analysis System, Kaiserslautern, Germany, pp. 135-155, 1994.
- [12] M.-Y. Chen, A. Kundu, and J. Zhou, "Off-Line Hand-written Word Recognition Using a Hidden Markov Model Type Stochastic Network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 16, pp. 481–496, May 1994.