RECOGNITION OF OFFLINE HANDWRITTEN WORDS AND ITS EXTENSION TO PHRASE RECOGNITION

by Gyeonghwan Kim

A dissertation submitted to the Faculty of the Graduate School of the State University of New York at Buffalo in partial fulfillment of the requirements for the degree of Doctor of Philosophy

April 1996

© Copyright 1996 by Gyeonghwan Kim

COMMITTEE MEMBERS

CHAIRMAN	Dr. Sargur N. Srihari
	Computer Science, SUNY at Buffalo
Members	Dr. Ramalingam Sridhar Electrical Engineering, SUNY at Buffalo
	Dr. Raj Acharya
	Electrical Engineering, SUNY at Buffalo
	Dr. Venugopal Govindaraju
	Computer Science, SUNY at Buffalo
Outside Reader	Dr. Ching Y. Suen
	Computer Science, Concordia University, Montreal Canada

This dissertation was defended on April 10, 1996

iii

To My Parents for their prayer

•

and My Wife

for her love

ACKNOWLEDGMENT

First of all, I am deeply grateful to my advisor Prof. S. N. Srihari for providing me the opportunity to conduct this research at CEDAR, a friendly and perfect environment. He has funded, encouraged, and guided me in various aspects.

I would like to thank Prof. R. Sridhar for his understanding and encouragement during the early days of my study in this country, and serving as a committee member.

Also, I would like to thank Prof. R. Acharya for serving as a committee member in spite of his busy schedule.

My special thanks also go to Dr. V. Govindaraju for his friendship, support, and help in various ways as a colleague and committee member. Without his help, my thesis could not be completed in such a wonderful shape.

I would like to extend my thanks to all members of CEDAR. Special thanks to members of HWAI group who helped me with bright suggestions: Evelyn Kleinberg, Sriganesh Madhvanath, and Uma Mahadevan for their so many challenges and ideas; Dr. John Favata and Dr. Mou-Yen Chen, designers of HWRs at CEDAR, providing me good motivations; Heike Jones, Eugenia Smith, and Rita Hamilton, secretaries of CEDAR, for their kind help.

Finally, I would like to thank to all of my family members in my country. I always could feel warmth from their support, prayer, and encouragement. And to my lovely son for giving me delight, joy, and peace everyday.

v

ABSTRACT

The main challenges in the implementation of an offline handwriting recognition system stem from the vast diversity in writing styles. An attempt to solve the problem raises several issues: (i) appropriate image representation and efficient image handling procedures using the representation, (ii) design of a matching algorithm, and (iii) scalability. This thesis describes the implementation of a handwritten word recognition system and its extension to recognition of phrases. Our methodology places no constraints on writing styles. Furthermore, it handles phrases with several words and a mixture of purely alphabetic and numeral characters.

Each module of our methodology optimizes for both performance as well as speed. Chain code based image representation is used for preprocessing, segmentation and feature extraction. An oversegmentation scheme is used to handle diverse writing styles. Ligatures and concavity features are used to locate all possible segmentation points. Chain code segment images are converted to global and local feature vectors in an efficient manner.

The concept of variable duration is introduced to make the matching algorithm fast and accurate. It uses character segmentation statistics obtained during the training phase. The variable duration represents the possible ways in which a character image can be split into segments. The best match between the input image and lexicon entries is made using dynamic programming. The matching window size is controlled by factors such as the variable duration of the character, the number of segments of the image, and the number of characters in a lexicon entry.

Two applications are presented to illustrate the extension of our methodology to phrases – (i) recognition of street name images, and (ii) reading bank check amounts. A trainable word segmentation approach using a neural network is developed to circumvent the problems in conventional methods (which mainly rely on the analysis of geometric relationship between adjacent connected components). Our approach uses additional information such as author's writing style which is captured by characterizing the variation of spacing and height between adjacent major character segments. In the check recognition application, we use a single recognition engine serially rather than the commonly used parallel methods. Experimental results prove the efficiency of our methodology.

vi

Contents

Ac	Acknowledgment v								
Abstract									
1	Intr	oduction	1						
	1.1	Taxonomy of handwriting recognition	2						
	1.2	Problem definition	5						
	1.3	Motivation and objectives	9						
	1.4	Thesis scope	11						
2	Rev	iew of Handwritten Word Recognition	13						
	2.1	Properties of handwritten words	13						
		2.1.1 Image quality	13						
		2.1.2 Human factors	16						
		2.1.3 Miscellaneous factors	18						
	2.2	Background research	19						
		2.2.1 Early work	19						
		2.2.2 Segmentation strategies	20						
		2.2.3 Matching methodologies	22						
		2.2.4 Features	23						
		2.2.5 Postprocessing	24						

		2.2.6 Online word recognition	25
	2.3	Related fields	25
		2.3.1 Speech recognition	25
		2.3.2 Optical character recognition	27
		2.3.3 Hardware	27
		2.3.4 Psychological aspects	28
	2.4	Summary	28
3	Ima	ge Processing for Handwriting Recognition	29
	3.1	Image representation	30
	3.2	Chain code based image processing	32
		3.2.1 Noise removal	32
		3.2.2 Slant angle estimation and correction	33
		3.2.3 Smoothing contours	36
		3.2.4 Stroke width estimation	40
	3.3	Experiments and Results	42
	3.4	Summary	43
4	Han	dwritten Word Recognition Engine	44
	4.1	Segmentation	48
		4.1.1 Segmentation procedure	48
		4.1.2 Segmentation statistics	50
	4.2	Features	50
	4.3	Recognition	52
		4.3.1 First stage matching	53
		4.3.2 Second stage matching	57
		4.3.3 Locating character boundaries	61
	4.4	Experiments and results	62

		4.4.1 Timing profile	63
		4.4.2 Recognition performance	63
		4.4.3 Effect of introducing the variable duration	63
		4.4.4 Failure analysis	63
	4.5	Summary	68
5	Stre	eet Name Recognition	69
	5.1	Previous work in word segmentation	72
	5.2	Problem description	74
	5.3	Methodology	76
		5.3.1 Segmentation	78
		5.3.2 Recognition	82
	5.4	Experiments and results	84
		5.4.1 Word segmentation	84
		5.4.2 Recognition	86
		5.4.3 Selection of input parameters for the neural network	88
	5.5	Summary	89
6	Rea	ding Bank Check Amounts	91
	6.1	System overview	93
	6.2	Segmentation	95
		6.2.1 Problem definition	95
		6.2.2 Word segmentation based on character segmentation	96
	6.3	Recognition of amounts	102
		6.3.1 Legal amount recognition	102
		6.3.2 Courtesy amount recognition	104
		6.3.3 Cross validation between the legal and courtesy amounts	104
	6.4	Experiments and results	108

Bi	Bibliography 1										127																
A	Trai	ining																									115
	7.2	Contr	ibuti	ons		٠		•	•••		•	•	• •	•	• •	•••	•		•	•	•	•	• •	•	•	•	. 114
	7.1	Sumn	nary	• •		•	•••	•			•	•		•	• •		•		•	•	•	•	• •	•	•	•	112
7	Con	clusio	ns																								112
	6.5	Sumn	nary	•••		•		•	•••			•					•	• •	•	•	•	•	• •	•	•	•	. 111
		6.4.2	Rec	ogni	tion	•	•••	•		•••	•	•		•	• •		•		•	•		•					. 108
		6.4.1	Seg	men	tatic	n		•	•••	•••	٠	•		•				• •	•	•	•	•		•	•	•	. 108

List of Tables

3.1	Smoothing operation type table: ${\boldsymbol{p}}$ - slope of the preceding compo-	
	nent, ${\bf c}$ - slope of the current component, and ${\bf n}$ - slope of the next	
	component	38
3.2	Steps of smoothing procedures for Figure 3.4(g): box represents win-	
	dow for classification	39
3.3	Timing analysis (note: feature extraction time includes time taken	
	for combining segments before feature extraction)	42
3.4	Impact of the preprocessing (in %)	43
4.1	Segmentation statistics	51
4.2	First stage of matching procedure	57
4.3	Second stage of matching procedure	58
4.4	Variable duration (min:max)	59
4.5	Boundary of each character	61
4.6	Timing profile in <i>msec</i> on Sparc 10	64
4.7	Recognition accuracy (in %)	64
4.8	Failure analysis	67
5.1	Test results of word segmentation	86
5.2	Recognition performance - numbers in () represent the number of	
	cases resulted in better recognition confidence comparing to that of	
	Method1 (S : Success; F : Failure)	88

5.3	Word segmentation performance with different inputs	89
6.1	Issues specifically related to reading each field	92
6.2	Lexicon entries used for the recognition of the legal amount	99
6.3	Interpretation of the flags used for estimating number of digits in	
	the amount. note: confidence values of the top choices are also	
	examined for the amounts with questionable number of digits	107
6.4	Recognition performance	109
6.5	Sources of the rejection	111
A.1	Clustering results of characters	117

•

List of Figures

1.1	Types of handwriting : (a) hand-printed, alpha-numeric, (b) cursive,	
	and (c) mixed	6
1.2	Coarticulation: the character " i " has different shape in a script	7
2.1	Properties of handwritten words: Image factors - (a) speckle noise,	
	(b) noise from feeder, and (c) incorrect operation of scanner	14
2.2	Properties of handwritten words: Image factors - (a) noise caused by	
	texture/color of paper, (b) thick pen stroke, (c) background image,	
	and (d) use of poor pen	15
2.3	Binarization results for a gray image (with 256 levels) shown in (a)	
	with different threshold values: (b) 80, (c) 132 (chosen by analysis of	
	the entire image), and (d) 180	16
2.4	Properties of handwritten words: Human factors - (a) slant (vertical)	
	and (b) skew (horizontal)	17
2.5	Properties of handwritten words: Human factors - (a) touching char-	
	acters and (b) fragment in strokes	17
2.6	Properties of handwritten words: Human factors - (a) spelling error	
	and (b) abbreviation	18
2.7	Properties of handwritten words: Miscellaneous factors - (a) under-	
	line and (b) leftovers from the other fields	19
3.1	Chain code representation: (a) slope, (b) structure	31

3.2	Noise handling: (a) small components including single pixel compo-	
	nents are eliminated, (b) an image rejected due to severe background	
	noise	33
3.3	Noise handling: (a) contour representation of an image, extraction	
	of upward (b) and downward (c) vertical lines, (d) correction of x-	
	coordinates based on the estimated angle, (e) connecting broken	
	points. (f) smoothing contours	35
3.4	(a - f) Types of smoothing operations - for each type, left sequence	
	turns into right sequence after the smoothing operation (\times rep-	
	resents elimination of the components and \bigcirc represents starting	
	point) and (h) shows how the algorithm is working for the circled	
	area of the example image of (g)	37
3.5	Estimating average stroke width: (a) a sample word image (b) mea-	
	suring y-distances for each x-coordinate and analyzing occurrences	
	of stroke widths of image (a)	41
4.1	Word recognition methodology	45
4.2	Matching : (a) segmentation points, (b) confidence of the match, and	
	(c) matching paths and confidence	47
4.3	Segmentation: (a) original test image (b) chain code representation	
	after slant normalization (c) splitting upper- and lower- contours (d)	
	ligatures (e) concavities (f) segmentation points	49
4.4	Splitting and merging of segments: (a) segmentation result (b) exam-	
	ples of splitting and merging, and corresponding notations	55
4.5	An example of matching procedure: (a) segmentation result of an	
	image (b) dynamic matching with variable duration for the image (a)	60

4.6	Performance and speed comparison: (a) performance - the system	
	with variable duration is represented by solid line and gives the best	
	recognition accuracy for all lexicon sizes, (b) speed of the system with	
	variable duration is between the two other systems.	65
4.7	Error Analysis: (a) background noise, (b) noise introduced by scan-	
	ning device, (c) oversegmentation, (d) unable to segment, (e) recogni-	
	tion error (confused with a similar entry in lexicon), and (f) recogni-	
	tion error (due to tough edge)	66
5.1	System overview: Input is a phrase of multiple words and a lexicon	
	of single words (extracted from the main part of a street name). Out-	
	put is a ranked lexicon sorted in the order of "goodness" of match	
	between the image and lexicon entries.	71
5.2	Inter-component gaps using geometric relationship between compo-	
	nents generated by (a) bounding box method (b) run length based	
	method and (c) convex hull based method	74
5.3	An example of lexicon expansion: (a) a numeric street name (b) ex-	
	panded lexicon entries for the image	76
5.4	Examples of street name images: (a) street line address, (b - c) the	
	same street name with different pre-directionals, (d - e) numeric	
	street name, (f - g) suffix, (h) directional suffix, and (i) miscellaneous	
	addons	77
5.5	Word segmentation: (a) an input street name image (b) slant normal-	
	ized (c) character based segmentation (d) base line and middle line	
	(e) merged segments (f) center line of each segments	79
5.6	Merging bonding boxes of character segments: (a) a box completely	
	belonged to another, (b) a box partially overlapped, and (c) a box	
	overlapped in horizontal direction without touching	80

÷

5.7	Structure of the neural network for word segmentation	81
5.8	Word segmentation result on image of Figure 5.5	82
5.9	Matching word segments with lexicon entries	83
5.10)A sample image and training patterns: (a) original image (b) result of	
	character segmentation and merging (c) training patterns	85
5.11	Examples of segmentation results from our system: (a - d) segmen-	
	tation points perfectly located; (e - g) failures; (h) over-segmentation	
	cases	87
6.1	The components of a typical bank check	93
6.2	Check reading system architecture	94
6.3	Examples of amounts in which it is difficult to locate the boundary	
	between dollar and cent amounts: (a) courtesy amounts (b) legal	
	amounts	96
6.4	Examples of legal amount make segmentation difficult: (a - b) not	
	enough spaces between words, (c) long horizontal line at the end of	
	the field	97
6.5	Word segmentation (based on character segmentation): (a) original	
	input image (b) character segmentation result after slant correction	
	(c) centerline of each segment after merging (d) word segmentation	
	result	98

6.6	Recognition of the segments: (a) segment number (b) A block at row	
	position 1 and column position 2 represents the recognition result	
	obtained by submitting to a word recognizer a combination of seg-	
	ments one and two which in turn correspond to the first and the	
	second parts of the word "hundred". The number before each word	
	represents the matching distance. A smaller value indicates a better	
	match. Boxes indicate the recognition results that are used to feed	
	the lexicon to the subsequent stage of courtesy amount recognition.	101
6.7	Recognition results for a hypothesis	102
6.8	Flags for estimating the order of magnitude of the legal amount	104
6.9	Matching between amount candidates from the legal amount recog-	
	nition and segments in the courtesy amount	105
6.10	Examples of amounts images : (a - d) correctly recognized and (e - h)	
	rejected cases	110
A. 1	Samples of clustering results for (a) character \mathbf{G}/\mathbf{g} and (b) character	
	N/n	116

Chapter 1

Introduction

Is paper going to disappear as computer and related technologies advance? During the history of machine reading there have been speculations on whether the need for paper will eventually disappear (Pavlidis, 1993). However, paper is still the most commonly used media for saving and transporting people's idea. Quite to the contrary, computers have increased the amount of paper rather than reduce it. As long as there is paper around there will be a need for machine recognition of text on paper. Volumes of computer generated documents have to be dealt with by computer programs. As computer technology advances the need for systems which can handle real problems increases.

Recognition of handwriting by a computer program addresses the above mentioned need. The objective of the task is to interpret images containing handwritten scripts and to produce the equivalent ASCII characters with or without a lexicon. To achieve this objective, various technologies developed and used in general pattern recognition are involved; namely, image processing, segmentation, features, classification, and postprocessing using contextual information. Research in this area has been intense for some years now and considerable progress has been made in specific applications, including reading addresses on mail pieces (Cohen et al., 1994; Cohen et al., 1991b; Govindaraju et al., 1993; Srihari, 1992), reading address blocks on tax forms (Belaid et al., 1995; Srihari et al., 1995b; Yuan et al., 1995), extracting census data on forms (Breuel, 1994b; Madhvanath et al., 1995), reading amounts from bank checks (Agarwal et al., 1995; Dimauro et al., 1993; Graf et al., 1995; Moreau et al., 1991; Paquet & Lecourtier, 1991; Simon et al., 1994), and interpreting general text (Srihari et al., 1995a).

When implementing a system targeted for real applications, trade-off between speed and accuracy (Baird, 1993; Tappert, 1991) is pivotal. Unfortunately, the issue of speed largely has been neglected by the research community.

This research pertains to the development of a handwritten word recognition system which is efficient in speed as well as accurate in recognition. The design issues that arise are: (i) implementation of efficient image processing procedures and (ii) developing appropriate matching algorithms.

The methodology of handwritten word recognition is applicable to interpretation of interpreting handwritten phrases as in street name recognition from mail pieces and text amounts in bank checks.

1.1 Taxonomy of handwriting recognition

Handwriting recognition is the task which interprets images containing objects written by hand and extracts the underlying intent. There are many factors which affect the structure and the performance of handwriting recognition systems. We describe some of the important distinctions to be made in characterizing the handwriting recognition problem.

• **online and offline** - These are based on the manner in which handwriting is captured for recognition. Online handwriting recognition means that the computer recognizes the writing while the user writes. This requires a transducer that captures the writing during the writing process. The most common of these devices is the tablet or digitizer. Offline handwriting recognition, on the other hand, is performed after the writing is completed. An optical scanner converts the image of the writing into a bit pattern (Tappert et al., 1990). Lack of temporal information in offline data is a significant difference between the two and this makes the segmentation and recognition of the offline handwriting difficult.

• units of handwriting - Units of meaningful handwriting can be classified into following four categories: (i) character, (ii) word, (iii) phrase and (iv) sentence (Srihari, 1993). Even though there are logically lower level units, such as pixels and strokes, they do not have meaning in terms of recognition, but they are implicitly used during the recognition of the meaningful units. Character is the primitive recognition unit. Basically characters consist of 52 upper/lower characters and 10 numerals. Depending on applications, several special characters can be included. Pixels and features are primitives that are commonly used in recognition. Unlike machine print, matching at the pixel level is almost impossible in handwritten character recognition. Word is defined as consisting of characters. Recognition is a three-step process based on the recognition results of constituent characters: (i) segmentation of a word image into characters, (ii) recognition of segmented characters, and (iii) contextual postprocessing. Postprocessing techniques differ depending on matching strategies used. *Phrase* is a legal combination of words. If the total number of phrases is not too large, it is possible to deal with phrase as a recognition unit, like a word, by ignoring spaces between words (wholephrase approach). On the other hand, in a word-based approach, each word is treated as a unit and word segmentation and postprocessing steps are

necessary. *Sentence* consists of words and phrases. Syntactic, statistical, or hybrid models are used to interpret sentences.

- writing styles Depending on writing habits and constraints given by guide lines (as in forms) writing style can vary (*e.g.*, boxed, spaced, cursive and mixed (Tappert, 1986)). Boxed and spaced writings, which are frequently observed in the forms designed for OCR processing to minimize problems in the segmentation and recognition, are the most constrained. Writers are required to ensure that characters do not overlap with/without help of preprinted boxes. All of these styles require that each character is written discretely. On the other hand, connected writing which lacks such constraints includes cases where a few characters although written discretely are touching. It can also be mixed with cursive style which uses different character shapes altogether.
- writer dependency This distinction is the form of enrollment users must make for the system to perform as designed. *Writer-dependent* systems require a potential user to provide sufficient data (can be substantial) for training. Such systems can offer the best performance by avoiding the need to capture variability of writing styles across writers. *Writer-independent* systems require no enrollment data. Various styles of handwriting collected from a large group of people are used for training.
- **context** Systems designed for recognizing words and phrases implicitly incorporate contextual variation in the data. *Context-dependent* systems use contextual information during the segmentation and recognition. *Contextindependent* systems usually perform segmentation and recognition separately.

• **lexicon** - *A priori* statistics captured from training data imposes high-level constraints during the recognition. A common approach is to restrict the output of the system to words present in the lexicon. This works well with applications dealing with small lexicons. However, it is difficult to construct a lexicon which deals with general text. *n*-gram language models are commonly used.

1.2 Problem definition

Reading unconstrained handwritten *words* is the main objective of this thesis. There are many issues that must be addressed in order to accomplish this challenging task. Moreover, this thesis emphasizes implementation of a fast recognizer that can be used in real time systems. Therefore, problems which have been ignored in research domain thus far are legitimate issues. In this section, we elaborate on these very issues.

• **diversity in writing style** - The challenge of implementing such a fast system stems mainly from the wide variety of writing styles. A number of categories can be used for classification of the writing styles such as cursive or discrete and extent of touching between characters. Depending on the assumption employed in the recognition phase, the classification can be an issue. For example, some recognizers assume input images are purely cursive. However, we deal with all types of writing. Therefore, we do not attempt a strict classification of writing styles (Chen et al., 1994; Favata, 1993; Tappert, 1986). A rough classification is used: *cursive, hand-printed* and *mixed* (Figure 1.1) taking into account only the shape of characters and not their connectivity per say.

IN AV BUFFALO Ruffalo (a) Albaner • • • (b)Bu ffalo (c)

Figure 1.1: Types of handwriting : (a) hand-printed, alpha-numeric, (b) cursive, and (c) mixed

A character may have different shapes depending on the shapes of the surrounding characters. This phenomenon is frequently observed in the speech processing and is called the *coarticulation* effect (Rabiner & Juang, 1993). An assumption in continuous speech that words are pronounced similar to the way they are pronounced in isolation is reasonable in most cases. However, sometimes, the pronunciation differences are quite radical rendering a boundary sound deleted or replaced with a totally different sound. Similar phenomenon is observed in handwriting, especially cursive script. For instance, in Figure 1.2 the shape of the character "*i*" is affected by the shape of the surrounding characters and is different from its second occurrence.

IN AV WM BUFFALO Buffalo (a) Albane • (b) Bu ffalo $\{c\}$

Figure 1.1: Types of handwriting : (a) hand-printed, alpha-numeric, (b) cursive, and (c) mixed

A character may have different shapes depending on the shapes of the surrounding characters. This phenomenon is frequently observed in the speech processing and is called the *coarticulation* effect (Rabiner & Juang, 1993). An assumption in continuous speech that words are pronounced similar to the way they are pronounced in isolation is reasonable in most cases. However, sometimes, the pronunciation differences are quite radical rendering a boundary sound deleted or replaced with a totally different sound. Similar phenomenon is observed in handwriting, especially cursive script. For instance, in Figure 1.2 the shape of the character "*i*" is affected by the shape of the surrounding characters and is different from its second occurrence.

Dennesselle

Figure 1.2: Coarticulation: the character "*i*" has different shape in a script.

image representation and processing - Selecting an appropriate image representation format for an application is often motivated by speed versus accuracy trade-offs. Word images captured by scanners store two-dimensional position information and intensity of objects. Traditional image processing techniques have dealt with the two-dimensional representations that need tremendous amount of processing time due to the nature of the structure. There have been efforts to use simple and reduced descriptions with minimum loss of information that facilitate easy matching. However, non-conventional representations may add additional complexity to the processing. Using a mixture of image representations is commononly used (Favata, 1993; Kimura et al., 1993).

1 . A.

• **segmentation** - The task of segmentation is complex due to the diversity of writing styles and variation of size and stroke width within a word image. Two main problems in the recognition methodology based on segmentation are (i) locating accurate segmentation points between character boundaries and (ii) dealing with various character shapes split by segmentation points. In case of hand-printed script (Figure 1.1(a)) segmentation is a relatively simple task. In case of cursive script (Figure 1.1(b)) and mixed script (Figure 1.1(c)) segmentation is relatively hard. Several approaches have been explored. Typically, words are first classified into a category of their type and subsequently different schemes are applied depending on the script type

(Favata, 1993; Ganzberger et al., 1994).

Location of exact segmentation points to split touching characters is difficult if it is independent of recognition. The task of segmentation is integrally coupled to the recognition methodology and is analogous to segmenting continuous speech (Rabiner & Juang, 1993).

- **matching methodologies** Since a word is essentially a sequence of characters, a natural approach to word recognition is to segment the word into characters and recognize the individual characters using optical character recognizers (OCR) with or without lexicon. In most applications it is reasonable to assume that a lexicon is provided. The lexicon can be either static or generated dynamically. The task of address interpretation on a mail piece is an example of an application where the lexicon is generated dynamically. The zip code on the address provides all possible city names or street names as the lexicon. Different zip codes provide different lexicons (Govindaraju et al., 1993). An example of an application with static lexicon is the reading of amounts on bank checks. The lexicon is fixed in this case to about 40 words. Depending on the segmentation scheme and lexicon handling strategy, the matching methodology can vary. Choosing the most appropriate matching technique is a research issue.
- **extension from words to phrases** Direct extension of the concept used in the word recognition scheme to the recognition of phrases is a research issue. Expecting an isolated word image as an input unit is not realistic in all real-world applications. Quite often, a field consists of more than a word. Locating word gaps and interpretation of word segments are typical problems must be addressed.

1.3 Motivation and objectives

Increasing demands from real applications has brought efficiency of algorithms to the forefront. So far research has been performed under the assumption that the speed issues can be ignored. Developing algorithms suitable for real applications is the basic motivation of this research. We focus on the following issues:

image handling - Two-dimensional image representations which have been
widely used in conventional image processing require tremendous amount
of processing time.

A number of image representations developed for description of connected components have been considered. Chain code, a popular representation, is adopted for image representation for all image handling procedures in our system.

• **design of a recognition engine** - Traditional approaches using the *segmentation-OCR-postprocessing* paradigms are basically an extension of character recognition. Segmentation of an input word image is followed by matching of segment(s) and templates of characters (OCR).

Empirical observation shows that the size of a lexicon which is the collection of words that can be matched with the image is often small in many practical applications. For example, in the application of interpretation of addresses on mail pieces, the zip code on the address provides all possible city names or street names as a reasonably sized lexicon. In the application of reading bank check amounts, the lexicon is fixed and size is less than 40 words. In cases with large lexicons, a fast lexicon filter based on wholistic features can reduce the number of lexicon entries significantly (Madhvanath & Govindaraju, 1993). A method of word recognition that uses word models (as opposed to character models) is considered in this research. The key concept underlying this approach is the early involvement of lexicon in the recognition process. The word image is compared with only those words which are present in the lexicon thus eliminating any need for postprocessing.

Because of the role of lexicon during matching, information collected during the training phase can be used to reduce the computational complexity of the matching phase. Statistics on character segmentation are collected during training and used for determining the size of the matching window.

• extension toward application of phrase recognition - The word recognizer introduced in this thesis is designed for the recognition of isolated words. When the system is used for the recognition of fields having more than a word, two approaches can be considered: (i) treat the entire field as a single unit by ignoring spaces between words both in the image and lexicon, and (ii) an external module provides isolated words after word segmentation as input to the system with a corresponding lexicon.

Location of word gaps is not a simple problem. Concepts used in word recognition can be directly extended and applied to the problem of phrase recognition. As in character segmentation, certain features are extracted and used for locating word gaps. In contrast to conventional approaches that use the geometric relationship between adjacent components, a method using results of character segmentation, which is the necessary step for the word recognizer, is developed. The concept of variable duration which is used in word recognition is extended and used in phrase recognition to avoid exhaustive matching. Recognition of the main body of a street name and reading bank check amounts, are developed to verify our methodology.

1.4 Thesis scope

The chapters in this thesis are organized as follows.

Chapter 2 reviews nature of handwritten words, basic concepts and techniques used in HWR, and relevant literature.

In Chapter 3, all preprocessing procedures implemented using chain code representation are described. This includes noise removal, slant detection and correction, smoothing contours, and estimation of average stroke width. Experimental results are presented to show the efficiency of the procedures.

Chapter 4 describes functional modules of our handwritten word recognition engine. Features for character segmentation and strategies for determining segmentation points are described. Feature vectors used in the matching procedure are defined and their extraction procedures are described. Concepts used in matching are described and performance of the recognition engine is measured in terms of both speed and the accuracy.

Chapter 5 presents a system for recognition of street name from address blocks as an extension of the single word recognition engine. Word segmentation using a neural network is introduced.

Chapter 6 describes a system for reading amounts from bank checks as another application of phrase recognition. Cross validation using recognition of the legal and the courtesy amounts is introduced. Recognition based word segmentation is applied to collect meaningful amounts. A digit string of the courtesy amount is treated as a word. Chapter 7 highlights the contribution of this research, draws conclusions and touches upon possible future work.

Chapter 2

Review of Handwritten Word Recognition

In this chapter, we describe general aspects of handwritten word recognition (HWR). This includes the properties of handwritten words, technologies used in handwriting recognition problems, and background research.

2.1 Properties of handwritten words

2.1.1 Image quality

Factors which determine image quality are closely related to each other.

• noise from imaging devices - Images are collected using digitizing scanners, and often transmitted via communication devices, such as satellites and transmitters. Random speckle noise is the commonly present form of noise (Figure 2.1(a)). Roller traces of automatic feeding mechanisms in scanners also contribute to noise (Figure 2.1(b)). Incorrect operation of scanner results in significant damage of the image (Figure 2.1(c)).



(a)



(b)



(c)

Figure 2.1: Properties of handwritten words: Image factors - (a) speckle noise, (b) noise from feeder, and (c) incorrect operation of scanner

• quality of writing media (pen and paper) - Characteristics of pen and paper also play a role in determining quality. Colors and background images of paper often cause problems along with the characteristics of scanning devices and algorithms. Figure 2.2(a) shows how texture of the paper affects the scanning process. After scanning, texture becomes severe background noise. This phenomenon is frequently observed in fancy envelopes and personal checks. Intelligent thresholding algorithms can account for simple patterns during the binarization process (Liu & Srihari, 1994). However, complex background images with several gray levels and colors cannot be eliminated by most global binarization process and remain with strokes of the handwritten text (Figure 2.2(c)). Pen tip width and pressure applied



Figure 2.2: Properties of handwritten words: Image factors - (a) noise caused by texture/color of paper, (b) thick pen stroke, (c) background image, and (d) use of poor pen

while writing determines the width of strokes and its variations. Figure 2.2(b) shows an image with relatively thick stroke compared to the image size. Furthermore, the texture of paper and quality of pen can create holes in strokes and blobs on contours (Figure 2.2(d)).

• **binarization** - Programs dealing with handwritten images are usually designed to process binary images (Liu & Srihari, 1994; Otsu, 1979). Intelligent and dynamic binarization is warranted which can account not only for the intensity of strokes but the intensity of background as well. Figure 2.3 shows a gray level image binarized at varying the thresholds. Figure 2.3(b) and Figure 2.3(d) show that binarization plays an important role in determining image quality. Figure 2.3(c) shows the result of binarization using



Figure 2.3: Binarization results for a gray image (with 256 levels) shown in (a) with different threshold values: (b) 80, (c) 132 (chosen by analysis of the entire image), and (d) 180

an algorithm which globally analyzes the intensity distribution of the entire image (Palumbo et al., 1986).

2.1.2 Human factors

Characteristics of handwriting are determined to large extent by the writer. In addition to the style, number of factors are also relevant such as slant, touch-ing/disconnected strokes, and spelling errors.

- **slant** Slant in writing is very common especially in cursive writing, perhaps because most writers learn to write with a slant in early school. We refer to the horizontal inclination of handwriting (from a baseline) as *skew* and the vertical inclination as *slant*. Figure 2.4(a) illustrates positive and negative slant. Figure 2.4(b) illustrates positive and negative skew. Besides the writer, skew can also be introduced from other sources as well. such as, when paper is fed to a scanner.
- **touching/disconnected strokes** Some cursively written characters and most discrete characters are made of more than one stroke, hence the possibility of disconnected strokes. Stroke(s) of a character can be mis-positioned or missing entirely. *i*-dots, *j*-dots, and *t*-bars usually account for the same.



Figure 2.4: Properties of handwritten words: Human factors - (a) slant (vertical) and (b) skew (horizontal)

Furthermore, exaggerated path of strokes cause nearby characters to touch (Figure 2.5(a)). On the other hand, strokes can also be fragmented during binarization resulting in disconnected pieces of stroke(s) (Figure 2.5(b)).

Since character is the basic unit in most recognition process, appropriate grouping and splitting of strokes is important.



Figure 2.5: Properties of handwritten words: Human factors - (a) touching characters and (b) fragment in strokes

• **spelling errors** - This is yet another factor that makes recognition difficult, especially so when paradigms of recognition take advantage of contextual information. Figure 2.6(a) is one such example. Non-standard abbreviations pose a similar problem (Figure 2.6(b)).



Figure 2.6: Properties of handwritten words: Human factors - (a) spelling error and (b) abbreviation

2.1.3 Miscellaneous factors

In most real applications, extraction of isolated handwritten units is non-trivial. Overlapping of guide lines and responses, and interference between writing on consecutive lines add to the complexity.

- **lines (underlines)** Pre-printed forms and envelopes often provide guide lines to eliminate the possibility of skew. However, the binary image does not discriminate between the guide lines and character strokes in terms of their intensity. Algorithms must be developed to remove the guide lines. The elimination process introduces artifacts, such as elimination of certain key strokes, most of the time (Figure 2.7(a)).
- **line separation (leftovers)** Presence of descenders and ascenders in writing create the possibility of overlap between lines. Line separation algorithms


Figure 2.7: Properties of handwritten words: Miscellaneous factors - (a) underline and (b) leftovers from the other fields

are not perfect most of the time (Figure 2.7(b)).

2.2 Background research

2.2.1 Early work

Considerable amount of attention has been paid to the handwritten script recognition problem, starting in the early 1960's. Most of the work was confined to online recognition.

In 1962 Eden, recognized as one of the pioneers in HWR (Eden, 1962), interpreted handwriting as a sequence of basic strokes: bar, hook, arc and loop. He proposed an analysis by synthesis method of recognition using simulation of cursive handwriting using primitive elements. One drawback of this method is the difficulty presented in the interpretation of non-unique sequences.

Mermelstein described a recognition method which identifies strokes in context (Mermelstein & Eden, 1964).

Earnest pointed out the limitation of character based methods. He describes direct extraction of properties from each word and comparison with the corresponding lexicon entries to reduce ambiguities. The recognition process normally yields a small set of words warranting further processing (Earnest, 1962).

Spinrad recognized hand printing using line features (Spinard, 1965).

Enrich and Koehler introduced a contextual postprocessor to determine the correct word from a set of best alternatives for each character generated by a character recognizer (Ehrich & Koehler, 1975).

Sayre describes a recognition process which includes contextual interpretation using bigrams and trigrams (Sayre, 1973).

Farag proposed the use of topological features and context (Farag, 1979).

Burr was mainly concerned with isolated cursive character recognition using warp-based shape matching, and its extension to word recognition by dictionary lookup (Burr, 1983).

Brown described a variety of preprocessing techniques for improving the recognition rate of a cursive script recognition system (Brown & Ganapathy, 1983).

Harmon provides a comprehensive survey of the field up until 1972 (Harmon, 1972).

2.2.2 Segmentation strategies

There have been three strategies in HWR for segmentation: (i) character based, (ii) holistic, and (iii) hybrid.

- **character based** This involves segmenting a word into characters and applying OCR directly to each segment. However, determining meaningful segmentation points without input from character recognition is almost impossible. Depending on the number of stages in segmentation and the extent of OCR use, segmentation strategies can be classified as following:
 - uniform fixed segmentation where a word is chopped at fixed intervals generating a series of segmentation points. Unfortunately, the "magical" interval can not be determined easily. Therefore, one is forced to resort

to over-segmentation and depend on appropriate grouping of segments during recognition.

- *heuristic segmentation* overcomes some of the drawbacks of the first method. Ligature and concavities in a word image are used as basic features to determine the possible segmentation points. Average stroke width and number of characters provide some cues. This method minimizes the number of segmentation points and thus reduces the complexity of search. However, by this method touching characters and certain character pairs which are hard to segment can not be split.
- multi-stage segmentation where recognition module and contextual information (such as statistics derived from the lexicon) are used in early segmentation. Lexicon size can be reduced during the early stages (*e.g.*, based on simple geometric features (Kimura et al., 1993)). Subsequent stages generate more precise segmentation points (Bozinovic & Srihari, 1989).
- **holistic** This approach avoids segmentation algorithm by treating words as a whole. Holistic recognition matches features of the whole image, such as the number of vertical strokes, ascenders, descenders, and loops, with those of words in a lexicon to identify the word. This approach overcomes the difficulties in segmenting a word. Harmon (Harmon, 1962) implemented a holistic handwriting reader. Farag proposed a word level recognition using the directions of stroke sequence as features (Farag, 1979). More recently, Hull (Hull, 1988) and Madhvanath (Madhvanath & Govindaraju, 1993) reported experimental results where holistic methods are used. However, se lecting appropriate holistic features to discriminate between lexicon entries is difficult when the lexicon size is large.

• **hybrid** - The third approach is a compromise of the above two methods. Holistic features are used in the early stages of lexicon reduction. Subsequently, segmentation based recognizers operate on the reduced lexicon. This approach has promise of high performance for large lexicon applications in terms of speed as well as accuracy.

2.2.3 Matching methodologies

Three broad methodologies described in literature are: (i) dynamic programming,(ii) Hidden Markov Models, and (iii) neural networks.

- Dynamic programming (DP) is widely used for solving sequential decision problems. The problem is one of finding the minimum cost, as well as the corresponding sequence of moves (Rabiner & Juang, 1993; Silverman & Morgan, 1990). DP techniques were introduced in the 70s and 80s for non-linear time alignment of speech patterns (dynamic time warping: DTW) (Myers & Rabiner, 1981; Ney, 1984; Sakoe, 1979). DTW can effectively minimize errors occurring during time alignment of two speech sequences. The basic idea of DTW is non-linear stretching and compressing a signal in time. In 1982, Tappert presented the recognition of cursive writing using elastic matching (DTW) using small training and test sets (Tappert, 1982). The recognizer used DTW to find the prototype sequence which best fits the unknown word, yielding the corresponding letter sequence as the estimated letter string. In 1983, Burr introduced DP techniques in isolated character recognition (Burr, 1983). More recently, Kimura used the dynamic matching in offline HWR (Kimura et al., 1993).
- Hidden Markov Models (HMMs) present powerful statistical tools for modeling patterns (Rabiner, 1989; Rabiner & Juang, 1993; Rabiner & Juang,

1986). HMM is a doubly stochastic process producing output through an observation which is another stochastic process. HMM is appropriate for linking character based recognition schemes and high-level knowledge driven schemes.

Several HWR methodologies have used HMM (Agazzi & Kuo, 1993b; Chen, 1993; Chen et al., 1994; Gilloux et al., 1993; Hu et al., 1994).

• **Neural network** techniques are becoming attractive because they can effectively handle high degree of uncertainty and variability by employing automatic learning algorithms. However, they deal mostly with isolated handwritten character recognition. (Hild & Waibel, 1993; Matan et al., 1992) are some exceptions.

2.2.4 Features

Finding optimal features is non-trivial in most pattern recognition areas. Overall cost, data rate, flexibility of representation, and quality of signal, must be taken into account (Kung, 1993).

Many feature analysis techniques have been developed and applied to handwriting recognition (Govindan & Shivaprasad, 1992; Suen et al., 1980):

- features based on global transformation and series expansion Transformations and series expansions, such as Fourier, Walsh, Harr, K-L expansion, etc., make features invariant to some global deformation like translation and rotation. However, such extraction techniques are computationally intensive.
- features based on statistical distribution of points The features include zoning, moments, n-tuples, characteristic loci, crossings and distances. These features are tolerant to distortion and take care of style variations

to some extent. They provide high speed and low complexity for implementation.

• **geometrical and topological features** - These are popular and represent global and local properties of characters. These include strokes, bays in various directions, end points, intersections of line segments, loops, stroke relations, and angular properties. These features have high tolerance to distortion, style variation, translation and rotation.

Combination of the above described features are commonly used in HWR systems. Chen *et al.* (Chen et al., 1995; Chen et al., 1994) use topological features, geometrical features, moment features, pixel distribution, and reference line features for their HMM type word recognizer. Hepp (Hepp, 1994) uses topological features such as hump, hole and cave. In holistic approaches, however, the feature extraction scheme is different. Position and number of ascenders and descenders, loops, and estimated word length (Guillevic & Suen, 1995; Madhvanath & Govindaraju, 1993) are extracted.

2.2.5 Postprocessing

Recognition ambiguities can be resolved by using language models that restrict the set of paths through a hypothesis graph to only those that are compatible with the language model. Language models are built from lexicons (Breuel, 1994a; Chen et al., 1994).

Systems hypothesize all possible words that match entries in the lexicon. In lexicon driven approaches, postprocessing is part of the recognition process it-self (Kim & Govindaraju, 1995; Kimura et al., 1993).

2.2.6 Online word recognition

Advent of compact and powerful computers, and accurate tablets have renewed interest in online handwriting recognition (Tappert et al., 1990). Recent advances combine input and output surfaces providing immediate feedback to writers. Their utility is further emphasized in applications of large-alphabet languages, like Chinese where keyboards can be cumbersome.

An advantage of online handwriting over offline is that temporal or dynamic information of the writing is saved. This information includes the number of strokes, the order of the strokes, ordered sequence of coordinate points, and the speed of the writing for each stroke. The temporal information improves recognition accuracy although it renders the recognition algorithm difficult because it presents a fine granularity not apparent in the offline images.

Due to the similarity between online and speech signals, algorithms developed for recognition of speech have been directly applied to the recognition of online handwriting (Agazzi & Kuo, 1993a; Chen et al., 1994; Hild & Waibel, 1993; Renals et al., 1992; Seni, 1995; Werbos, 1991).

2.3 Related fields

There are a number of related fields which bear direct or indirect relevance to HWR.

2.3.1 Speech recognition

HMMs and neural networks have been widely used in speech recognition. HMM is a stochastic finite state machine, consisting of a set of states and corresponding transitions between states (Rabiner & Juang, 1993). The underlying stochastic

process is not directly observable (i.e. hidden), but still affects the observed sequence of acoustic features. Neural networks typically have a layered feed forward architecture with an input layer, hidden layers and an output layer. Each layer computes a set of liner discriminant functions followed by a nonlinear thresholding function (e.g., sigmoid). However, general neural network architectures, which have been designed for classification of static patterns, are not appropriate for continuous speech recognition. Several network architectures have been developed to overcome the drawback. Recurrent Neural Network (RNN) takes the acoustic input and has a recurrent connection from the output of the state later back to its input (Kung, 1993). The network is trained using back propagation through time. Time-Delay Neural Network (TDNN) is a feed-forward approach in which the loops in RNN are replaced by explicit use of several preceding activation values (Haffner & Waibel, 1992; Waibel et al., 1989). Hybrid of HMMs and neural networks has been proposed to take advantage of each architecture (Bourlard & Wellekens, 1990; Morgan & Scofield, 1991; Morgan & Bourlard, 1995). The neural network estimates emission probabilities for HMMs which deal efficiently with the temporal aspect of speech. Rabiner and Juang (Rabiner & Juang, 1993) provide an introduction to the modern practice of speech recognition.

To apply speech recognition techniques in HWR several issues have to be considered. First how does one produce a consistent sequence of features from the pattern. Online handwriting contains temporal information rendering suitable for direct applying speech recognition techniques. However, for offline handwriting it is not as straight forward. Several attempts have been made to extract temporal information from offline handwriting (D. S. Doermann and A. Rosenfeld, 1995; Govindaraju et al., 1992). However, it is not possible to extract the same extent of information that online handwriting has. Seeking good representation of the sequence of the signal which can be applicable to such systems, for example LPC features in speech signal, is another issue. Various types of feature are defined and applied in many systems.

2.3.2 Optical character recognition

OCR deals with recognition of optically processed characters (Mantas, 1992). In terms of features used, character recognition techniques can be broadly classified into two types (Suen et al., 1980; Wakahara, 1993). *Template matching* is a technique in which an input character is directly compared to a standard set of prototypes stored. This type of technique suffers from sensitivity to noise and is not adaptive to differences in writing style. *Feature analysis and matching techniques* match in feature space. Significant features are extracted from a character and compared to the feature descriptions of ideal characters.

Govindan and Shivaprasad (Govindan & Shivaprasad, 1992), and Mori *et al.* (Mori et al., 1992) trace the development of OCR and offline handwriting recognition systems including description of general processing techniques and commercially available systems. Suen (Suen, 1986) gives overview of the field.

2.3.3 Hardware

Several hardware technologies are related to the development of HWR systems. Computing power of computers is constantly increasing. Development of accurate and fast input and output devices, and communication technologies have renewed interest in HWR. Digitizers and scanners provide sufficient resolution and sampling speeds. Advances in combining digitizers and flat displays bringing input and output onto the same surface have had a positive impact and helped create the notion of *electronic ink*. Electronic ink is the instantaneous display of the trace of the motion of the stylus tip directly under the stylus (Tappert et al., 1990). In addition, advances of data storage, mobile tele-communication, and VLSI technologies directly affect the development of HWR (Essex, 1995). HWR systems need to store large volumes of data and relevant databases. Cost of data storage is constantly reducing and new high density media is being developed (Re-writable CD-ROMs (Bender & Popovic, 1995)).

2.3.4 Psychological aspects

There is no recognition system that can perform as well as humans. Understanding the human vision system and how we recognize objects has been a research issue (Marr, 1982). A trial to improve character recognition performance by analyzing the characteristics of characters in a handwriting written by the same author is an example of applying human factors in HWR (Kawatani, 1995; Kawatani & Miyamoto, 1993).

2.4 Summary

This chapter has described the following two aspects of designing HWR systems: (i) properties of handwritten word images and (ii) methodologies and technologies that can be used in building HWR systems.

HWR systems should be able to handle a wide variety of writing styles and image quality. In addition, recent advances in the computer and communication technologies and related fields are making new demands in many real application fields.

Chapter 3

Image Processing for Handwriting Recognition

Implementation of a fast system for handling handwritten script is a challenging task which is in increasing demand for real applications (Govindaraju et al., 1994: Srihari et al., 1995a). The common design issues in developing a fast system are: (i) implementation of efficient image processing routines (Braun et al., 1995; Brown & Ganapathy, 1983) and (ii) developing appropriate matching algorithms (Chen et al., 1994; Favata, 1993; Hild & Waibel, 1993). We present our image handling methodology in this chapter and the matching methodology in Chapter 4.

Conventional image processing techniques deal with two-dimensional images. Significant processing time is spent in image handling of objects using straight forward two-dimensional representations. Image representation methods have been developed to minimize the amount of data (without loss of information) and to increase processing speed. Examples of such representation techniques are chain code (Freeman, 1961), line adjacency graphs (LAG) (Pavlidis, 1978; Pavlidis, 1982), and symmetric axis transformation (SAT) (Blum & Nagel, 1978). However,

29

introducing new representation techniques requires modification of existing image handling procedures.

In this thesis, we use chain code based image representation. We describe suitable image handling procedures, such as preprocessing, segmentation, and feature extraction that can apply to this representation. Preprocessing deals with noise removal, slant correction and smoothing of contours. Noise introduced during data capture is eliminated by comparing the size of connected components with an estimate of average stroke width of the word image. Slant angle is estimated by averaging orientation angles of vertical strokes and shifting the x-coordinates of components accordingly. Smoothing eliminates jaggedness of contours.

In our recognition engine described in Chapter 4, there are two major functional modules which contain image handling procedures: segmentation and feature extraction. Detailed explanations on the image handling procedures of those two modules are described in Section 4.1 and Section 4.2 respectively. Focus of this chapter is on the preprocessing modules.

3.1 Image representation

Connected component labeling is a grouping process that joins pixels into a region. In general, HWR applications deal with connected components instead of pixels to minimize operations on large sized image arrays. Several methods have been adopted for the representation of connected components to optimize computer resources and processing time. *Line adjacency graph* (LAG) is a technique employing run length encoding to represent objects (Pavlidis, 1982). A run length encoding of a binary image is a list of horizontal runs of foreground pixels. For each run, the location of the starting pixel and either its length or the location of the last pixel are saved. The nodes of such a graph correspond to intervals.



Branches connect nodes if corresponding intervals are on adjacent lines their projections along the scan direction overlap, and their pixels have the same intensity. *Symmetric axis transformation* (SAT), also called medial axis, is a representation that goes orthogonal to the boundary, while the conventional description traverses the boundary (Blum & Nagel, 1978). *Chain code* is a linear structure that results from quantization of the trajectory traced by the centers of adjacent boundary elements in an image array (Freeman, 1974; Lu & Dunham, 1991). Each data node in the structure represents one of eight grid nodes that surround the previous data node. This is a compact representation allowing efficient processing. Figure 3.1(a) shows the slope convention used in our implementation of the chain code representation.

The basic unit of chain code data is an 8-byte structure defined as follows:

```
structure COMPONENT
{
    short Y,X;
    unsigned char status, mode, slope, curv;
```

For efficient manipulation, certain properties of the chain code contour are stored in an array (Figure 3.1(b)). *Data* field in the array contains positional and slope information of each component of the traced contour. Properties stored in the *information* fields are: coordinates of bounding box of a contour, number of components in the corresponding data fields, area of the closed contour, and a flag which indicates whether the contour is interior or exterior. Chain code representation of an entire image consists of n such contour arrays cascaded, where n is the number of closed contours. Properties of a contour are defined by assigning different numbers in the status and mode fields of the structure, COMPONENT. Selection of a contour is achieved by searching the information field.

3.2 Chain code based image processing

Preprocessing procedures based on the chain code representation for handwritten words are described in this section.

3.2.1 Noise removal

Scanning devices and transmission media introduce noise during the image capture process. Noise can be easily identified by computing the size of connected components. Single pixel components are the first removed during the chain code conversion procedure. Second, components which have a small size with respect to the average stroke width are removed. Width and height of components are used for identifying such components. Figure 3.2(a) shows an example in which small components (compared to the average stroke width) as well as single pixel components which are eliminated. It should be noted that *i-dot*s are not eliminated because we use the average stroke width to identify noise.



Figure 3.2: Noise handling: (a) small components including single pixel components are eliminated, (b) an image rejected due to severe background noise

Quality of writing devices and texture of paper often introduce fragments that could be regarded as noise. Some images have severe background noise due to the color of paper and artifacts from scanning device which cannot be eliminated easily. We reject such images in the early stages of processing based on the number of connected components in an image, and number of components on the exterior contour and size of the connected components. Figure 3.2(b) shows an image having severe background noise and the corresponding contour representation.

3.2.2 Slant angle estimation and correction

Handwritten words often have a slant which needs to be corrected so that segmentation and recognition can proceed. Vertical line elements from contours are extracted by tracing chain code components using a pair of one dimensional filters. Each filter is an eight element array of different weights used for detecting vertical lines having opposite directions. A convolution operation between the filter and five consecutive components is applied iteratively by sliding the filter one component at a time. Coordinates of the start and end point of each line element extracted provides the slant angle. Global slant angle is the average of all angles, weighted by their length in the vertical direction because the longer line gives more accurate angle than the shorter one. Figure 3.3(b) and Figure 3.3(c) show the vertical lines collected from an example image shown in Figure 3.3(a) using Algorithm 3.1.

Algorithm 3.1.

```
v\_up[8] = \{2, 1, 0, 1, 2, 3, 4, 3\}
v\_down[8] = \{2, 3, 4, 3, 2, 1, 0, 1\}
for i = strat\_component to end\_component step 1
sum\_slope\_up = 0;
sum\_slope\_down = 0;
for j = i - 2 to i + 2 step 1
sum\_slope\_up = sum\_slope\_up + v\_up[ j.slope ];
sum\_slope\_down = sum\_slope\_down + v\_down[ j.slope ];
end
if sum\_slope\_up > THRESH_V then i.status = MARK_V_UP fi
if sum\_slope\_down > THRESH_V then i.status = MARK_V_DOWN fi
end
```

Tangent of the estimated slant angle (θ) is used to correct for the slant. Eq. (3.1) gives the adjustment made to the *x* and *y* coordinates.

$$x_n = x - y \tan(\theta)$$

$$y_n = y$$
(3.1)







(d)





Figure 3.3: Noise handling: (a) contour representation of an image, extraction of upward (b) and downward (c) vertical lines, (d) correction of x-coordinates based on the estimated angle, (e) connecting broken points. (f) smoothing contours where x_n and y_n are coordinates after correction. Quantization error incurred during the correction process introduces "*jigs*" and broken points in the horizontal direction (Figure 3.3(d)). An interpolation method is applied to connect the broken points (Figure 3.3(e)). The slant normalization algorithm assumes that the slant angle is between -45° and 45°. A smoothing algorithm (Section 3.2.3) is applied to the contours which yields the final result shown in Figure 3.3(f).

3.2.3 Smoothing contours

Smoothing eliminates small blobs on a contour and illegal combinations of chain code slopes introduced during slant correction (Section 3.2.2). A sliding three-component window is applied to all components. All combinations of three adjacent chain code slopes are analyzed and classified into seven types based on the treatment warranted.

- type 0: needs no correction.
- type 1: remove 1st component and adjust slope of 2nd component (Figure 3.4(a)).
- type 2: remove 2nd component and adjust slope of 3rd component (Figure 3.4(b)).
- type 3: remove 2nd and 3rd components (Figure 3.4(c)).
- type 4: remove 1st and 2nd components and adjust slope of 3rd component (Figure 3.4(d)).
- type 5: remove all three components (Figure 3.4(e)).
- type 6: adjust slopes of 2nd and 3rd component (Figure 3.4(f)).

Depending on the type, components can be removed from the chain code structure, and coordinates and slopes of components can be updated to maintain the



Figure 3.4: (a - f) Types of smoothing operations - for each type, left sequence turns into right sequence after the smoothing operation (× represents elimination of the components and () represents starting point) and (h) shows how the algorithm is working for the circled area of the example image of (g)

p	C				I	1				p	C	<u> </u>]	1			
		0	1	2	3	4	5	6	7			0	1	2	3	4	5	6	7
0	0	0	0	2	4	3	2	0	0	1	0	0	0	2	2	3	0	0	$\overline{\mathbf{a}}$
	1	0	0	0	0	4	3	2	6	H	1	0	0	0	$\overline{6}$	$\frac{1}{2}$	3	0	$\frac{3}{6}$
	2	1	1	1	1	4	5	3	2	Ħ	2	0	0	0	0	2	4	3	$\frac{1}{2}$
	3	1	1	1	1	4	4	5	4	Ī	3	2	6	0	0	0	4	4	3
ĺ	4	4	4	4	4	4	4	4	4		4	3	1	1	1	4	4	5	4
	5	2	3	5	4	1	1	1	1		5	4	4	4	4	4	4	4	4
	6	2	2	3	5	0	0	0	0		6	1	2	3	4	5	1	1	1
	7	0	6	2	3	0	6	0	0		7	0	6	2	3	4	4	0	0
2	0	0	0	2	2	3	5	0	0	3	0	1	1	1	2	3	4	5	4
	1	0	0	0	6	2	3	2	6		1	0	0	0	6	2	3	4	4
	2	0	0	0	0	2	4	3	2]	2	0	0	0	0	2	2	3	2
	3	2	6	0	0	0	0	4	3		3	0	6	0	0	0	6	2	3
	4	3	2	1	1	1	1	4	5		4	3	2	0	0	0	0	2	4
	5	5	3	4	1	1	1	4	4		5	4	3	2	6	0	0	0	4
	6	4	4	4	4	4	4	4	4		6	5	4	3	2	1	1	4	4
	7	1	1	4	3	5	4	1	1		7	4	4	4	4	4	4	4	4
4	0	4	4	4	4	4	4	4	4	5	0	4	4	5	4	3	1	1	1
	1	1	1	1	1	4	3	5	4]	1	4	4	4	4	4	4	4	4
	2	0	0	0	0	2	2	3	5]	2	5	4	1	1	1	2	3	4
	3	2	6	0	0	0	6	2	3]	3	4	4	0	0	0	6	2	3
	4	3	2	0	0	0	0	2	4		4	3	4	0	0	0	0	2	2
	5	4	3	2	6	0	0	0	0		5	2	3	2	6	0	0	0	6
	6	4	5	3	2	1	1	1	1]	6	2	4	3	2	0	0	0	0
	7	4	4	5	3	4	1	1	1		7	0	4	4	3	2	6	0	0
6	0	1	1	4	5	3	2	1	1	7	0	0	0	2	4	3	2	0	0
	1	1	1	4	4	5	3	4	1]	1	0	0	0	4	4	3	2	6
	2	4	4	4	4	4	4	4	4		2	1	1	4	4	5	4	3	1
	3	5	4	1	1	1	1	4	3		3	4	4	4	4	4	4	4	4
	4	3	5	0	0	0	0	2	2		4	3	4	5	4	4	1	1	2
	5	2	3	2	6	0	0	0	6		5	2	3	4	4	0	0	0	6
	6	2	4	3	2	0	0	0	0		6	2	2	3	4	0	0	0	0
	7	0	0	4	3	2	6	0	0		7	0	6	2	3	2	6	0	0

Table 3.1: Smoothing operation type table: ${\bf p}$ - slope of the preceding component, ${\bf c}$ - slope of the current component, and ${\bf n}$ - slope of the next component

relation between adjacent pixels. The subsequent smoothing operation is performed based on updated slopes. If any component under testing is eliminated from the structure the start and end points of the window are also updated. For efficient classification a three-dimensional table which stores types and offset is used. Table 3.1 shows list of types for all combinations of slopes of three consecutive components. In the table, **p**, **c**, and **n** represent slopes of the previous, current, and next components. For example, if **p** = 2, **c** = 4, and **n** = 6, the table returns 4 as the type.

Table 3.2 describes steps in the smoothing process for the upper contour of the circled area of Figure 3.4(g). Figure 3.4(h) shows components and the corresponding slopes of the area before (left) and after (right) applying the algorithm.

slopes	type	action
77057576	0	no modification is needed
7[705]7576	2	2nd component is removed and slope of 3rd compo- nent is updated to 6
7[767]576	0	no modification is needed (type classification is per- formed using the updated component and the follow- ing component)
7767576	6	slopes of 2nd and 3rd components need to be up- dated to 6
7766676	0	no modification is needed. (type classification is per- formed based on the updated slopes)
7766676	0	no modification is needed

Table 3.2: Steps of smoothing procedures for Figure 3.4(g): box represents window for classification

3.2.4 Stroke width estimation

Stroke width can vary locally depending on writing devices and pressure within a script. It is reasonable, therefore, to talk of *average* stroke width which can be used for other image processing procedures in an adaptive manner. For example, size of small objects that could be regarded as noise is determined based on this parameter. Stroke width is also used for locating ligatures in the segmentation procedure (Section 4.1). To estimate average stroke width, chain code contours of a word image are divided horizontally. By tracing contours from the left-most point to the right-most point, following y-distances are computed for each x-coordinate as shown in Figure 3.5 (b): (i) distance between upper and lower trace of the outer contour, (ii) distance between the upper trace of the inner contour and upper trace of the outer contour, and (iii) distance between lower trace of the inner contour and lower trace of the outer contour.

Based on the number of occurrences of each distance value for the image shown in Figure 3.5(a) the histogram of Figure 3.5(b) is obtained. The peak of the histogram gives a good estimate of average stroke width (Eq. (3.2)).

$$N_d$$
 = number of occurrences of distance value d
 $P_d = \arg \max_d (N_d)$
(3.2)

To account for the fluctuation in pen movement, Algorithm 3.2 is used to estimate the actual stroke width in the local area.

Apacuse

(a)



(b)

Figure 3.5: Estimating average stroke width: (a) a sample word image (b) measuring y-distances for each x-coordinate and analyzing occurrences of stroke widths of image (a)

Functional	Tim	e consum	Description	
Module	10	100	1000	(section number)
Pre-processing	10	10	11	3.2.1, 3.2.2, 3.2.3
Segmentation	15	15	15	3.2.4, 4.1
Feature extraction	34	34	37	4.2
Sub-total	59	59	63	
(image/total)	(28.9%)	(15.4%)	(9.1%)	
Matching	145	324	633	
Total	204	383	696	

Table 3.3: Timing analysis (*note*: feature extraction time includes time taken for combining segments before feature extraction)

Algorithm 3.2.

 $d = P_d;$ **while** ($d < MAX_STROKE_WIDTH \land N_d > P_d/2$) d = d+1; **end** $actual_stroke_width = d+1;$

3.3 Experiments and Results

To evaluate the impact of the chain code based preprocessing procedures described in this chapter, we used a test set of 3,000 word images, including city names, firm names, personal names, street names, and state names, collected from live mail pieces digitized at 212dpi. Average size of the images was 217×82 .

	with I	ore-pro	cessing	w/o pre-processing				
	10	100	1000	10	100	1000		
top 1	96.8	88.5	74.0	91.1	75.8	58.9		
top 2	98.6	93.4	83.9	96.1	84.4	67.1		
top 20		99.0			97.2			
top 50			98.6			93.1		

Table 3.4: Impact of the preprocessing (in %)

Experiments focused on two aspects: (i) efficiency of the image handling modules in terms of speed and (ii) recognition performance.

We have categorized the image handling modules into three groups as shown in Table 3.3 to show the processing times.

Table 3.4 shows the impact of our preprocessing modules by comparing the recognition performance of systems with and without the preprocessing modules described. It should be noted that performance degradation is significant particularly for systems with large lexicons. Performance drops of 5.8%, 14.4%, and 20.4% were observed for lexicon sizes of 10, 100, and 1000, respectively.

3.4 Summary

Implementation of efficient and robust image handling modules is a critical in real time handwriting recognition systems. We have described chain code based image representation and implemented the image handling procedures using this representation in our HWR. Experimental results indicate that the implementation is very efficient in terms of speed while the performance improvement is encouraging.

Chapter 4

Handwritten Word Recognition Engine

A handwritten word recognition engine is developed to recognize isolated unconstrained words. Figure 4.1 illustrates the methodology that we describe in this thesis. There are two processing phases: *training* and *recognition*. Input images go through the steps of chain code generation, preprocessing, segmentation and feature extraction in both phases.

- 1. **Chain code generation** step converts the binary image input into a chain code representation by coding the boundary contours of components in the image while preserving the positional and directional information of adjacent pixels (Freeman, 1961). An array is defined for efficient representation and manipulation of data. Single pixel components (speckle noise) are detected and removed. Subsequent image handling steps work with chain code data.
- Preprocessing step includes noise removal, slant correction and smoothing (Brown & Ganapathy, 1983; Schürmann et al., 1992). The procedures are described in detail in Chapter 3.



Figure 4.1: Word recognition methodology

- 3. **Segmentation** step returns the segmentation points to be used for grouping one or more segment(s) to form meaningful characters. The segmentation points are determined using a combination of ligatures and concavity features on the contour. Average stroke width of an image is estimated and used in an adaptive fashion to determine the features. The number of segmentation points is kept to a minimum while ensuring that a segmentation point exists to split joint characters and the maximum number of segmentation points per character is four.
- 4. Feature extraction step generates feature vectors for combination of segments that are hypothesized to be characters. Global and local features are defined and extracted from 3×3 sub-images of the segment(s) represented by their chain code.

- 5. **Training phase** uses a training set consisting of binary word images different from the test set. The reference of character segmentation points in the input image are determined manually. Segmentation statistics such as how often a particular character is split into how many segments (1 ... 4) is stored. Feature vectors are extracted for the manually referenced character segments and used by a clustering algorithm to find cluster centroids of characters to be used in recognition. Additional details of the training procedure are described in Appendix A.
- 6. **Recognition step** uses segmentation statistics, character cluster centroids, a dictionary, and feature vectors derived from the test image. A dynamic matching scheme is used to compare features of a segment or a combination of consecutive segments with the cluster of centroids of a character in a lexicon entry. This procedure is used to rank the lexicon entries.

Figure 4.2 illustrates the recognition procedure. A sample word image, "word", is split by the segmentation program (Figure 4.2(a)). The example illustrates the matching scheme between the image of Figure 4.2(a) and a lexicon entry "word". In Figure 4.2(c), arcs between adjacent nodes represent the possibility of grouping the image segments (Figure 4.2(a)) between those segmentation points as a character hypothesis. For example, the arc from node 2 to node 5 represents the confidence of match between cluster centroids of "o" (from training) and the feature vector extracted from the point 2 to 5 in the image. Based on the matching scores the best path, represented by solid arcs, is found using dynamic programming. The confidence of the match is scored by dividing the accumulated cost along the path by the number of characters of the lexicon entry in consideration.



lowest cost path = 5.0 + 6.1 + 3.8 + 4.9 = 19.8confidence of match = 19.8/4 = 4.95

(b)



(c)

Figure 4.2: Matching : (a) segmentation points, (b) confidence of the match, and (c) matching paths and confidence

4.1 Segmentation

A segmentation algorithm should be general and robust to handle various styles of writings and thickness of strokes. We make the following assumptions: (i) the number of segments per character must be at most 4, and (ii) all touching characters should be separated.

4.1.1 Segmentation procedure

Handwriting is generated by the movement of a pen from left to right along an axis in the horizontal direction, giving rise to what are called *ligatures* (Simon, 1992). The ligatures are strong candidates for segmentation points in cursive scripts. Ligatures are extracted as follows. A word image is split horizontally (Figure 4.3(c)). If the distance between *y*-coordinates of the upper-half and lower-half of the outer contour for a *x*-coordinate is less than or equal to the average stroke width, then the *x*-coordinate is marked as an element of a ligature. This procedure is repeated for subsequent *x*-coordinates. Ligatures in the extremities are eliminated to reduce the number of potential segmentation points. Figure 4.3(d) shows ligatures obtained for an example word.

On the other hand, segmentation points between discretely written (as opposed to cursive) touching characters cannot be hypothesized by ligatures alone (for example, the touching of "S" and "y" in Figure 4.3). Alternatively, concavity features in the upper contour and convexities in the lower contour are used in conjunction with ligatures (Figure 4.3(e)). Heuristics are applied to reduce the number of potential segmentation points. For example, if ligatures and concavities overlap, concavities are ignored. The final segmentation points are stored with the corresponding segment information: the indices of related contours and coordinates of the bounding box of the segment.

(a)



(f)

Figure 4.3: Segmentation: (a) original test image (b) chain code representation after slant normalization (c) splitting upper- and lower- contours (d) ligatures (e) concavities (f) segmentation points

4.1.2 Segmentation statistics

Segmentation statistics of each character represent the possible ways in which a training character image can be split into segments (Chen et al., 1994).

Eq. 4.1 gives the computational form for computing the duration probability (dur(j,i)).

$$dur(j,i) = Pr(j | T_i)$$

$$= \frac{\text{no. of times that training character } (T_i) \text{ is segmented into } j \text{ segments}}{\text{no. of training characters of } (T_i)}$$

$$(4.1)$$

where $T_i = 0, 1, \dots, 9, A(a), B(b), \dots, Z(z), i = 1, 2, \dots, 36$, and j = 1, 2, 3, 4.

Our assumption allows at most 4 segments per character during the segmentation phase. However, statistics reveal that the upper limit of 4 segments is valid for only a few characters (such as "*m*" and "*w*"), and that for most characters it is less than 4. This information is used to advantage in speeding up the matching step of the recognition phase and improving recognition accuracy. Confusions are reduced by matching within the window size controlled by the statistics, number of characters in a lexicon entry, and number of segments of the word image. Table 4.1 shows the segmentation statistics. For instance, it can be seen that 82% of "*O*" are not split while 18% are split into two segments.

4.2 Features

We have designed a feature extractor that converts chain code images into feature vectors in a simple and fast manner. The same feature extraction procedure is used in both the training and recognition phase as shown in Figure 4.1. In the training phase, extracted features are provided to a clustering procedure so that patterns of similar shapes can be represented by a codeword in the codebook

	number of segments				
character	1	2	3	4	
0	0.81646	0.17655	0.00662	0.00000	
1	0.99882	0.00118	0.00000	0.00000	
2	0.36199	0.58329	0.05080	0.00390	
3	0.62680	0.35355	0.01848	0.00115	
4	0.09486	0.77207	0.12768	0.00477	
5	0.44688	0.43248	0.11035	0.00891	
6	0.07897	0.90418	0.01684	0.00000	
7	0.63985	0.32673	0.03341	0.00000	
8	0.75980	0.23124	0.00825	0.00068	
9	0.75553	0.23855	0.00517	0.00073	
a	0.50623	0.44725	0.04574	0.00075	
b	0.36559	0.44802	0.16845	0.01792	
с	0.83636	0.14965	0.01398	0.00000	
đ	0.35555	0.52888	0.10666	0.00888	
e	0.79811	0.17643	0.02378	0.00165	
f	0.61696	0.32432	0.05312	0.00559	
g	0.39736	0.50789	0.08157	0.01315	
h	0.11460	0.68576	0.16081	0.03881	
i	0.94677	0.05015	0.00307	0.00000	
j	0.54902	0.31372	0.13725	0.00000	
k	0.18085	0.62234	0.17819	0.01861	
1	0.86100	0.12050	0.01800	0.00050	
m	0.00000	0.08317	0.70406	0.21276	
n	0.03304	0.75373	0.20456	0.00865	
0	0.77521	0.20979	0.01498	0.00000	
р	0.42901	0.50925	0.05864	0.00308	
Q	0.55555	0.22222	0.22222	0.00000	
r	0.64878	0.31138	0.03577	0.00406	
s	0.66912	0.28571	0.04331	0.00184	
t	0.69116	0.27010	0.03078	0.00794	
u	0.01084	0.95664	0.03116	0.00135	
v	0.05166	0.89298	0.05535	0.00000	
w	0.00000	0.02443	0.86278	0.11278	
x	0.18650	0.61904	0.18650	0.00793	
У	0.08144	0.78733	0.13122	0.00000	
z	0.44444	0.44444	0.11111	0.00000	

	Table 4.1:	Segmentation	statistics
--	------------	--------------	------------

(Duda & Hart, 1973). Segmentation statistics are obtained during the training phase. In the recognition phase, features of segment(s) of a test image are compared to the codewords to find the best match.

Chain code features: 74 chain code based features are used. Two global features: aspect ratio and stroke ratio of the entire template are used (Eq. 4.2). Segment(s) are divided into 3×3 sub-images and local features are collected from each sub-image. Distribution of the 8 directional slopes for each sub-image form the 72 local feature vectors ($8 \times 3 \times 3$) as shown in Eq. 4.3.

global features

$$F_{g_i} = sigmoid(\frac{H_i - V_i}{V_i}) \quad \text{for } i = 1,2$$
(4.2)

where

 $H_1 = X_{max} - X_{min},$ $V_1 = Y_{max} - Y_{min}$ for aspect ratio $H_2 = N_{horizontal_stroke},$ $V_2 = N_{vertical_stroke}$ for stroke ratio

local features

$$F_{l_{ij}} = \frac{s_{ij}}{N_i S_j}$$
 for $i = 1, 2, \dots 9$ and $j = 0, 1, \dots 7$ (4.3)

where

 $s_{ij} = \text{number of components with slope } j \text{ from sub-image } i$ $N_i = \text{number of components from sub-image } i$ $S_j = \max_i \left(\frac{s_{ij}}{N_i}\right)$

4.3 Recognition

The objective is to find the best match between a word in the lexicon and the image. Depending on the metric used for classification and the method of lexicon handling, various matching schemes are possible. A dynamic matching approach which has been extensively studied by the speech recognition community (Myers

& Rabiner, 1981; Ney, 1992; Sakoe, 1979) is used. Instead of passing on combination of segments to a generic OCR, the lexicon is brought into play early in the process. A combination of adjacent segments (up to a maximum of 4 based on our segmentation criteria) are compared to only those character choices which are possible at the position in the lexicon entry being considered.

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 the matching. For each match, the minimum distance value of each comparison is retained. In the second phase of the match, a global optimum path is obtained by using dynamic programming based on the saved minimum distances of the first matching phase.

4.3.1 First stage matching

Let an input word image be denoted as

$$T = \{s_0, s_1, \cdots, s_{N_r-1}\}$$
(4.4)

in which N_s is the number of segments in the input image and s_k represents the k-th segment.

Let the lexicon entries be referred to as R_i .

$$R_i = \{c_i(0), c_i(1), \cdots, c_i(N_c(i) - 1)\} \quad 1 \le i \le N_i$$

$$(4.5)$$

where $N_c(i)$ is the number of characters in the *i*-th lexicon entry, N_l is the number of lexicon entries, and $c_i(j)$ represents the *j*-th character of the *i*-th lexicon entry.

Let us define the notion of merged segments as follows:

$$S(b,e) = \{s_b \oplus s_{b+1} \oplus \dots \oplus s_e\}$$

= $\{s_\nu\}_{\nu=b}^e$ for $0 \le b < N_s, 0 \le e < N_s, \text{and } b \le e$ (4.6)

where *b* is a beginning segment, *e* is the ending segment, and \oplus is the merging operator (Figure 4.4).

From the training data of clusters, codewords of all characters, C_m , can be obtained:

$$\mathbf{C}_{m} = \{\mathbf{x}_{m}(i)\}_{i=1}^{N_{m}}$$
(4.7)

where $m = 0, \dots, 9, a, \dots, z, N_m$ is number of codewords of character *m*, and $\mathbf{x}_m(i)$ represents the *i*-th codeword of character *m*.

Therefore, the basic problem of matching can be described as follows: For a given input word image, T, find a lexicon entry, R_i , which has the minimum distance from T:

$$D^{\star} = \min_{1 \le i \le N_l} D(R_i, T)$$

$$R^{\star} = \arg \min_{R} D(R_i, T)$$
(4.8)

First, matching between each character of a given lexicon entry ($c_i(j)$, $0 \le j < N_c(i)$) and an arbitrary portion of the segment(s), specified by a beginning test segment, *b*, and an ending test segment, *e*, of the test input *T* is performed.

$$\hat{D}(v, b, e) = d(\mathbf{F}(S(b, e)), \mathbf{x}_{c_i(j)}(v)) \quad \text{for } 0 \le v < N_{c_i(j)}$$
(4.9)

where $d(\cdot, \cdot)$ is distance between two feature vector and $\mathbf{F}(\cdot)$ represents feature extraction operation for segment(s). Eq. 4.9 gives the distance between each codeword of character $c_i(j)$ and all possible pairs of test segments starting at *b* and ending at *e*. The best match for all *v* is obtained by Eq. 4.10.

$$\bar{D}(b,e) = \min_{1 \le v \le N_{c_i(J)}} [\hat{D}(v,b,e)]$$
(4.10)
SAMALAL

(a)

SOLJADLELLOL *s*₁ *s*₂ *s*₃ *s*₄ *s*₅ *s*₆ *s*₇ *s*₈ *s*₉ *s*₁₀ *s*₁₁ *s*₁₂ s_0 $S(0,3) = s_0 \oplus s_1 \oplus s_2 \oplus s_3$ $S(0,2) = s_0 \oplus s_1 \oplus s_2$ $S(0,1) = s_0 \oplus s_1$

 $S(1,2) = s_1 \oplus s_2$

 $S(1,3) = s_1 \oplus s_2 \oplus s_3$



 $S(1,4) = s_1 \oplus s_2 \oplus s_3 \oplus s_4$

(b)

Splitting and merging of segments: (a) segmentation result (b) ex-Figure 4.4: amples of splitting and merging, and corresponding notations

For meaningful alignment of segment(s) against a character of a lexicon entry and minimization of computational complexity, we apply constraints on the segmentation criteria, number of characters in a lexicon entry, and number of segments in a test image. These factors together determine the size of the permissible matching window. Accordingly, the range of ending segments for $c_i(j)$ id decided as follows.

$$e_{max}^{j} = \min(N_{s} - (N_{c}(i) - j), (j+1) \cdot M_{d} - 1)$$

$$e_{min}^{j} = \max(N_{s} - (N_{c}(i) - (j+1)) \cdot M_{d} - 1, j) \quad \text{for } j = 1, \cdots, N_{c}(i) - 1$$
(4.11)

where $1 \le d \le M_d$ for each *j* and M_d is the maximum number of segments per character, *i.e.* 4.

The range is determined according to our segmentation assumptions: a character in R_i consumes at least one segment and a character can be segmented as many as M_d segments. The interval between e_{min}^j and e_{max}^j is used to determine the matching window size for the reference character $c_i(j)$. For a given range of ending segments, a starting segment, b, is determined based on the notion that bshould be bigger than e_{min} of the previous character and at most M_d segments can be combined for the match. Table 4.2 shows the computation results of D(b, e) in Eq. 4.10 for the example word image in Figure 4.4. The horizontal axis represents the ending segment number and the vertical axis represents the j-th character of a lexicon entry. Four rows for each character are shown representing the duration of segments, 1 to 4, from bottom to top. For example, a component in third column (ending segment number of 2) and second row (duration of 2) of second character (j = 1, "y" in this particular example) represents $\overline{D}(1,2) = 12.41$. Similarly, the component with ending segment number of 5, j = 2("r") and duration of 4 represents $\overline{D}(2,5) = 11.59$. Eq. 4.9 and Eq. 4.10 are computed within the matching window, thus significantly reducing the computational complexity.

					endi	ng s	segmen	-					
i_	_d0	1_	2	3	4	5	66	7	- 8	9	10	11	12
() (s)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 5.10 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 7.33\\ 0.00\end{array}$	$0.00 \\ 8.15 \\ 0.00 \\ 0.00$	$7.72 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
(y)	$\begin{array}{c} 4 & 000 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$0.00 \\ 0.00 \\ 0.00 \\ 10.33$	$0.00 \\ 0.00 \\ 12.41 \\ 10.05$	$0.00 \\ 10.99 \\ 4.75 \\ 7.51$	9.73 10.07 10.69 10.73	$11.59 \\ 9.64 \\ 10.22 \\ 0.00$	$9.67 \\ 9.30 \\ 0.00 \\ 0.00$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
(r)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$0.00 \\ 0.00 \\ 0.00 \\ 4.07$	$\begin{array}{c} 0.00 \\ 0.00 \\ 9.29 \\ 6.01 \end{array}$	$0.00 \\ 11.41 \\ 9.28 \\ 4.96$	$10.48 \\ 9.39 \\ 7.56 \\ 6.72$	8.52 8.08 7.40 3.82	$10.26 \\ 6.94 \\ 7.08 \\ 5.62$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
3 (a)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$0.00 \\ $	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 6.60\end{array}$	$0.00 \\ 0.00 \\ 8.22 \\ 7.17$	0.00 8.41 7.28 4.55	$7.55 \\ 5.41 \\ 4.65 \\ 4.49$	$7.54 \\ 5.78 \\ 7.10 \\ 5.03$	6.43 6.74 5.33 3.72	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
4 (c)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$0.00 \\ 0.00 \\ 0.00 \\ 7.63$	$\begin{array}{c} 0.00 \\ 0.00 \\ 6.44 \\ 4.91 \end{array}$	$ \begin{array}{c} 0.00 \\ 7.11 \\ 8.26 \\ 3.88 \end{array} $	$10.07 \\ 9.20 \\ 8.82 \\ 2.75$	$ \begin{array}{r} 6.43 \\ 9.45 \\ 6.89 \\ 2.71 \\ \end{array} $	$6.77 \\ 11.12 \\ 5.83 \\ 3.34$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
5 (u)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$ \begin{array}{c} 0.00 \\ 0.00 \\ 0.00 \\ 6.86 \end{array} $	$ \begin{array}{c} 0.00 \\ 0.00 \\ 4.82 \\ 6.26 \end{array} $	0.00 5.90 5.56 5.85	$\begin{array}{c} 6.68 \\ 7.32 \\ 4.91 \\ 5.89 \end{array}$	5.51 7.24 2.91 6.39	6.77 5.62 5.85 7.36	0.00 0.00 0.00 0.00	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
(S)	$\begin{array}{c} 4 & 0.00\\ 3 & 0.00\\ 2 & 0.00\\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$0.00 \\ 8.55 \\ 5.18 \\ 3.89$	5.57 8.68 5.24 5.24	8.39 9.70 5.39 6.98	11.37 8.17 3.63 4.27	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$
7 (ㅋ)	$\begin{array}{c} 4 & 0.00 \\ 3 & 0.00 \\ 2 & 0.00 \\ 1 & 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	0.00 0.00 0.00 0.00	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	10.30 6.97 7.95 3.31

Table 4.2: First stage of matching procedure

4.3.2 Second stage matching

In the second stage, the individual matching scores of the first stage are combined to compute the accumulated cost (distance) over the entire lexicon entry. This is accomplished by using dynamic programming. Distance of the best path ending at segment e for each j-th character is computed by Eq. 4.12.

$$\bar{D}_{j}(e) = \min_{1 \le b \le e} [\bar{D}(b, e) + \bar{D}_{j-1}(b-1)]$$
(4.12)

Based on the recursion of Eq. 4.12, we can formulate the second stage of this dynamic matching procedure for determining the overall best path as follows.

• Step 1 : Initialization

$$\bar{D}_0(e) = \bar{D}(0,e) \quad 0 \le e \le M_d - 1$$
(4.13)

					endir	ng s	segment	5					
i	0	1	2	3	4	5	6	7		2	10	11	12
$\begin{array}{c} 0(s) & 5\\ 1(y) & 0\\ 2(r) & 0\\ 3(a) & 0\\ 4(c) & 0\\ 5(u) & 0\\ 6(s) & 0\\ 7(e) & 0 \end{array}$.10 .00 .00 .00 .00 .00	$\begin{array}{c} 7.33 \\ 15.43 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \end{array}$	$ \begin{array}{r} 8.15\\ 17.38\\ 19.50\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00 \end{array} $	$\begin{array}{c} 7.72 \\ 12.07 \\ 23.39 \\ 26.10 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 14.83\\ 17.03\\ 27.72\\ 33.73\\ 0.00\\ 0.00\\ 0.00\\ 0.00\end{array}$	$\begin{array}{c} 0.00\\ 17.79\\ 19.63\\ 21.58\\ 32.54\\ 40.59\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 17.02\\ 20.16\\ 21.63\\ 25.45\\ 38.56\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 21.77\\ 22.81\\ 24.43\\ 31.30\\ 0.00\\ 0.00\\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 23.46\\ 25.52\\ 30.32\\ 35.19\\ 0.00 \end{array}$	$\begin{array}{c} 0.00 \\ 0.00 \\ 0.00 \\ 26.80 \\ 27.34 \\ 35.56 \\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 30.05\\ 34.33\\ 0.00\\ 0.00\\ \end{array}$	$\begin{array}{c} 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 0.00 \\ 30.97 \\ 0.00 \end{array}$	$\begin{array}{c} 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 0.00\\ 34.28 \end{array}$

Table 4.3: Second stage of matching procedure

• Step 2 : Recursion

$$\begin{split} \bar{D}_{1}(e) &= \min_{b} [\bar{D}(b,e) + \bar{D}_{0}(b-1)] \qquad e_{min}^{1} \leq e \leq e_{max}^{1} \\ \bar{D}_{2}(e) &= \min_{b} [\bar{D}(b,e) + \bar{D}_{1}(b-1)] \qquad e_{min}^{2} \leq e \leq e_{max}^{2} \\ \cdots \\ \bar{D}_{j}(e) &= \min_{b} [\bar{D}(b,e) + \bar{D}_{j-1}(b-1)] \qquad e_{min}^{j} \leq e \leq e_{max}^{j} \\ \cdots \\ \bar{D}_{N_{c}-1}(e) &= \min_{b} [\bar{D}(b,e) + \bar{D}_{N_{c}-2}(b-1)] \qquad e_{min}^{N_{c}-1} \leq e \leq e_{max}^{N_{c}-1} \end{split}$$
(4.14)

• Step 3 : Final Minimum Distance

$$D^{\star} = \bar{D}_{N_c-1}(N_s - 1) \tag{4.15}$$

Table 4.3 shows the computation results of dynamic programming based on the first stage matching shown in Table 4.2 and test input word image shown in Figure 4.4.

Further minimization of computational complexity is achieved by using the statistics of Table 4.1. We have assume that a character can be split into at most four segments. However, according to the duration probability shown in Table 4.1, most characters are composed of less than four segments. Also, the statistics reveal that for some characters there is no possibility of being split into less than two segments. Furthermore, the frequently used characters, such as "i"

character	0	1	2	3	4	5
duration	1:2	1:1	1:3	1:3	1:3	1:3
character	6	7	8	9	a	b
duration	1:3	1:3	1:2	1:2	1:3	1:3
character	С	d	e	ſ	g	h
duration	1:3	1:3	1:3	1:3	1:3	1:3
	1				·····	
character	i	j	k	l	m	n
character duration	i 1:2	$\frac{j}{1:3}$	k 1:3	l 1:2	m 2:4	n 1:3
character duration character	i 1:2 0	j 1:3 p	<u>к</u> 1:3 q	l 1:2 r	m 2:4 s	n 1:3 t
character duration character duration	i 1:2 0 1:3	j 1:3 p 1:3	k 1:3 q 1:3	l 1:2 r 1:3	m 2:4 s 1:3	n 1:3 t 1:3
character duration character duration character	i 1:2 0 1:3 u	j 1:3 p 1:3 v	k 1:3 q 1:3 w	l 1:2 r 1:3 x	m 2:4 s 1:3 y	n 1:3 t 1:3 z

Table 4.4: Variable duration (min:max)

and "*l*", are split into less than three segments 98% of the time. Table 4.4 shows the variable duration based on Table 4.1. The lower limit as well as the upper limit of the number of segments is specified and used for determining the size of matching window in the subsequent matching procedures. It should be noted that introducing the variable duration results in improvement of recognition accuracy as well as speed as evidenced by the experimental results in Section 4.4.

To incorporate the notion of duration probability as part of finding the best match (Lee & Rabiner, 1989), the local distance shown in Eq. 4.9 is modified as follows:

$$\hat{D}(v,b,e) = d(\mathbf{F}(S(b,e)), \mathbf{x}_{c_i(j)}(v)) - \lambda \cdot dur(e-b+1, c_i(j)) \quad \text{for } 0 \le v < N_{c_i(j)}$$
(4.16)

where $dur(\cdot, \cdot)$ is the duration probability defined in Eq. 4.1 and λ is estimated based on characteristics of the feature set. Experiments show that the incorporation of the duration probability improves recognition accuracy as well.

Figure 4.5(b) illustrates comparisons between possible combination of segments of an input image (Figure 4.5(a)) and reference feature vectors of codewords



Figure 4.5: An example of matching procedure: (a) segmentation result of an image (b) dynamic matching with variable duration for the image (a)



Table 4.5: Boundary of each character

to find the best match. An arrow represents a match between the codewords of a character in the vertical axis and a segment or a combination of segments covered by the arrow in the horizontal axis. Interval between the beginning segment and the ending segment of matchings at the position of a character defines the size of the matching window.

In Figure 4.5(b), the maximum number of segments covered by an arrow is determined by the variable duration of the character under testing. For example, for the character "i", an arrow covers at most two segments based on the variable duration of the character. On the other hand, for the character "s", an arrow covers at most three segments. Dashed arrows in the figure represent comparisons that are unnecessary due to the lower limit of the variable duration. In the example, the lower limit of the variable duration of the character "w" is 2 and the comparison between segment 0 and the codewords of the character is avoided. It results in elimination of subsequent comparisons for the subsequent characters.

4.3.3 Locating character boundaries

As a byproduct of our approach, we can also obtain character boundaries in a word image from Table 4.3. Table 4.5 shows the boundary of each character. While Table 4.2 shows the character 's' best matched with the first segment using local information, Table 4.5 shows 's' made of two segments using global information. Locating character boundaries is particularly useful in generating character databases for training and testing.

4.4 Experiments and results

To evaluate the speed and recognition accuracy of the system, 3,000 postal words (digitized at 212 dpi), including city names, firm names, personal names, street names, and state names are used. Given a test word image, corresponding dictionaries are randomly generated with size of 10, 100, and 1,000 words. The true word is always present in the lexicon.

In addition to the 74 feature set described in Section 4.2, a 38 feature set is used to build an alternate faster system. The 38 feature set is a subset of the 74 feature set. It consists of 2 global features (Eq. 4.2) and 36 local features (Eq. 4.17).

$$F_{l_{ij}} = \frac{s_{ij}}{S_i} \tag{4.17}$$

where

i = sub-image number $(1, 2, \dots 9)$, *j* = slope *mod* 4, *s_{ij}* = number of components with *j*, from sub-image *i*, and *S_j* = $\max_{i} s_{ij}$

All components are classified into 4 directional categories, horizontal, vertical and two diagonal. Table 4.6 and Table 4.7 show the improvements obtained by using variable maximum number of segments and duration probability. In Table 4.6 and Table 4.7, the first column (var) of each lexicon size represents the tests with the variable matching window size. The second column (d=4) and the third

column (d=3) represent the tests with fixed matching window sizes of 4 and 3, respectively.

4.4.1 Timing profile

Table 4.6 shows the time taken by each major module shown in Figure 4.1 using the two feature sets. As the size of the matching window decreases, feature extraction time also decreases. Time taken for recognition is determined by both the lexicon size and the matching window size. It should be further noted that the image processing time is small, mainly due to the advantages of chain code representation.

4.4.2 Recognition performance

Table 4.7 shows the recognition performance of each configuration of the system. When the matching window size is fixed to 4, performance is not as good as when the window size is adaptive. The fast system with the reduced 38 feature set produces comparable accuracy specially when dealing with a small lexicon.

4.4.3 Effect of introducing the variable duration

Figure 4.6 provides a graphical illustration of the results for the 74 feature system. According to Table 4.1 more than 20% of characters "m" and more than 10% of characters "w" split into 4 segments. Both "m" and "w" characters are frequent in English. Table 4.7 and Figure 4.6(a) reflects the drop in performance due to images with "m"s and "w"s when the window size is 3.

4.4.4 Failure analysis

Three causes of error have been identified:

Lexicon	duration	va	ar	d =	= 4	d =	= 3
Size	feature	74	38	74	38	74	38
10	Chain Gen.	22	22	21	22	21	22
	Preprocessing	10	10	10	9	9	9
	Segmentation	15	15	15	15	15	15
	Feature Ext.	34	32	34	31	24	22
	Recognition	145	67	172	80	136	63
	Total	226	146	252	157	205	131
100	Chain Gen.	22	22	22	22	22	22
	Preprocessing	10	10	10	10	9	10
	Segmentation	15	15	15	15	15	15
	Feature Ext.	34	32	34	32	25	22
	Recognition	324	147	379	147	308	226
	Total	405	226	460	226	379	295
1000	Chain Gen.	27	27	26	26	26	26
	Preprocessing	11	11	10	11	10	10
	Segmentation	15	15	16	15	16	15
	Feature Ext.	37	35	37	34	26	24
	Recognition	633	387	720	439	610	383
	Total	723	475	809	525	688	458

Table 4.6:Timing profile in msec on Spare 10

							1		
Lexicon	duration feature		Va	ar	d =	= 4	d =	= 3	
Size			74	38	74	38	74	38	
10	Тор	1	96.80	95.80	96.63	95.40	95.83	95.10	
		2	98.63	98.26	98.76	98.23	98.30	97.80	
100	Тор	1	88.23	85.86	86.83	84.63	86.83	84.96	
		2	93.36	91.93	92.86	91.20	92.36	90.80	
		20	98.93	98.96	99.00	99.03	98.40	98.76	
1000	Тор	1	73.80	68.70	71.90	65.20	72.40	67.90	
		2	83.20	80.20	81.40	78.20	82.20	78.90	
	ł	50	98.70	97.90	98.50	97.50	97.60	97.20	

Table 4.7: Recognition accuracy (in %)



Figure 4.6: Performance and speed comparison: (a) performance - the system with variable duration is represented by solid line and gives the best recognition accuracy for all lexicon sizes, (b) speed of the system with variable duration is between the two other systems.



Figure 4.7: Error Analysis: (a) background noise, (b) noise introduced by scanning device, (c) oversegmentation, (d) unable to segment, (e) recognition error (confused with a similar entry in lexicon), and (f) recognition error (due to tough edge)

ſ	Туре	Recognition	Segmentation	Image quality
ſ	No. of images	51 (53.1%)	35 (36.5%)	10 (10.4%)

Table 4.8: Failure analysis

- 1. *Image quality*: Figure 4.7(a) shows an example of severe background noise and Figure 4.7(b) shows typical scanning noise not handled by our preprocessing routines.
- 2. Segmentation: Figure 4.7(c) shows character "O" split into four segments. The matching window size of the character "O" is three (Table 4.4) hence the matching procedure ignores the last segment of the character reducing the matching confidence. This type of error is closely related to binarization which introduces fragmentation of strokes. Figure 4.7(d), on the other hand, shows the stroke of "L" touching the "e" causing a segmentation problem.
- 3. If the lexicon has very "similar" entries there is scope for confusion. "*Chuck*" is selected as the first choice and "*crooks*"(truth) is the second choice (Figure 4.7(e)). Roughness in chain code contour usually reduces the confidence of matches (Figure 4.7(f)).

We examined 96 images classified as failures (74 feature based system and lexicon size of 10). Table 4.8 shows the different causes of failure.

Chen and Kundu (Chen & Kundu, 1993) report results on the same data set and with the same lexicons as we have used. They report top choice recognition performances of 93.2%, 85.2%, and 64.6% for lexicons of sizes 10, 100 and 1,000 respectively. This compares with the performance numbers: 96.8%, 88.23%, and 73.8% reported in this chapter. Also, the top choice recognition accuracy of 88.29% obtained by us compares with 72.3% reported in (Chen et al., 1994).

4.5 Summary

A fast chain code based word recognition system has been implemented. The speed is $100 \sim 200 \text{ msec}$ on a single Sparc 10 platform with a 10 word dictionary. The top choice performance is 97%. Development of efficient methods for preprocessing, segmentation, and feature extraction result in speed improvements. Use of the variable duration in the recognition procedure improved performance as well as speed. The size of the feature set, type of features, and the concept of restricted matching window size all contribute to improvements in speed. Processing time increases at the rate of log(N) as lexicon size N increases.

Chapter 5

Street Name Recognition

Recent research in recognition of general handwritten text has revealed a number of complex problems (Govindaraju et al., 1994). While recognition of isolated units of writing, such as a character or a word, has been extensively studied in literature (Bozinovic & Srihari, 1989; Favata, 1993; Madhvanath & Govindaraju, 1993; Simon, 1992), 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*, which concerns extracting words from a phrase, (ii) *word recognition*, which concerns handwritten word recognition algorithms with and without lexicon and (iii) *postprocessing*, which concerns use of linguistic constraints to improve recognition performance

Our focus in this chapter is on the first two stages. The importance of correct word segmentation, should be evident from the sequential architecture of a 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 5.1 illustrates an overview of the system described in this chapter. 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 preprocessing 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 chapter, 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 5.1, the desired output is 3 sub-images containing the words: "S", "Michigan", and "Ave", respectively. Recognition concerns the matching between the lexicon and the 'main' part of the image. Based on the goodness of the match, a ranked list of the lexicon entries is obtained as the output.

In this chapter, 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 by extracting information from character segments. Computational complexity, a serious concern while matching lexicon entries and all possible combinations of segments, is addressed by the 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 (Table 4.1) generated by our character segmentation algorithm applied to thousands of



Figure 5.1: System overview: Input is a phrase of multiple words and a lexicon of single words (extracted from the main part of a street name). Output is a ranked lexicon sorted in the order of "goodness" of match between the image and lexicon entries.

handwritten words (Kim & Govindaraju, 1995).

The test bed used in this work is a dataset of street name images extracted from mail pieces. Recognizing the street name is an integral part of an address interpretation system (Govindaraju et al., 1993). 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: "*Michigan*" in the example of Figure 5.1). Matching between the lexicon and the input image (street name) is performed by a dynamic programming based word recognizer (Kim & Govindaraju, 1995).

Organization of this chapter is as follows. Section 5.1 describes previous work on general word segmentation. Section 5.2 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 5.3 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 5.4 is about experiments and results. Summary of the work and some concluding remarks are presented in Section 5.5.

5.1 Previous work in word segmentation

Most of the previous work on recognition of handwriting assumes that words are already isolated or that isolation is trivial (*e.g.*, words are separated by large amounts of white space) (Baird & Thompson, 1990), or that words are written in boxes whose location is known (Duda & Hart, 1968). Few published reports describe methods of separating words in unconstrained handwriting (Cohen et al., 1991a; Cohen et al., 1991b). Brady showed how certain filters can highlight the word gaps in grey-level machine-printed text images (Brady, 1981). The paper focused on an explanation for certain psychological text-interpretation studies. The method was not tested on handwritten text. There are limited number of research reports that present a comprehensive discussion on this topic (Mahadevan & Nagabushnam, 1995; Seni & Cohen, 1994).

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.

The commonly adopted computational approach to word separation using spatial distance cues consists of the following steps: (i) determine the connected components in the given line, (ii) compute the distance (or gap) between pairs of adjacent components, (iii) sort the gaps in descending order of magnitude, and (iv) classify the gaps into inter-word gaps and inter-character gaps by choosing a threshold. Gaps greater than the threshold are deemed to be word separation points.

Computing the distance between adjacent components ((ii) above) has been an open issue which warrants further research. The objective is to obtain an estimate of the inter-component gaps as perceived by humans.

The first and most straightforward estimation method 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 (Seni & Cohen, 1994), 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. The third technique approximates the gaps between components by the distance between their *convex hulls* (Mahadevan & Nagabushnam, 1995).

Figure 5.2 illustrates the gaps estimated by the above techniques between two connected components. The bounding box method fails to report any gap in the example (Figure 5.2(a)). The convex hull method reports a reasonable estimate of the gap in the example (Figure 5.2(c)).

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

Craw

Crawford st

Crawford st

Figure 5.2: Inter-component gaps using geometric relationship between components generated by (a) bounding box method (b) run length based method and (c) convex hull based method

between words than between characters. Second, the perceived space between components cannot be easily estimated by a 1-dimensional scalar metric.

5.2 Problem description

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. For example, if the phrase in question is extracted from the legal amount on a bank check then the lexicon for each word in the phrase can be a list of about 40 words (*e.g.*, one, two, three, \cdots , ten, twenty, \cdots , hundred, \cdots). If the application is one of reading the employee name from a form of personal information, the lexicon is a list of names of all the employees. In the case of street names on mail pieces, the lexicon is a list of street names in a given ZIP code. Thus the size and nature of a lexicon can vary depending on the application.

Generating a lexicon for the street name recognition task is somewhat involved. As can be seen in Figure 5.4, 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 entry 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, such as Figure 5.3(a), wherein to account for a single street name the lexicon must contain 18 different forms of the name (Figure 5.3(b)) 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 5.4(i)). Certainly, the lexicon entries do not anticipate their presence. Furthermore crossing out of writing (Figure 5.4(i)) is an additional source of concern.



(a)

(b)

Figure 5.3: An example of lexicon expansion: (a) a numeric street name (b) expanded lexicon entries for the image

S9thSt South9thSt 9thSt SNinthSt

SNinth SouthNinth Ninth

5.3 Methodology

A street line can be logically divided into two parts: *street number* and *street name*. Sample images of phrases extracted from street lines of addresses are shown in Figure 5.4. 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 5.4(g)). Typically, street name recognition involves segmenting the phrase into the constituent words and recognizing the individual words. Several approaches have been described in the literature (Govindaraju et al., 1993; Govindaraju et al., 1994; Mahadevan & Nagabushnam, 1995; Seni & Cohen, 1994; Simon et al., 1994) to address this issue.

In the subsequent sections we present an elegant method of extending an



Figure 5.4: Examples of street name images: (a) street line address, (b - c) the same street name with different pre-directionals, (d - e) numeric street name, (f - g) suffix, (h) directional suffix, and (i) miscellaneous addons

existing word recognition engine that operates on single words to handle phrases of several words. There are two major parts to the methodology: (i) segmentation and (ii) recognition.

5.3.1 Segmentation

A character segmentation algorithm locates segmentation points between characters in a word is explained in Section 4.1. Segmentation points are determined using features like ligatures and concavities.

We propose a word segmentation method (that divides a phrase into words) based on the character segmentation method. 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 Figure 5.5(a). It is to be noted that the gap between "n" and "A", is significantly less than gap 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 5.5(b)). The original character segmentation algorithm segments the image into characters and sub-characters (Figure 5.5(c)).

Four reference lines have been traditionally used to extract spatial features: upper bound, upper half, lower half, and lower bound (Bozinovic & Srihari, 1989; Seni & Cohen, 1994). However, all these lines cannot be accurately located, specially when dealing with phrases (as opposed to single words) and unconstrained

Trenton the Trenton

(a)

(b)





Figure 5.5: 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

handwriting. Therefore, we choose to define two reference lines: base line and middle line (Figure 5.5(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.

- Only the boxes positioned certain distance (h in Figure 5.6) above the base line are considered, where h = average_box_height/3.
- If a box belongs completely to another box, the box is merged (Figure 5.6(a)).

- If part of the area of a box overlaps with another box, the box is merged (Figure 5.6(b)).
- If a box overlaps in the horizontal direction without touching, the box is merged (Figure 5.6(c)).



Figure 5.6: Merging bonding boxes of character segments: (a) a box completely belonged to another, (b) a box partially overlapped, and (c) a box overlapped in horizontal direction without touching

For example, in Figure 5.5(d), a horizontal stroke over the vertical stroke of letters "*T*" and "*r*" is merged into the vertical stroke of "*T*". After the merge is completed it forms a larger bounding box and no further merging is possible between the box and letter "*r*".

Figure 5.5(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 5.5(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 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



Figure 5.7: Structure of the neural network for word segmentation

 C_i and C_{i+1} and H_i be the height of B_i and defined as follows:

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

$$H_{i} = \begin{cases} Y_{u}(C_{i}) - Baseline(X(C_{i})) & \text{if } Baseline(X(C_{i})) \ge Y_{b}(C_{i}) \\ Y_{u}(C_{i}) - Y_{b}(C_{i}) & \text{otherwise} \end{cases}$$
(5.1)

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 . *Baseline*(\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 5.7. Following 8 parameters are selected for inputs of the neural network: 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}), a flag for inter-component transition (F_{tr}), and

Irent

Figure 5.8: Word segmentation result on image of Figure 5.5

a flag for the last segment (F_l). It should be noted that the height information as well as interval information of segments is included.

Figure 5.8 shows the word segmentation result of the image of Figure 5.5 using the neural network.

5.3.2 Recognition

To explain the matching procedure in this application, let us assume that the example image shown in Figure 5.5(a) is split into three segments as shown in Figure 5.9 instead of two (Figure 5.8). Also let us assume that there are two lexicon entries, "trenton" and "main".

Each word segment is submitted to the word recognition engine along with the number of character segments. Matching between a word segment (or a combination of word segments) and the lexicon entries is performed for all possible combinations of the word segments. The final matching result for each lexicon entry is represented by the best match. In Figure 5.9 the matching confidence value of "trenton" and "main" are 4.325 and 6.746 respectively.

A way of minimizing the computational complexity based on the concept introduced in the word recognition engine is considered. As explained in Chapter 4, variable duration of each character is obtained from training phase (Table 4.4). In this application, we define the lower and upper bounds of variable duration for each lexicon entry using Eq. (5.2).

					OUTPUT			
					segment(s)	# seg	trenton (7:21)	main (5:12)
			trenton main		Tren	6		6.746
segment	seg(ch)		¥		Trenton	10	4.325	7.871
Tren	6	+	WORD RECOGNITION	+	Trenton Are	15	5.399	
ton	4 		ENGINE		ton	4	• • •	
/tve	5				ton Ave	9	6.944	7.021
					Ave	5		9.111



$$lex_{min}(i) = \sum_{j} dur_{min}(lex_entry[j]) \quad \text{for } j = 0 \text{ to } N_c(i) - 1$$

$$lex_{max}(i) = \sum_{j} dur_{max}(lex_entry[j]) \quad \text{for } j = 0 \text{ to } N_c(i) - 1$$
(5.2)

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, $dur_{min}(\cdot)$ and $dur_{max}(\cdot)$ return the lower and upper bounds of variable duration of a character defined in Table 4.4, and $lex_{min}(i)$ and $lex_{max}(i)$ are the lower and upper bounds of variable duration of the *i*-th lexicon entry.

Using Eq. (5.2) the lower and upper bounds of the variable duration of lexicon entries are defined as following:

trenton: min =
$$t(1) + r(1) + e(1) + n(1) + t(1) + o(1) + n(1) = 7$$

max = $t(3) + r(3) + e(3) + n(3) + t(3) + o(3) + n(3) = 21$
main: min = $m(2) + a(1) + i(1) + n(1) = 5$
max = $m(4) + a(3) + i(2) + n(3) = 12$

The following rule is used to skip unnecessary computation.

for each lexicon entry
$$\label{eq:scalar} \begin{array}{l} \mbox{if } (\ WS_{num} < lex_{min} \ \mbox{or} \ WS_{num} > lex_{min} \) \ \mbox{reject} \\ \\ \mbox{else perform matching} \end{array}$$

where, WS_{num} is the number of character segments for the word segment being tested. In the example, 5 cases (marked by "—" at the output table) are rejected based on this rule at an early stage. Furthermore, if a part of a word segment has been compared with a particular character in a lexicon entry, future comparison with the same character is avoided by saving all matching pair results.

5.4 Experiments and results

5.4.1 Word segmentation

Street name images collected from live mail pieces were used for training and testing the neural network. 518 images were used for the training. 7,802 training patterns were collected using the images by visual examination with desirable word segmentation points marked manually. Figure 5.10 shows a sample image and the corresponding 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 a handwritten address interpretation system (Govindaraju et al., 1993). Table 5.1 shows the word segmentation results on the two sets of images. 82-86% test images are perfectly

Toledo Jane



i	I_i	H_i	I_{i+1}	H_{i+1}	F_{tr}	I _{av}	H_{av}	F_l	Truth
1	0.714	0.469	0.616	0.922	1.000	0.479	0.497	0.000	0.000
2	0.616	0.922	0.470	0.391	0.000	0.479	0.497	0.000	0.000
3	0.470	0.391	0.551	0.406	0.000	0.479	0.497	0.000	0.000
4	0.551	0.406	0.260	0.734	1.000	0.479	0.497	0.000	0.000
5	0.260	0.734	0.324	0.281	0.000	0.479	0.497	0.000	0.000
6	0.324	0.281	0.357	0.422	0.000	0.479	0.497	0.000	0.000
7	0.357	0.422	1.000	1.000	0.000	0.479	0.497	0.000	0.000
8	1.000	1.000	0.389	0.406	1.000	0.479	0.497	0.000	1.000
9	0.389	0.406	0.292	0.375	1.000	0.479	0.497	0.000	0.000
10	0.292	0.375	0.422	0.188	0.000	0.479	0.497	0.000	0.000
11	0.422	0.188	0.292	0.234	0.000	0.479	0.497	0.000	0.000
12	0.292	0.234	0.535	0.266	0.000	0.479	0.497	0.000	0.000
13	0.535	0.266	0.000	0.000	0.000	0.479	0.497	1.000	0.000

(c)

Figure 5.10: A sample image and training patterns: (a) original image (b) result of character segmentation and merging (c) training patterns

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 correct segmentation rate of 91%.

Figure 5.11(a - d) show examples of perfect segmentation, given the tolerance of the recognition methodology that follows. These may be considered to be errors if alternate recognition methods were to follow. Figure 5.11(e - g) are typical errors discovered in the test.

	no. of images	perfect	oversegmentation	error
set l	518	446	26	46
(training set)		(86.1%)	(5.0%)	(8.9%)
set2	107	88	8	11
(test set)		(82.2%)	(7.5%)	(10.3%)

Table 5.1: Test results of word segmentation

5.4.2 Recognition

Among the 107 street name images, 11 images which were classified as word segmentation error, were excluded during the recognition performance evaluation phase. Table 5.2 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. Method3 uses the word segmentation method described in this paper. 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



(a)









(f)





(c)



Snow Strut BERGEN New Which BEBGED (d) (h)

Examples of segmentation results from our system: (a - d) segmen-Figure 5.11: tation points perfectly located; (e - g) failures; (h) over-segmentation cases

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. It should be noted that the ratio of lexicon sizes used in Method1 against Method2 and Method3 is 4.6:1. Average time (on Sparc 10) consumed by each method is shown in the lower half of Table 5.2. It needs

		Met	hod 1	Meth	od2	Meth	iod3
	S	F	S	F	S	F	
accu	racy	79	17	83	13	81	15
	2			(60)		(56)	
	image		.5.0	89	.0	122	2.8
time	time recog		523.8		1087.7		1.6
(msec)	total	63	8.8	117	6.7	704.4	

Table 5.2: Recognition performance - numbers in () represent the number of cases resulted in better recognition confidence comparing to that of Method1 (S : Success; F : Failure)

to be noted that the number of cases resulted in better recognition confidence comparing to that of Method1 (the numbers in () in Table 5.2) is significantly increased, as well as increase of total number of correctly recognized cases.

5.4.3 Selection of input parameters for the neural network

To prove the effectiveness of the features selected for the word segmentation some tests were performed by limiting the number of input parameters to the neural network. Two simple neural networks were created with a 5-3-1 structure: *Temp_net1* and *Temp_net2*. *Temp_net1* has features I_i , I_{i+1} , F_{tr} , I_{av} , and F_l and *Temp_net2* has features I_i , H_i , F_{tr} , I_{av} , and H_{av} as inputs. *Temp_net1* is to prove the effectiveness of height information and *Temp_net2* is to prove the effectiveness of parameter selection from two consecutive character segments. As we can see

		perfect	oversegmentation	error
set1	original	446 (86.1%)	26 (5.0%)	46 (8.9%)
(518)	Temp_net1	398 (76.8%)	35 (6.8%)	85 (16.4%)
L	Temp_net2	393 (75.9%)	29 (5.6%)	96 (18.5%)
set2	original	88 (82.2%)	8 (7.5%)	11 (10.3%)
(107)	Temp_net1	81 (75.7%)	8 (7.5%)	18 (16.8%)
L	Temp_net2	82 (76.6%)	6 (5.6%)	19 (17.8%)

Table 5.3: Word segmentation performance with different inputs

in the Table 5.3, lack of height information (*Temp_net1*) and use of an interval and height information from a single segment (*Temp_net2*) cause significant drop in segmentation performance comparing to the original results presented in the Table 5.1.

5.5 Summary

In this chapter, 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. Experimental results prove the effectiveness of parameters selected to be input to a neural network.

The recognition algorithm described in this chapter minimizes common problems in keyword recognition in street names. Lexicon expansion is used to take care of all possible variations of prefixes and suffixes.

The speed of the system is comparable to a system which treats a street name as a single word.

We have experimentally shown the effectiveness of our approach in terms of computational speed and accuracy. 83% phrase recognition accuracy was achieved on a test set comprising of street name images extracted from mail pieces.
Chapter 6

Reading Bank Check Amounts

It is estimated that Americans write about $50 \sim 70$ billion bank checks a year (John, 1992; Zimmerman, 1995). The banks in Unites States therefore have to deal with processing many million checks everyday. Each check is handled up to six times at the bank. Recent advances in automated recognition of bank checks have attempted to minimize the manual effort involved.

Automation of bank check recognition has several salient aspects. Wide variations in handwriting style exist not only in terms of cursive or discrete but also in the manner in which amounts are represented. For example, "*one thousand and five hundred dollars*" and "*fifteen hundred dollars*" are different ways of representing the amount \$1,500. Moreover, the absence of constraints on the number of digits that can appear in the numeral amount and the need for high recognition performance (Moreau et al., 1991; Paquet & Lecourtier, 1991) compound the difficulty of the task. The one advantage that this application offers is the redundancy between the legal field (text) and the courtesy field (numeral) as they both convey the same information (Cohen, 1994).

Bank check recognition systems can be divided into two categories: (i) systems that read only one of the amounts – legal or courtesy (Agarwal et al., 1995; Dodel

Issues	Legal	Courtesy
Amount Extraction	0	0
Touching digits		0
Touching characters	0	
Decimal point detection		0
Presence of dollar and cent		
amounts in the same field		0

Table 6.1: Issues specifically related to reading each field

& Shinghal, 1995; Guillevic & Suen, 1995; Paquet & Lecourtier, 1991; Simon et al., 1994), and (ii) systems that read both amounts and use a combination of the results (Cohen, 1994; Moreau et al., 1991). In this chapter, we introduce a method for assigning the amount value of bank checks by using *both* the legal field and the courtesy field. Input to the system is a binary image of a bank check. Output is a list of choices for the value of the check. A lexicon driven handwritten word recognizer described in Chapter 4 is the core of the system. Table 6.1 outlines the problems specific to the recognition of the legal and the courtesy amounts.

Organization of this chapter is as follows. Section 6.1 presents an overview of bank check recognition and outlines general challenges of the task. Section 6.2 describes a strategy for locating word breaks in the legal amounts. Section 6.3 outlines the recognition procedures for the legal and the courtesy amounts using a word recognizer and a cross validation method. Section 6.4 details our experimental results and error analysis. Summary of the work and some concluding remarks are presented in Section 6.5.



Figure 6.1: The components of a typical bank check

6.1 System overview

Figure 6.1 illustrates the components of a typical bank check. The objective of the check reading system is to obtain the value of the check by reading the legal and the courtesy amounts. In the subsequent sections we describe the main challenges of this task: (i) strategies for reading the two fields and cross validation of one field with the other, and (ii) segmentation of the text phrase in the field into its constituent words.

Figure 6.2 illustrates the system architecture. The legal amount is recognized by the lexicon driven word recognizer. A list of hypothesized amounts is passed onto the courtesy amount recognizer. A numeral string word recognizer uses the list of hypotheses as its lexicon. It should be noted that the numeral string recognizer operates in exactly the same way as the legal phrase recognizer. In the final phase, the confidence of match from the recognition of both the fields is used to minimize error rate.

The legal and courtesy fields are located by an analysis of the check image. Cues, such as presence of two long horizontal lines near the legal amount, and the dollar sign ('\$') and a box near the courtesy amount are used.

Preprocessing includes noise removal, underline removal, box removal, slant





correction, and converting the sub-images into chain code (Freeman, 1961) representation. Underline and box images are removed by the analysis of connected components represented by run-length coding (Pavlidis, 1982). Handwriting strokes crossing the lines are saved. Small components (compared to the average stroke width) are regarded as noise and removed. Vertical strokes are extracted and the slant angle is estimated by averaging the angles of the strokes based on their vertical length. Slant correction is performed according to the estimated angle directly on the chain code representation. Detailed preprocessing procedures are explained in Chapter 3.

The legal amount is regarded as a text phrase of one or more words. Word segmentation is performed to divide the phrase into its constituent words which can be subsequently recognized.

Detection of the boundary between dollar and cent amounts is another serious issue that must be addressed during the recognition of both fields. The different ways of representing amounts makes exact location of the boundaries difficult (Figure 6.3).

6.2 Segmentation

The goal of this module is to separate the words in the legal amount. The input to this module is the set of connected components forming the image of a legal amount. The output is a list of hypotheses, where each hypothesis is a partition of the original image into sub-images. A sub-image can correspond to exact words or sub-words.

6.2.1 Problem definition

Word segmentation in the legal field of bank checks is made difficult by the nature of the data and patrons' writing habits. The number of words in the fields is



Figure 6.3: Examples of amounts in which it is difficult to locate the boundary between dollar and cent amounts: (a) courtesy amounts (b) legal amounts

not known and hard to predict. Furthermore, patrons tend to clump the words together so that additional words are not illegally squeezed (Figure 6.4(a, b)). And a long horizontal line is drawn at the end of the amount to use up the entire space in the field (Figure 6.4(c)) which makes the segmentation difficult.

6.2.2 Word segmentation based on character segmentation

In the following subsections we describe our approach to word segmentation using spatial distances as well as recognition. This method captures the author's writing style 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 ligature is

Threaher dred eighty and 35/100 mly-(a) Two thousand and fourty five only Minity five and too m (c)

Figure 6.4: Examples of legal amount make segmentation difficult: (a - b) not enough spaces between words, (c) long horizontal line at the end of the field

encoded into the segmentation scheme. Furthermore, we make the additional assumption that character segments must have at least one significant vertical stroke (Figure 6.5(c)) as explained in Section 5.3.1.

Our approach has proved to be efficient and powerful in the application of bank checks for the following reasons: (i) size of the lexicon is small, (ii) character segmentation is the necessary step for word recognition, and (iii) word recognition is used to drive the word segmentation process.

Consider the legal amount image of Figure 6.5(a). The image is converted into a contour representation and the slant is estimated and corrected. A long horizontal stroke at the end of the amount (which is drawn to use up the entire space in the amount field) hinders the extracting of information of the author's writing style. The stroke is removed based on the density of components obtained by computing the ratio between the number of contour pixels and the length of the components. The character segmentation algorithm segments the image into



Figure 6.5: Word segmentation (based on character segmentation): (a) original input image (b) character segmentation result after slant correction (c) centerline of each segment after merging (d) word segmentation result

characters and sub-characters (Figure 6.5(b)).

Two reference lines are defined as explained in Section 5.3.1. Based on the two reference lines, segments are examined and merged. Figure 6.5(c) shows bounding boxes after the merge.

A vertical centerline is located based on the distribution of pixels inside the box. The intervals between two adjacent boxes characterize the author's writing

en Zy Ss	five ten fifteen twenty seventy dollars	four nine fourteen nineteen sixty thousand	three eight thirteen eighteen fifty hundred only	two seven twelve seventeen forty ninety and	one six eleven sixteen thirty eighty dollar
----------------	--	---	--	---	---

Table 6.2: Lexicon entries used for the recognition of the legal amount

style in terms of inter-character and inter-word spacing. The intervals are normalized, and average and standard deviation are computed based on which word gaps are defined (Figure 6.5(d)). The letter "d" in the word "hundred", which is classified as a segment, is eventually combined with the word "hundred" during the recognition phase.

Word segments are submitted to the word recognition engine along with the lexicon shown in Table 6.2 and recognition results are collected as shown in Figure 6.6. The numbers in the left and top side of the figure correspond to the segment numbers of Figure 6.6(a). They define the beginning and the ending segment numbers of combined segments. Matching confidences are used for determining the global matching score for each meaningful amount.

To minimize computational complexity we take the advantage of the word recognition engine. For each entry, the number of characters and the maximum number of segments allowed are compared with the number of character segments in the word segment before the match is performed. Furthermore, if a part of the word segment has been compared with a particular character in a lexicon entry, future comparison with the same character is avoided by saving the matching results. For example, if the lexicon entry "*two*" has been compared to the word segment "*three*", a part of computation is not duplicated when the entry "*three*" is compared to the word segment because the character "*t*" in the

first position is common.

Exhaustive search of all possible amounts that can be derived from the recognition results minimizes the risk of missing the true word in the list. However, this increases the size of the lexicon processed by the courtesy amount recognition. Following heuristics are used to determine which results are to be used to form a list as the lexicon to the subsequent stage of courtesy amount recognition (boxes in Figure 6.6). The decision is made based on only the confidence of the top choice.

- 1. If the distance score is less than 3.0 (very high confidence) then select that column alone from the corresponding row.
- 2. If the difference of confidences of the best and the second best columns is within 20% of the best one then select both columns (*e.g.*, the second row in Figure 6.6(b)).
- 3. If the difference in step 2 is greater than 20% then select only the best column (*e.g.*, the first row in Figure 6.6(b)).
- 4. Groups with very low confidence (9.0 in our scale) are ignored in the selection.

Based on the above rules, only a few strings are hypothesized. It should be noted that we do not allow an overlap between two adjacent segments during the generation of hypotheses. We obtained the following 4 hypothesized strings for the example illustrated in Figure 6.6.

- 1. [(0,0): Three [(1,2): hundred [(3,3): twenty [(4,4): nine [(5,7): &80]
- 2. [(0,0): Three [(1,2): hundred [(3,3): twenty [(4,4): nine [(5,6): &8
- 3. [(0,0): Three [(1,1): hundre (2,3): dtwenty [(4,4): nine (5,7): &80
- 4. [(0,0): Three [(1,1): hundre [(2,3): dtwenty [(4,4): nine [(5,6): &8

	Thre	e huen	dred	turent	ty n	~~~?	9 	33 Tor	Ż
	0	1	2	3		4	5	6 7	,
				(a)					
	0	1	2	3	4	5		6	7
3 4 5 5 5	.6: three 9.0: .1: five 6.0: .1: four 6.0: .1: one 6.1: .4: nine 6.2:	thousand 9.0: thirteen 6.1: fourteen 6.1: twelve 6.2: nineteen 6.3:	thousand 9.0; fourteen 9.0; thirteen 9.0; nineteen 9.0; hundred 9.0;	and 9.0 dollar 9.0 dollars 9.0 eight 9.0 eighteen 9.0	: and 9. dollar 9. dollars 9. eight 9. eighteen 9.	0: and 0: dollar 0: dollars 0: eight 0: eighteen	9.0: 9.0: 9.0: 9.0: 9.0: €	and 9.0; doilar 9.0; doilars 9.0; eight 9.0; eighteen 9.0;	and dollar dollars eight eighteen
1	4.8: 5.1: 5.2: 5.3: 5.6:	nine 4.1: hundred 4.7: dollars 4.9: eleven 5.0: thirty 5.2:	hundred 6.6: dollars 6.6: eleven 7.1: thousand 7.1: seven 7.2:	nineteen 9) seventeen 9.3 thirteen 9.3 thousand 9.0 eighteen 9.0	:seventeen 9. : dollar 9. : and 9. : eight 9. : dollars 9.	0: and 0: dollar 0: dollars 0: eight 0: eighteen	9.0: 9.0: 9.0: 9.0: 9.0:	and 9.0: dollar 9.0: dollars 9.0: eight 9.0: eighteen 9.0:	and dollar dollars eight eighteen
2		9.0: 9.0: 9.0: 9.0: 9.0:	and 5.2: dollar 5.5: dollars 6.1: eight 6.2: eighteen 6.6:	twenty 6.6 seventy 7.3 ninety 7.3 eighty 7.3 sixteen 7.6	:seventeen 6. : nineteen 7. : eighteen 7. : fourteen 7. : thousand 7.	7:seventeen 4: eighteen 4: fourteen 7: nineteen 3: thousand	9.0:3e 9.0: f 9.0: 5 9.0: 5	eventeen 9.0:s Sourteen 9.0: Housand 9.0: dollars 9.0: dollar 9.0:	eventeen dollar and eight dollars
3			4.8: 6.0: 5.1: 6.2: 6.4:	twenty 7.0 ninety 7.0 seventy 7.0 forty 7.1 thirteen 7.2	: seventy 7. :seventeen 7. : fourteen 7. : nineteen 7. : thousand 7.	2:seventeen 2: fourteen 5: nineteen 6: thousand 6: thirteen	8.1:3e 9.0: f 9.0: r 9.0: t 9.0: t	eventeen 9.0:s Jourteen 9.0: Sineteen 9.0: Stirteen 9.0: Shousand 9.0:	eventeen dollar and eight dollars
7				5.7 6.4 6.5 6.6 7.3	: sixty 7. : forty 7. : fifty 7. : nine 7. : only 7.	1: sixteen 4: sixty 4: eight 5: ninety 5: seven	7.1: m 7.4: s 7.5: 7.5: 7.6: t	ninetoen 7.3: sighteen 7.6: sixteen 7.7: tifteen 7.7:s hirteen 7.3:	nineteen eighteen tifteen eventeen sixteen
5					9. 9. 9. 9. 9.	0: and 0: dollar 0: dollars 0: eight 0: eighteen	8.0: 8.1: 8.2: 8.2: 8.2:	two 7.5: dcllar 7.6: thiry 7.7: forty 7.7: three 7.7:	dollara fifty dollar thirty twelve
6							7.2: 9.0: 9.0: 9.5: 9.5:	six 6.9: ten 6 9: one 7.2: and 7.3: two 7.5:	31x two tive ine sixty
7								a'9: 3'9: 3'9: 3'9: 3'9: 3'9:	and dollar dollars eight eighteen

(b)

Figure 6.6: Recognition of the segments: (a) segment number (b) A block at row position 1 and column position 2 represents the recognition result obtained by submitting to a word recognizer a combination of segments one and two which in turn correspond to the first and the second parts of the word "hundred". The number before each word represents the matching distance. A smaller value indicates a better match. Boxes indicate the recognition results that are used to feed the lexicon to the subsequent stage of courtesy amount recognition.

6.3 Recognition of amounts

Our approach to check recognition differs from traditional methods in two significant ways. First, our emphasis on both the legal amount and the courtesy amount is balanced. Previous approaches have largely focused on courtesy amount alone because of the greater reliability of recognizing numerals. However, we have developed an accurate word recognition engine which performs equally well on alpha words and digit strings (Kim & Govindaraju, 1995). Second, our combination strategy is serial rather than the more commonly used parallel methods. The recognition engine first processes the legal amount and outputs a list of hypotheses which are submitted to the courtesy amount recognizer as a lexicon. This strategy of serial combination using a single word recognition engine, we believe, is unique to our methodology to check recognition.

6.3.1 Legal amount recognition



One of the recognition results of word segments collected from Figure 6.6 is shown in a lattice in Figure 6.7. The lattice is traversed, and meaningful *dollar* amounts are generated. We use a transition table to denote the legal bigram transitions and an exhaustive tree search algorithm to determine all possible dollar amounts. In Figure 6.7, a subset of the possible paths are indicated. Some example dollar amounts generated by the paths are: "three hundred twenty nine", "three hundred ninety nine", "five hundred twenty nine", and "five hundred ninety nine".

In addition to generating the possible paths, we also generate prefixes of valid legal amounts. For example the amounts "*three*" and "*three hundred*" are considered valid outputs. The reason for including the prefixes is that we do not know the exact boundaries of the dollar amount in the image. Thus in certain cases, the sub-images corresponding to the cents part are recognized as valid words, and the longest valid path through the lattice may not be correct.

We keep a record of occurrences of words during the traversal to estimate the order of magnitude of the amount. Six bit flags are assigned for this purpose as shown in Figure 6.8. Based on the recognized words from the legal amounts, corresponding flags are set. Only the top choice of each segment(s) with high confidence (*i.e.*, less than 6.0 in our scale) is considered while setting the flags. In our example, "*three*" set the flag **S1** which means a single digit is in the first position. Next possible transition to "*hundred*" sets the flag **H**. For all hypothesized amounts, this process is repeated and saved along with the amount. Two types of flags are provided for the single and the double digits to identify whether it is detected in the first position of the legal amount. Rules for interpreting the flags are explained in the following section. The confidence in recognition of a dollar amount is computed as the average of the recognition confidences of each word in the legal field. The list of amounts along with the confidences and flags is then submitted to the courtesy amount recognizer as the lexicon.





6.3.2 Courtesy amount recognition

The list of amounts recognized from the legal amount are converted into the corresponding digit strings and used as lexicon entries for courtesy amount recognition. The alpha character segmentation algorithm applied to the legal amount is applied to handle touching digits as well. Therefore a segment obtained by the segmenter can be a digit or a part of a digit. Since we do not attempt to locate the boundary between dollar and cents amount (in this thesis), the combined segments are submitted to the recognizer by adding one segment at a time. Figure 6.9 shows the matching results between segment(s) against the lexicon entries. Top two choices from each match are saved for the purpose of cross validation.

6.3.3 Cross validation between the legal and courtesy amounts

Based on the recognition results of both amounts using the methods described above, we try to solve the following so-called *chicken-egg* problem. While estimating the number of digits in the check amount, the legal amount is more tractable than the courtesy amount because of the fact that components of the legal amount contain the information both implicitly and explicitly. For example, the words





"hundred" and "thousand" provide the information explicitly.

The flags saved during the phase of parsing the legal amount are used for determining the number of digits. First, the lexicon for the courtesy amount recognition is sorted based on confidences and then the flags of top N (10 in our case) choices are examined. The top N choices reflect the characteristics of the amount in terms of the number of digits. N is chosen depending on the number of word segmentation hypotheses and the number of words per segment involved in amount generation. Flags for all N amounts are ORed and the result is analyzed based on the definition of each flag. Table 6.3 shows how the flags are interpreted. The 10 choices from the legal amount recognition are listed below:

rank	amount	confidence	flags	(hex)
0	3	3.5851	1	0
1	300	3.8606	1.	4
2	5	4.1358	0	i
3	500	4.1360	4	
4	320	4.1695	10	5
5	3000	4.3170	1(C
6	520	4.3531	6	
7	3020	4.4738	12	2
8	390	4.5613	14	1
9	5000	4.5923	0	

After performing the *OR* operation the flags, we obtain 16 (hexadecimal) as the final set flags. The true value (329) is not present in the list implying: (i) there should be a keyword "*hundred*" (**H** is set), and (ii) a single digit number is placed at the first position of the amount (**S1** is set).

Strict rules are applied in the interpretation of the flags. If there is any indication of confusion in terms of estimating the number of digits, the image is rejected. As we can see in the example, the role of flags \mathbf{T} and \mathbf{H} is crucial. Let us assume that the combined segment "*hundred*" is recognized as the second choice and the \mathbf{H} flag is not set. The result flags after *OR*ing will be 12 instead of 16 causing a reject.

Consider the true amount \$300 and that only the flag **S1** is set by recognition of "*hundred*" as the second choice in the position. In the courtesy amount recognition, there is a possibility for both 3 and 300 as top choices in each digit position. Therefore, in this case (*i.e.*, only the flag **S1** or **D1** is set) we examine the top choices in the corresponding digit number to see if the following two digits are 0s. If the following two are 0s we reject the case since the flag **H** is not set and there is a confusion between 3 and 300.

flags	no. of digits	examples
XX1XXX	4 or 5	fifteen thousand, three thousand
011XXX	4	three thousand
101XXX	5	fifteen thousand
001XXX	4 or 5	see note below
XXXIXX	3 or 4	fifteen hundred, three hundred
01X1XX	3	three hundred
10X1XX	4	fifteen hundred
00X1XX	4	see note below
010000	1	three
10000X	2	fifteen
11000X	1 or 2	see note below
others	reject	

Table 6.3: Interpretation of the flags used for estimating number of digits in the amount. **note:** confidence values of the top choices are also examined for the amounts with questionable number of digits

Once the number of digits is estimated, the following three items are used to finalize the amount: (i) number of digits estimated, (ii) confidence values of top two choices of the courtesy amount consistent with the number of estimated digits and (iii) confidence of the legal amount which is selected as top choice in the step (ii).

If the top choice of the courtesy amount has high confidence the corresponding legal amount with reasonable confidence is accepted. On the other hand, if the confidence of the top choice of the courtesy amount is marginal but the confidence of the corresponding legal amount is high, then the difference between the top two choices of the courtesy amounts is examined. If the distance is large the image is accepted.

6.4 Experiments and results

Because bank checks have confidential information, such as names, address, and account numbers, we were unable to use real bank checks. Instead, as part of our data collection effort we artificially created blank copies of checks (Figure 6.1) and requested students and employees in the State University of New York at Buffalo to fill out the legal amount and courtesy amount with imaginary numbers. 235 check images were collected from 91 participants, and scanned at 200 dpi. No image was rejected by visual inspection. We expected our system to intelligently reject inconsistent checks (*e.g.*, legal and courtesy amounts do not match).

The word recognition engine was trained on data from addresses in mail pieces. To provide reference feature vectors to the matching procedure, a codebook was constructed. The codewords are trained using 20,529 character images extracted from handwritten images of city names.

6.4.1 Segmentation

Performance of the word segmentation algorithm was measured using the recognition results of the legal amount. We applied the word segmentation algorithm of Section 6.3. The list of amount candidates was examined to see whether the true value was present. 80.0% of cases contained the true amount in the list. Average size of the list was 107.

6.4.2 Recognition

235 check images were submitted to the system. Some of them have inconsistency in the amounts, unexpected amounts (*e.g.*, zero dollar), and very poor image quality. In cross validation phase of both amounts, there are two steps: *step 1* estimating the number of digits in the legal amount, and *step 2* - using the confidence values of the legal and courtesy amounts, and distance between the top

[accept	reject		
step 1	183	52		
	(77.9 %)	(22.1 %)		
step 2	error rate	acceptance rate		
	0.0 %	43.8 % (103)		
	1.8 %	47.2 % (111)		

Table 6.4: Recognition performance

two choices of the courtesy amount.

Table 6.4 shows the performance of each step. Of the 183 images accepted in step 1, 132 cases contain the true amount in the lexicon for the courtesy amount recognition. Step 2 is expected to reject the cases which do not contain true amounts based on confidences.

Because of the nature of the application, we had a strict control on the error rate. Lower half of Table 6.4 shows the performance of our system for two different error rates.

We analyzed 29 cases which have true amounts after the legal amount recognition but rejected during cross validation (132 - 103 = 29). Reasons for rejection in step 2 are classified as in Table 6.5.

Figure 6.10 shows some examples of amounts processed by our system indicating that we can handle various writing styles (Figure 6.10 (a - d)), poor quality (Figure 6.10 (e)), incorrect syntax (Figure 6.10 (f)), unexpected format (Figure 6.10 (g)), and word segmentation error (Figure 6.10 (h)).

(g)

Figure 6.10: Examples of amounts images : (a - d) correctly recognized and (e - h) rejected cases

source	instances
low confidence of the legal amount recognition	10
low confidence of the courtesy amount recognition	20
not enough gaps between two top choices	4

Table 6.5: Sources of the rejection

6.5 Summary

In this chapter, we have introduced a strategy for recognition of bank check amounts. A word segmentation algorithm (driven by recognition) was designed to overcome the problem of separating words in the legal amount field. The characteristics of the word recognizer we have developed and the small size of the lexicon make the segmentation process efficient. Based on top choice confidences, a few word segmentation hypotheses are made. The order of magnitude of the amounts are estimated from the legal amount and forms the first step of cross validation. In the second step rules based on confidence values of both the amount fields and the distance between the top two choices, are applied to reject inconsistent cases. In the experiment, we correctly read 43.8 % of check images with 0 % error rate.

Chapter 7

Conclusions

7.1 Summary

In this thesis, we describe a handwritten word recognition engine and its enhancement to phrase recognition tasks. Developing efficient algorithms for real applications is the basic motivation of the development. Every functional module in the recognition engine is developed to meet both the speed and accuracy goals. Our focus is on the following issues which bear special significance in the implementation of real-time systems.

• chain code based image handling - We employ a chain code based image representation and introduce image handling procedures commonly used in most handwriting recognition systems including preprocessing, segmentation, and feature extraction. Preprocessing deals with noise removal, slant estimation and correction, and smoothing of contours. Average stroke width of an image is estimated and used in an adaptive manner to determine the size of objects. The stroke width is in turn used to identify noise and potential segmentation points. Features are directly extracted from the chain code based images. In order to account for the wide variety of writing styles,

2 global and 74 local features are defined. All possible segmentation points are determined based on ligatures and concavity features. Experimental rcsults indicate that our implementation is extremely efficient and accurate.

- **design of an efficient matching algorithm** A method of word recognition that uses word models (as opposed to character models) is considered in this thesis. The key concept underlying this approach is the early involvement of lexicon in the recognition process. Empirical observations show that lexicon size is often small in most applications. Word images are compared with only the words present in the lexicon thus eliminating any need for postprocess-ing, a time consuming procedure. Statistics on character segmentation are collected during the training and used for determining the size of matching window.
- extension of the recognition algorithm to phrase recognition Expecting an isolated word image as an input unit is not always possible in real applications. Fields in forms generally have more than one word which leads to the issue of word segmentation. A word segmentation strategy based on character segmentation (a necessary step in our recognition engine) is introduced. A neural network is designed to locate word gaps. Features extracted from the character segmentation results are provided as input to the neural network. The concept of variable duration developed for the recognition engine is extended to lessen the computational complexity of the matching procedures. Two applications are introduced in this thesis to prove the efficiency of the extension: recognition of street names from mail pieces and reading bank check amounts. To alleviate the problem of diverse representations of street name images, the main body of the image is located and recognized using dynamically generated small lexicons. In the task of bank check amount recognition, cross validation is performed based on the recognition

results of the legal and the courtesy amounts. Our combination strategy for the cross validation is serial using a single recognition engine as opposed to the more common approach of using parallel methods. Experiments show promising results in terms of both speed and accuracy.

There are still several research issues in handwriting recognition. There is still no recognition system performs as well as humans.

7.2 Contributions

We believe that the work presented in this thesis has made the following contributions in the field of handwriting recognition.

- The concept of variable duration is derived and successfully applied in a recognition engine which deals with unconstrained handwritten words. Furthermore, the concept is expanded to the recognition of phrases.
- All image handling procedures for the recognition engine are efficiently implemented using the chain code image representation.
- Segmentation and recognition algorithms which can deal with all writing styles are developed.
- A novel word segmentation approach, which is based on the character segmentation results instead of conventional geometric distance metric between components has been developed.
- Two applications for the recognition of phrases is implemented and the generality and scalability of the recognition engine has been illustrated.

Appendix A

Training

To provide reference feature vectors in the matching procedure, we use a clustering method based on character level data and build a codebook. The codewords are trained using 21,054 character images which are extracted from handwritten images of US city names using the same segmentation method used in recognition. Number of segments for each character image is saved and used to build a table which contains duration probabilities. For numerals, 18,466 digit images are used. Segmentation statistics of numerals are collected by directly applying the segmentation algorithm to the isolated digit images.

K-means clustering algorithm with a fixed signal-to-noise ratio and fixed maximum number of clusters is used. Figure A.1 shows samples of the clustering results for characters G/g and N/n as examples. The clustering algorithm gathers characters which have similar shapes in terms of features and feature vectors of the clustered characters form codewords to represent the cluster. A table containing all codewords of characters in the codebook is used in the recognition phase.

Table A.1 shows training images which are nearest to the centroids of clusters after training.



Figure A.1: Samples of clustering results for (a) character \mathbf{G}/\mathbf{g} and (b) character \mathbf{N}/\mathbf{n}

Table A.1: Clustering results of characters

`

character	samples nearest to centroids
0	0000
1	///////////////////////////////////////
	11/11///11////
2	222222222
	2, , , , , , , , , , , , , , , , , , ,
	2222222222
3	33333333333
	2 , 2 3 2 3 2 3 7 9
	ZZ2 ZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZZ
	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
4	444444444444
	444444444
	444444444

character	samples nearest to centroids
5	5-55-555555
	555555555555555555555555555555555555555
	5555555555
6	66066666666
	066666666
7	7777777777
	~7~77777777
	7-21-77777
8	8888888888888
	888888888
	888888888888

character	samples nearest to centroids
9	99999999999999999999999999999999999999

character	samples nearest to centroids
а	2a Aa A A Q a A A
	adaaa Aaaa adaaca Aaaa
b	Lbzeerbess
	BBBBBBBBBBB
	BBBBBBBBB
с	200002cccc
	CercCerCle
d	dona 2dddd
	D. Dana Dala
	a la la a la da

character	samples nearest to centroids
e	22 RER CRER
	2 CEELEELE _ 2 CEELEELE
f	AFIANAFFIZ
	475251731
	THEF37H
g	99263878JX
	1996628396
	9896893397

character	samples nearest to centroids
h	& AHAH & WAHAM
	hkah++++Ha
	ANANKEW
i	121RALII
	212 LILINE
	1221AIAizII
j	ENRER CHLPT
k	& LAKKKKKA2
	K-KR & D X K K-K
	KAXKKKKK

character	samples nearest to centroids
1	L_LLLLL del
	XXILZIRILC
	LIJIZZIZZ
m	m m m M m M m m m m
	mm mm mm m m me Mm
	mmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmmm
n	nannan unn.
	mmmMNNMnn
	ERN PRAZINAM
0	20-00-00-0
	0000000000000
	FNGDA DOJGA

character	samples nearest to centroids
р	PPPPPPPPP
	PP+PPPPpppp
	~ plpppp
q	7 P Q R
r	~ FARR2 ~ 2 RA
	+r. RRANNO
	RZENCARTIR
S	250ASS2250
	SAAASAA
	25tras\$53

character samples nearest to centroids 一十一又龙一七节龙 t Re+ fly J. M. u. a. D. u. a. d u Uma MONTANER V N m W m m m m m m V.JVU.N.M.M.M v massmallel Y de an as a V The man me

character	samples nearest to centroids			
w	will up to Mar at us 20 m			
	as WWW in W in 120 cm			
	in What is the W the We and			
x	YXXYXXXXX			
	xpxXXXXXKv			
	x 2 X Y 4 Y X A PA			
у	447444437			
	MATTYYYYY			
	HER STRANK YY			
Z	723			
Bibliography

- Agarwal, A., Granowetter, L., Hussein, K., & Gupta, A. (1995). Detection of Courtesy Amount Block on Bank Checks. In *Proceedings of the Third International Conference on Document Analysis and Recognition, Montreal, Canada* (pp. 748–751).
- Agazzi, O. E. & Kuo, S. (1993a). Joint Normalization and Recognition of Degraded Document Images Using Pseudo-2D Hidden Markov Models. In International Conference on Document Analysis and Recognition, Tsukuba Science City, Japan (pp. 155–158).
- Agazzi, O. E. & Kuo, S. (1993b). Pseudo Two-Dimensional Hidden Markov Models for Document Recognition. *AT&T Technical Journal*, (pp. 60–72).
- Baird, H. & Thompson, K. (1990). Reading chess. *IEEE Transactions on Pattern* Analysis and Machine Intelligence, 12, 552–559.
- Baird, H. S. (1993). Recognition Technology Frontiers. *Pattern Recognition Letters*, 14, 327–334.
- Belaid, Y., Belaid, A., & Turolla, E. (1995). Item Searching in Forms: Application to French Tax Form. In Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada (pp. 744–747).

- Bender, S. L. & Popovic, M. (1995). Rewritable Optical is the Best of Both World. Imaging Magazine, 4.
- Blum, H. & Nagel, R. (1978). Shape Description Using Weighted Symmetric Axis Features. *Pattern Recognition*, 10, 167–180.
- Bourlard, H. & Wellekens, C. J. (1990). Links Between Markov Models and Multilayer Perceptrons. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(12), 1167–1178.
- Bozinovic, R. M. & Srihari, S. N. (1989). Off-line Cursive Script Word Recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 11(1), 68–83.
- Brady, M. (1981). Toward a Computational Theory of Early Visual Processing in Reading. *Visible Language*, 15(2), 138–215.
- Braun, A., Caesar, T., Gloger, J. M., & Mandler, E. (1995). Preprocessing Raw Binary Images by Means of Contours. In Proceedings of the Third International Conference on Document Analysis and Recognition, Montreal, Canada (pp. 640–643).
- Breuel, T. M. (1994a). A System for the Off-Line Recognition of Handwritten Text. *IDIAP MEMO 94-02*.
- Breuel, T. M. (1994b). Design and Implementation of a System for the Recognition of Handwritten Responses on US Census Forms. In *Proceedings of the conference on Document Analysis System, Kaiserslautern, Germany* (pp. 109–134).
- Brown, M. K. & Ganapathy, S. (1983). Preprocessing Techniques for Cursive Word Recognition. *Pattern Recognition*, 16(5), 447–458.
- Burr, D. J. (1983). Designing a Handwriting Reader. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-5(5), 554–559.

- Chen, M.-Y. (1993). Handwritten Word Recognition using Hidden Markov Models. PhD thesis, Dept. of ECE, State University of New York at Buffalo.
- Chen, M.-Y. & Kundu, A. (1993). An Alternative to Variable Duration HMM in Handwritten Word Recognition. *Pre-Proceedings, Third International Workshop* on Frontiers in Handwriting Recognition, Buffalo, NY, (pp. 82–91).
- Chen, M.-Y., Kundu, A., & Srihari, S. N. (1995). Variable Duration Hidden Markov Model and Morphological Segmentation for Handwritten Word Recognition. *IEEE Transactions on Image Processing*, 4(12), 1675–1688.
- Chen, M.-Y., Kundu, A., & Zhou, J. (1994). Off-Line Handwritten Word Recognition Using a Hidden Markov Model Type Stochastic Network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5), 481–496.
- Cohen, E. (1994). Computational theory for interpreting handwritten text in constrained domains. *Artificial Intelligence*, 67, 1–31.
- Cohen, E., Hull, J., & Srihari, S. N. (1991a). Understanding Handwritten Text in a Structured Environment : Determining ZIP codes from addresses. *International Journal of Pattern Recognition and Artificial Intelligence*, 5(1&2), 221–264.
- Cohen, E., Hull, J., & Srihari, S. N. (1994). Control Structure for Interpreting Handwritten Address. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(10), 1049–1055.
- Cohen, E., Hull, J. J., & Srihari, S. N. (1991b). Understanding Handwritten Text in a Structured Environment: Determining ZIP Codes from Addresses. In P. Wang (Ed.), *Character and Handwriting Recognition: Expanding Frontiers* (pp. 221–264). World Scientific.

- D. S. Doermann and A. Rosenfeld (1995). Recovery of Temporal Information from Static Images of Handwriting. International Journal of Computer Vision, 15(1/2), 143–164.
- Dimauro, G., Grattagliano, M. R., Impedovo, S., & Pirlo, G. (1993). A System for Bankchecks Processing. In International Conference on Document Analysis and Recognition, Tsukuba Science City, Japan (pp. 454–459).
- Dodel, J.-P. & Shinghal, R. (1995). Symbolic/Neural Recognition of Cursive Amounts on Bank Cheques. In Proceedings of the Third International Conference on Document Analysis and Recognition, Montreal, Canada (pp. 15–18).
- Duda, R. O. & Hart, P. E. (1968). Experiments in Recognition of Hand-printed Text. In AFIPS Conference Proceedings (pp. 1139–1149).
- Duda, R. O. & Hart, P. E. (1973). Pattern Classification and Scene Analysis. John Wiley & Sons, Inc.
- Earnest, L. D. (1962). Machine Recognition of Cursive Writing. In C. M. Popplewell (Ed.), *Information Processing* (pp. 462–466). Amsterdam: North-Holland Pub. Co.
- Eden, M. (1962). Handwriting and Pattern Recognition. *IRE Transactions on Information Theory*, IT-8, 160–166.
- Ehrich, R. W. & Koehler, K. J. (1975). Experiments in the Contextual Recognition of Cursive Script. *IEEE Transactions on Computers*, c-24(2), 182–194.
- Essex, D. (1995). Rewriting Handwriting Recognition. Byte, 20, 150-151.
- Farag, R. (1979). Word-Level Recognition of Cursive Script. *IEEE Transactions on Computers*, c-28(2), 172–175.

- Favata, J. T. (1993). Recognition of Cursive, Discrete and Mixed Handwritten Words Using Character, Lexical and Spatial Constraints. PhD thesis, Dept. of Computer Science, State University of New York at Buffalo.
- Freeman, H. (1961). Techniques for the Digital Computer Analysis of Chain-Encoded Arbitrary Plane Curves. In Proceedings of the National Electronics Conference, volume 17 (pp. 412–432).
- Freeman, H. (1974). Computer Processing of Line-Drawing Images. Computing Surveys, 6(1), 57–97.
- Ganzberger, M. J., Rovner, R. M., Gillies, A. M., & Hepp, D. J. (1994). Matching database record to handwritten text. In Proc. of the SPIE symposium on electronic imaging science and technology (Document Recognition), San Jose, CA, volume 2181 (pp. 66–75).
- Gilloux, M., Leroux, M., & Bertille, J.-M. (1993). Strategies for Handwritten Words Recognition Using Hidden Markov Models. In International Conference on Document Analysis and Recognition, Tsukuba Science City, Japan (pp. 299–304).
- Govindan, V. K. & Shivaprasad, A. P. (1992). Character Recognition A Review. Pattern Recognition, 23(7), 671–683.
- Govindaraju. V., Shekhawat, A., & Srihari, S. N. (1993). Interpretation of Handwritten Addresses in US Mail Stream. In *The third International Workshop on Frontiers in Handwriting Recognition, Buffalo, New York* (pp. 197–206).
- Govindaraju, V., Srihari, R. K., & Srihari, S. N. (1994). Handwritten Text Recognition. In Proceedings of the conference on Document Analysis System, Kaiserslautern, Germany (pp. 157–171).

- Govindaraju, V., Wang, D., & Srihari, S. N. (1992). Using Temporal Information In Off-Line Word Recognition. In Fifth Advanced Technology Conference, Washington D.C. (pp. 529–543).
- Graf, H. P., Burges, C., Cosatto, E., & Nohl, C. R. (1995). Analysis of Complex and Noisy Check Images. In *Proceedings of IEEE International Conference on Image Processing, Washington D.C.*, volume III (pp. 316–319).
- Guillevic, D. & Suen, C. Y. (1995). Cursive Script Recognition applied to the Processing of Bank Cheques. In *Proceedings of the Third International Conference on Document Analysis and Recognition, Montreal, Canada* (pp. 11–14).
- Haffner, P. & Waibel, A. (1992). Multi-State Time Delay Neural Networks for Continuous Speech Recognition. In J. E. Moody, S. J. Hanson, & R. P. Lippmann (Eds.), Advances in Neural Information Processing System 4 (pp. 135–142). San Mateo, CA: Kaufmann.
- Harmon, L. D. (1972). Automatic Recognition of Print and Script. *Proceedings of the IEEE*, 60(10), 1165–1176.
- Harmon, L. D. (Aug. 24, 1962). Handwriting Reader Recognizes Whole Words. *Electronics*, 35, 29–31.
- Hepp, D. J. (1994). Recognition of handprinted and cursive words by finding feature correspondences. In *Proc. of the SPIE symposium on electronic imaging science and technology*, volume 2181.
- Hild, H. & Waibel, A. (1993). Connected Letter Recognition with a Multi-State Time Delay Neural Network. In S. J. Hanson, J. D. Cowan, & C. L. Giles (Eds.), Advances in Neural Information Processing System 5 (pp. 712–719). San Mateo, CA: Kaufmann.

- Hu, J., Brown, M. K., & Turin, W. (1994). Handwriting Recognition with Hidden Markov Models and Grammatical Constraints. In *The fourth International Workshop on Frontiers in Handwriting Recognition, Taipei, Taiwan* (pp. 195– 205).
- Hull, J. J. (1988). A Computational Theory of Visual Word Recognition. PhD thesis, Dept. of Computer Science, State University of New York at Buffalo.

John, S. S. (1992). Checking Checks. Technology Review, 95, 15-16.

- Kawatani, T. (1995). Character Recognition Performance Improvement Using Personal Handwriting Characteristics. In *Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada* (pp. 98–103).
- Kawatani, T. & Miyamoto, N. (1993). Verification of Personal Handwriting Characteristics for Numerals and Its Application to Recognition. *Pattern Recognition Letters*, 14, 335–343.
- Kim, G. & Govindaraju, V. (1995). Handwritten Word Recognition for Real-Time Applications. In Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada (pp. 24–27).
- Kimura, F., Shridhar, M., & Narasimhamurthi, N. (1993). Lexicon Directed Segmentation - Recognition Procedure for Unconstrained Handwritten Words. In *The third International Workshop on Frontiers in Handwriting Recognition, Buffalo, New York* (pp. 122–131).
- Kung, S. Y. (1993). *Digital Neural Networks*. Englewood Cliffs, N.J.: Prentice-Hall, Inc.
- Lee, C. & Rabiner, L. R. (1989). A Frame-Synchronous Network Search Algorithm for Connected Word Recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 37(11), 1649–1658.

- Liu, Y. & Srihari, S. N. (1994). Document image binarization based on texture analysis. In *Proc. of the SPIE symposium on electronic imaging science and technology*, volume 2181.
- Lu, C.-C. & Dunham, J. G. (1991). Highly Efficient Coding Schemes for Contour Lines Based on Chain Code Representations. *IEEE Transactions on Communications*, 39(10), 1511–1514.
- Madhvanath, S. & Govindaraju, V. (1993). Holistic Lexicon Reduction. In The third International Workshop on Frontiers in Handwriting Recognition, Buffalo, New York (pp. 71–81).
- Madhvanath, S., Govindaraju, V., Ramanaprasad, V., Lee, D.-S., & Srihari, S. N. (1995). Reading Handwritten US Census Forms. In Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada (pp. 82–85).
- Mahadevan, U. & Nagabushnam, R. C. (1995). Gap Metrics for Word Separation in Handwritten Lines. In Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada (pp. 124–127).
- Mantas, J. (1992). An Overview of Character Recognition Methodologies. *Pattern Recognition*, 19(6), 425–430.
- Marr, D. (1982). VISION A Computational Investigation into the Human Representation and Processing of Visual Information. W. H. Freeman and Company, San Francisco.
- Matan, O., C. J. C. Burges, Le Cun, Y., & Denker, J. S. (1992). Multi-Digit Recognition Using A Space Displacement Neural Network. In J. E. Moody, S. J. Hanson, & R. P. Lippmann (Eds.), *Advances in Neural Information Processing System 4* (pp. 488–495). San Mateo, CA: Kaufmann.

- Mermelstein, P. & Eden, M. (1964). Experiments on Computer Recognition of Connected Handwritten Words. *Information and Control*, 7, 255–270.
- Moreau, J.-V., Plessis, B., Bougeios, O., & Plagnaud, J.-L. (1991). A Postal Check Reading System. In Proceedings of the First International Conference on Document Analysis and Recognition, St. Malo, France (pp. 758–766).
- Morgan, D. P. & Scofield, C. L. (1991). Neural Networks and Speech Processing. Kluwer Academic Publishers.
- Morgan, N. & Bourlard, H. (1995). Continuous Speech Recognition. *IEEE Signal Processing Magazine*, 12(3), 25–42.
- Mori, S., Suen, C. Y., & Yamamoto, K. (1992). Historical Review of OCR Research and Development. *Proceedings of the IEEE*, 80(7), 1029–1058.
- Myers, C. S. & Rabiner, L. R. (1981). A Level Building Dynamic Time Warping Algorithm for Connected Word Recognition. *IEEE Transactions on Acoustics*, *Speech, and Signal Processing*, 29(2), 284–297.
- Ney, H. (1984). The use of a One-Stage Dynamic Programming Algorithm for Connected Word Recognition. IEEE Transactions on Acoustics, Speech, and Signal Processing, ASSP-32(2), 263–271.
- Ney, H. (1992). A Comparative Study of Two Search Strategies for Connected Word Recognition: Dynamic Programming and Heuristic Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(5), 586–595.
- Otsu, N. (1979). A Threshold Selection Method from Gray-Level Histograms. *IEEE Transactions on Systems, Man, and Cybernetics*, SMC-9, 62–66.

- Palumbo, P. W., Swaminathan, P., & Srihari, S. N. (1986). Document image binarization: Evaluation of algorithms. In *Proc. of the SPIE symposium on applications of digital image processing IX*, volume 697.
- Paquet, T. & Lecourtier, Y. (1991). Handwriting Recognition: Application on Bank Cheques. In Proceedings of the First International Conference on Document Analysis and Recognition, St. Malo, France (pp. 749–757).
- Pavlidis, T. (1978). A Minimum Storage Boundary Tracing Algorithm and Its Application to Automatic Inspection. *IEEE Transactions on System, Man, and Cybernetics*, SMC-8(1), 66–69.
- Pavlidis, T. (1982). Algorithms for Graphics and Image Processing. Rockville, MD: Computer Science Press.
- Pavlidis, T. (1993). Recognition of printed text under realistic conditions. *Pattern Recognition Letters*, 14, 317–326.
- Rabiner, L. R. (1989). A Tutorial on Hidden Markov Models and Selected Applications in Speech Recognition. *Proceedings of the IEEE*, (pp. 257–286).
- Rabiner, L. R. & Juang, B. (1993). *Fundamentals of Speech Recognition*. Englewood Cliffs, New Jersey: Prentice Hall Signal Processing Series.
- Rabiner, L. R. & Juang, B. H. (1986). An Introduction to Hidden Markov Models. IEEE ASSP Magazine, (pp. 4–16).
- Renals, S., Morgan, N., Bourlard, H., Franco, H., & Cohen, M. (1992). Connectionist Optimisation of Tied Mixture Hidden Markov Models. In J. E. Moody, S. J. Hanson, & R. P. Lippmann (Eds.), *Advances in Neural Information Processing System 4* (pp. 167–174). San Mateo, CA: Kaufmann.

- Sakoe, H. (1979). Two-level DP-matching A Dynamic Programming-Based pattern matching algorithm for connected word recognition. *IEEE Transactions on Acoustics, Speech, and Signal Processing,* 27, 588–595.
- Sayre, K. M. (1973). Machine Recognition of Handwritten Words: A Project Report. Pattern Recognition, 5, 213–228.
- Schürmann, J., Bartneck, N., Bayer, T., Franke, J., Mandler, E., & Oberlander, M. (1992). Document Analysis - From Pixels to Contents. *Proceedings of the IEEE*, 80(7), 1101–1119.
- Seni, G. (1995). Large Vocabulary Recognition of On-line Handwritten Cursive Words. PhD thesis, Dept. of Computer Science, State University of New York at Buffalo.
- Seni, G. & Cohen, E. (1994). External Word Segmentation of Off-Line Handwritten Text Lines. *Pattern Recognition*, 27(1), 41–52.
- Silverman, H. F. & Morgan, D. P. (1990). The Application of Dynamic Programming to Connected Speech Recognition. *IEEE ASSP Magazine*, (pp. 6–25).
- Simon, J. C. (1992). Off-line Cursive Word Recognition. *Proceedings of the IEEE*, 80(7), 1150–1161.
- Simon, J. C., Baret, O., & Gorski, N. (1994). 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).
- Spinard, R. J. (1965). Machine Recognition of Hand Printing. Information and Control, 8, 124–142.
- Srihari, S. N. (1992). High-Performance Reading Machines. *Proceedings of the IEEE*, 80(7), 1121–1132.

- Srihari, S. N. (1993). From Pixels to Paragraph: the Use of Contextual Models in Text Recognition. In International Conference on Document Analysis and Recognition, Tsukuba Science City, Japan (pp. 416–423).
- Srihari, S. N., Govindaraju, V., & Favata, J. (1995a). Unconstrained Handwritten Text Recognition. In Symposium on Document Image Understanding Technology, Oct. 24-25, Maryland (pp. 226–235).
- Srihari, S. N., Shin, Y.-C., Ramanaprasad, V., & Lee, D.-S. (1995b). Name and Address Block Reader System for Tax Form Processing. In Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada (pp. 5–10).
- Suen, C. Y. (1986). Character Recognition by Computer and Application. In T. Y. Young & K.-S. Fu (Eds.), Handbook of Pattern Recognition and Image Processing (pp. 569–586). Academic Press.
- Suen, C. Y., Berthod, M., & Mori, S. (1980). Automatic Recognition of Handprinted Characters - The State of the Art. *Proceedings of the IEEE*, 68(4), 469–487.
- Tappert, C. C. (1982). Cursive Script Recognition by Elastic Matching. *IBM Journal of Research and Development*, 26(6), 765–771.
- Tappert, C. C. (1986). An Adaptive System for Handwriting Recognition. In H. Kao, G. Galen, & R. Hoosain (Eds.), *Graphonomics: Contemporary Research in Handwriting* (pp. 185–198). Elsevier Science Publishers B. V. (North-Holland).
- Tappert, C. C. (1991). Speed, Accuracy, and Flexibility Trade-Offs in On-Line Character Recognition. In P. Wang (Ed.), Character and Handwriting Recognition: Expanding Frontiers (pp. 79–94). World Scientific.

- Tappert, C. C., Suen, C. Y., & Wakahra, T. (1990). The State of the Art in On-line Handwritten Recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence, 12(8), 787–808.
- Waibel, A., Hanazawa, T., Hinton, G., Shikano, K., & Lang, K. J. (1989). Phoneme Recognition Using Time-Delay Neural Networks. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 37(3), 328–339.
- Wakahara, T. (1993). Toward Robust Handwritten Character Recognition. Pattern Recognition Letters, 14, 345–354.
- Werbos, P. J. (1991). Links Between Artificial Neural Networks (ANN) and Statistical Pattern Recognition. In I. K. Sethi & A. K. Jain (Eds.), Artificial Neural Networks and Statistical Pattern Recognition Old and New Connections (pp. 11–31). Elsevier Science Publishers Co.
- Yuan, J., Tang, Y. Y., & Suen, C. Y. (1995). Four Directional Adjacency Graphs (FDAG) and Their Application in Locating Fields in Forms. In *Third International Conference on Document Analysis and Recognition (ICDAR-95), Montreal, Canada* (pp. 752–755).
- Zimmerman, K. (1995). Banks Shoot for a Customer Service Impact with Imaging. Imaging Magazine, 4.