

**HIERARCHICAL CHARACTER RECOGNITION
AND ITS USE IN HANDWRITTEN WORD/PHRASE
RECOGNITION**

by

Jaehwa Park

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

November 1999

© Copyright 1999

by

Jaehwa Park

COMMITTEE

Chairman	Sargur N. Srihari, Ph. D. Distinguished Professor, Department of Computer Science and Engineering, State University of New York at Buffalo
Members	Peter D. Scott, Ph. D. Associate Professor, Department of Computer Science and Engineering, State University of New York at Buffalo Venu Govindaraju, Ph. D. Research Associate Professor, Department of Computer Science and Engineering, State University of New York at Buffalo
Outsider Reader	Daniel P. Lopresti, Ph. D. Bell Laboratories, Lucent Technologies, Inc. Murray Hills, NJ

This dissertation was defended on Wednesday, November 10, 1999

*To my parents
and
my wife, Jongwon*

ACKNOWLEDGMENT

I am deeply grateful to my advisor, Professor Sargur N. Srihari, for his encouragement and guidance, and for providing the opportunity for conducting this research at CEDAR.

I would like to thank Professor Peter D. Scott for his valuable comments and for serving as a committee member.

My special thanks go to Professor Venu Govindaraju for his friendship, support and help as a committee member and a colleague.

I would also thank my outside reader, Dr. Dan Lopresti, for sharing his valuable time in spite of his busy schedule and for his kind comments.

I would like to send special thanks to my colleagues, Dr. Ajay Shekhawat and Evelyn Kleinberg for their friendship, bright suggestions, and valuable help. I am grateful to CEDAR staff members, especially to Eugenia Smith, who is always ready to solve various administrative matters. I would like to extend my thanks to all members of CEDAR. I received a lot of help, comments, and suggestions from all the members of CEDAR during my research.

Finally, I would like to thank all my family members for their support, encouragement and prayers, especially my parents for making it possible, my two sons for giving me delight, joy and hope, and my lovely wife, Jongwon, for her sacrifice.

ABSTRACT

Off-line handwritten word/phrase recognition systems generally have monotonically cascaded architecture through several processing steps. In these architectures, the recognition engine follows a static model with a fixed feature space. Built-in resources are exhaustively used at each stage of the serial engine regardless of input complexity. However, the perception of a word is fundamentally an interactive process in human cognitive models; i.e., both conceptually-driven and data-driven processes work simultaneously and processing is recursively adapted to what is perceived in time. For optimality and efficiency, a system that achieves maximum performance with minimum processing effort is desirable, and autonomous adaptation to input is one of the solutions to the goal.

A recursive computational model for handwritten character/word/phrase recognition that has some similarities to the human cognitive approach is proposed. Two concepts, (i) altering recognition action using feedback and (ii) evaluating and regulating terminating conditions actively, are introduced for dynamic and interactive recognition. A hierarchical classification method is presented with dynamic usage of hierarchical feature space that preserves the benefits of the multi-resolution model. A lexicon-driven word recognizer which operates dynamically and has different degrees of classification ability is also presented. A concept of lexicon complexity derived from "matching transform distance" is utilized as a decision metric, which measures the difficulty of the given lexicon set with respect to the classification ability of the character recognizer.

Recognition, decision making and recursive updating from the closed loop architecture is designed for a recursive recognition scenario. This proposed model

recursively enhances the degree of classification until decision conditions are satisfied. A prime stroke analysis which detects the writer's spacing style is formulated to extend word recognition to phrase recognition system. The prime stroke period is used as a metric to limit the combination of strokes for word recognition and to generate phrase hypotheses. The proposed models achieve 98% and 96% of top choice correct rates in character and word recognition performance tested on NIST digit set and city name word images collected from the USPS mail stream respectively. Application of our phrase recognition to the recognition of street lines in USPS mail pieces resulted in a 10% improvement comparing to the performance of a word mode recognition.

Contents

Committee	iii
Acknowledgment	v
Abstract	vi
1 Introduction	1
1.1 Definitions	2
1.2 Motivation	4
1.3 Problem Issues	7
1.4 Proposed Model	8
1.5 Dissertation Organization	12
2 Related Studies	14
2.1 Character Recognition	14
2.1.1 Feature Extraction and Selection	15
2.1.2 Classification Method	16
2.1.3 Combination of Classifiers	18
2.2 Word Recognition	19
2.2.1 Analytical Word Recognition	19

2.2.2	Holistic Word Recognition	21
2.2.3	Word Recognition Combination	22
2.3	Phrase Recognition	22
2.3.1	Word Segmentation	23
2.3.2	Phrase Recognition Models	25
2.4	Cognitive Models	26
2.4.1	Word Perception Model	26
2.4.2	Epistemic Utility	28
3	Hierarchical Character Recognition	30
3.1	Hierarchical Classification	33
3.1.1	Probability in vector quantized hyper feature space	34
3.1.2	Simulation of multi-resolution	37
3.1.3	Efficient processing of the recursive tree	39
3.2	Implementation	43
3.2.1	Image Representation	43
3.2.2	Normalization	46
3.2.3	Sub-images	46
3.2.4	Features	47
3.2.5	Training	50
3.3	Experiment	51
3.3.1	Recognition Performance	53
3.3.2	Dynamic Characteristics	56
3.4	Discussion	60
4	Stroke Analysis	63

4.1	Segmentation Model	64
4.2	Normalization	68
4.3	Stroke Segmentation	71
4.4	Word Break Detection	75
4.4.1	Primitive Classification	76
4.4.2	Estimating Prime Frequency	79
4.5	Experiment	80
5	Recursive Word Recognition	82
5.1	Applying Utility Theory	84
5.1.1	Revision of Utility	86
5.2	Simulating Multi-Resolution in A Distance Metric	87
5.2.1	Multi-resolution Feature Space	88
5.2.2	Distance in A Recursive Tree	90
5.2.3	Word Recognition Accuracy	93
5.3	Lexicon Complexity	94
5.3.1	Edit Cost in A Segment Frame	97
5.3.2	Matching Distance	99
5.4	Implementation	101
5.4.1	Stroke segmentation	102
5.4.2	Features	103
5.4.3	Character Matching	106
5.4.4	Word Matching	108
5.5	Experiment	111
6	Dynamic Phrase Recognition	118

6.1	Prescreening	119
6.2	Hypotheses Generation	122
6.3	Hypotheses Verification	128
6.4	Experiment	129
7	Conclusions	133
7.1	Summary	133
7.2	Contribution	135
	Bibliography	136

List of Tables

3.1	Experiment types and size of training and testing sets	52
3.2	Recognition results using a quin tree in maximum depth of 4	53
3.3	Performance comparing by changing sub-image generation rule using quad and quin trees.	54
3.4	Performance in various classification methods using setting of EXP 3; The proposed method operated on a dynamic mode of $\log(\kappa) = 40$ and feature vectors for the neural network of 168 sized input layer were generated using a quad tree of depth 3 (1+4+16)	55
4.1	Segmentation Performance	81
6.1	Estimated number of characters	124
6.2	Raw lexicon	125
6.3	Expanded lexicon	125
6.4	Lexeme hypotheses - the symbol # indicates remaining image segments corresponding to additional informations	126
6.5	Phrase hypotheses - the number within parenthesis indicates index of lexeme hypothesis	127
6.6	Information of experiment image set	129
6.7	Experiment environment - all numbers are average values	130

List of Figures

1.1	Example of handwritten (a) characters, (b) words, and (c) phrases . .	3
1.2	General architecture of handwritten word/phrase recognition system	5
1.3	Handwritten words in different writing style	6
1.4	Multi-level interactive processing model	9
1.5	A practical model of recursive recognition system	12
3.1	Block diagram of hierarchical classifier	31
3.2	An example of using sub-images	33
3.3	Processing steps of class probability sampling	35
3.4	Mapping from binary feature vectors to vector codes when $N_x = 3$.	36
3.5	Sub-images generated by a quin tree division rule	38
3.6	An example of confidence update recursion; (a) starting from root, (b) after 1 recursion pass, (c) after 2 recursion passes	40
3.7	Methodology of simulating a multi-resolution feature space by re- cursive processing of <i>selected</i> sub-images	44
3.8	Contour representation; (a) contour slope division, (b) vertical weight map and (c) horizontal weight map	45

3.9	Contour representation and preprocessing; (a) input image, (b) pixel based contour representation, (c) piece-wise linearized contour and (d) piece-wise linear contour representation after normalization . . .	47
3.10	Character image and sub images of quin tree (class A) and quad tree (class 9) : (a) original image, (b) piecewise linearized contour, sub-regions at (c) layer 1, (d) layer 2, (e) layer 3, and (f) pseudo re-constructed contour image from center points of final layer	48
3.11	Gradient feature generation; (a) gradient division and (b) gradient angle measure	50
3.12	Dynamic recognition performance. (a) top choice correct rate and (b) second choice correct rate	57
3.13	Average resource usage. (a) node usage in each searching mode and (b) normalized node usage vs. performance	58
3.14	Performance by class	59
3.15	Confidence and separability changes - in EXP 1 with quin tree. (solid line shows a case of successful and dashed line shows a case of failure that true answer is not top ranked)	62
4.1	Examples of various spacing style: (a)-(c) character based spacing, (d)-(f) stroke based spacing.	65
4.2	Prime spatial period (a) and types of stroke segments (b)	67
4.3	Normalization steps	70
4.4	Normalization processing example, (a) a street line image, (b) critical points and piecewise linearized contours, (c) vertical contour components (thick lines) and filtered local minimum points (square dots), (d) slant and skew normalized contours	72
4.5	Stroke segmentation steps	73

4.6	Stroke segmentation processing example: (a) normalized contours, (b) stroke width sensing, (c) extracted convexities and concavities (d) extracted ligature zone (e) final segmentation points (f) segmented primitives.	74
4.7	Word segmentation steps	75
4.8	Primitive classification	76
4.9	Word segmentation processing example: (a) slant and skew normalized contours, (b) segmented primitives, (c) stroke interval estimation, (d) generated word segments	78
5.1	Recursive recognition	83
5.2	An example of lexicon complexity	95
5.3	Stroke segmentation: (a) normalized contours, (b) segmented primitive array and (c) prime stroke period estimation	102
5.4	Sub-image division in a grouped character segment	104
5.5	Character matching	105
5.6	Word matching based on dynamic programming with nullity	110
5.7	Accuracy support changes: (a) success case and (b) failure case	113
5.8	Adjusted distance changes: (a) utility values and (b) adjusted distance	114
5.9	Recognition performance with and without liability	115
5.10	Dynamic recognition characteristics - comparison between recursive system with active searching and non-recursive system with passive adding of sub-nodes	116
5.11	Dynamic recognition characteristics - recognition performance changes by various acceptance criterion	117
6.1	Dynamic matching phrase recognition	120

6.2	Word image segments	124
6.3	Performance of phrase recognition algorithm; Top choice correct rate	131
6.4	Reject to error characteristic of phrase recognition	132

Chapter 1

Introduction

The main object of Handwriting Recognition is interpretation of data which describes handwritten objects. On-line handwriting recognition deals with a data stream which is coming from a transducer while the user writes, and off-line handwriting recognition deals with a data set which has been obtained from a scanned handwritten document. The goal of handwriting recognition is to interpret the contents of the data and to generate a description of that interpretation in the desired format.

For the past few decades, intensive research has been done to solve this problem in related areas such as image processing, pattern recognition, cognitive science etc. Various approaches, system architectures and methodologies have been proposed to deal with application diversity. To date, challenging problems are being encountered and solutions to these are broadly targeted to improve accuracy and efficiency.

Remarkable progress has been made in practical automatic recognition and processing systems. Interpretation of addresses on mail pieces [1] [2] [7], recognizing address blocks on tax forms [5] [8], reading amounts from bank checks [3]

[4] and reading general text [6] [84] are some examples of the applications.

Conventional methodology to develop off-line handwriting recognition systems is based on monotonically cascaded architectures through several processing steps. Static computational models and rigid usage of resources are the basic approaches over each step. Global performance maximization within acceptable cost is the key idea to reach the optimal status. For optimality and efficiency, a system that achieves maximum performance with minimum processing effort is desirable. However, trade-off between efficiency and accuracy is inevitable when a system targeted for real application is designed.

In human cognition, perception is fundamentally an interactive process i.e., both conceptually-driven and data-driven processes work simultaneously and processing is recursively adapted to what is perceived in time. Following the human perception model, a paradigm of autonomous adaptation to input is one of the solutions to reach the optimality goal.

With this motivation, a computationally efficient solution to the problem of recognizing handwritten characters based on a hierarchical process that has some similarities to the human cognitive process is proposed. This method has been applied in recognizing handwritten words and phrases.

1.1 Definitions

Here we describe some important terminology and the system input and output types. Units of meaningful handwriting can be categorized into the following four types: character, word, phrase and sentence [6].

The dictionary definitions of the four units are: *character*: a graphic symbol as a hieroglyph or alphabet letter or digit used in writing or printing; *word*: a written or

0	1	2	3	4	5	6	7	8	9	#		
A	B	C	D	E	F	G	H	I	J	K	L	M
N	O	P	Q	R	S	T	U	V	W	X	Y	Z

(a)

Rochester
Buffalo
Zimmerman
BANK

(b)

Rossmore Trust Co
Englewood Post Office
Overland H.S
King Soopers w/c #827
ADP

(c)

Figure 1.1: Example of handwritten (a) characters, (b) words, and (c) phrases

printed character or combination of characters representing a spoken word; *phrase*: a word or group of words forming a syntactic constituent with a single grammatical function; *sentence*: a word, clause, or phrase or a group of clauses or phrases forming a syntactic unit. Handwritten character, word, phrase and sentence are defined as those that are written by a human hand.

There are two different system approaches, namely, *on-line* and *off-line*. An on-line system does the recognition process at same time while the units are captured directly from a writer. This approach requires a devices to transduce the units. In off-line application, a handwritten document is typically scanned from paper and made available in the form of a pixelized image to a recognition processor. The major difference is that on-line data has time-sequence contextual information but off-line data does not. This difference generates a significant divergence in processing architectures and methods.

1.2 Motivation

General processes of off-line handwritten word recognition are *preprocessing*, *recognition*, *postprocessing* and *decision making* as shown in Figure 1.2. *Preprocessing* is primarily related to image processing operations such as normalization to remove irregularities of handwriting. *Recognition* is the application of classification algorithms. *Postprocessing* uses independent contextual information to enhance recognition results. *Decision making* chooses the final answers to reach a desired level of performance.

Most system architectures are monotonically cascaded through these steps. Starting from a scanned image to a transcribed output, the steps are independent and the processes are sequential. Recognition engines are operated in a static

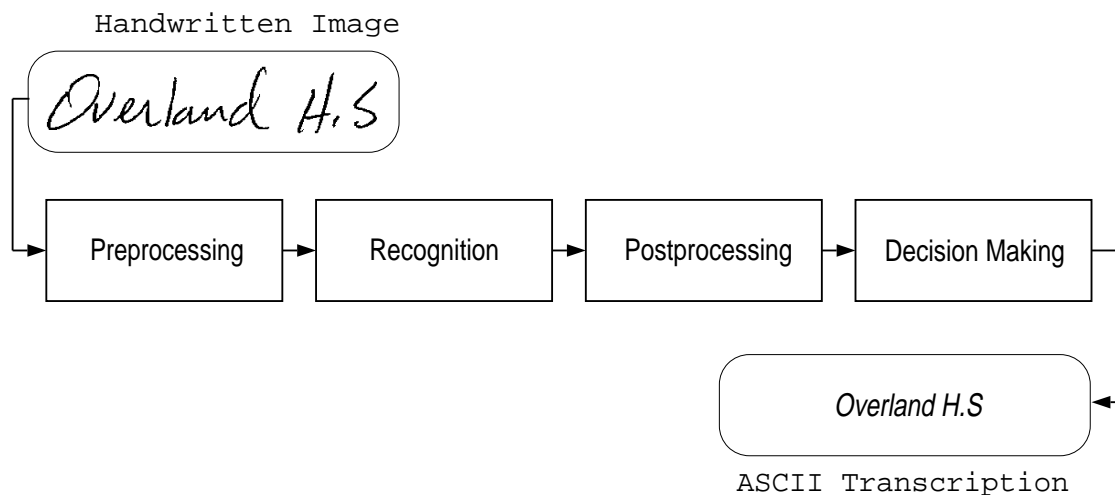


Figure 1.2: General architecture of handwritten word/phrase recognition system

model where a uniform classification algorithm is applied on all classes with a fixed feature space that is optimally selected from a training set. Also, resources are used in inflexibility at each stage of the serial engine without consideration of environmental changes. Decision making is designed independently to provide one and only one decision output which is represented as the true answer. The idea is to maximize the global performance by maximization at each stage. However, such a paradigm is prone to being non-optimal because individual stages usually do not have all the information.

Cognitively speaking, humans do not follow such monotonic paradigms. They can alter their course of action dynamically and select the utilities relative to their goals. The word perception model in the cognitive approach is assumed to be fundamentally an interactive process [85]. That is, top-down or conceptually

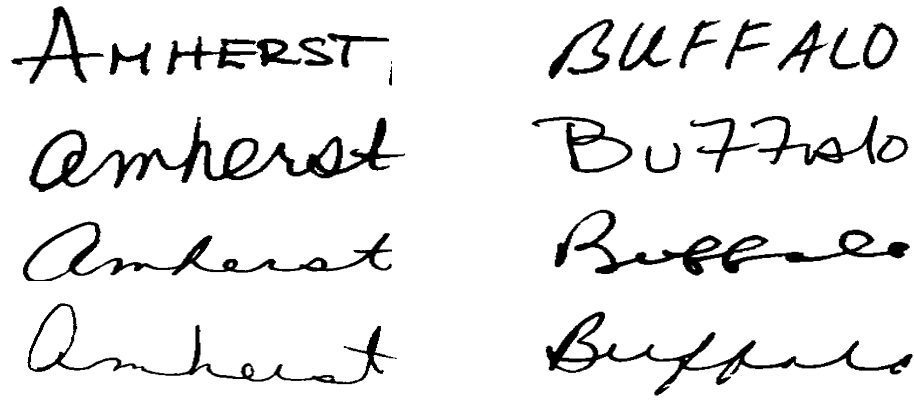


Figure 1.3: Handwritten words in different writing style

driven processing works simultaneously and in conjunction with bottom-up or data driven processing to provide a sort of multiplicity of constraints that jointly determine what we perceive.

As shown in Figure 1.3, required accuracy of identifying visual features which is adequate enough to classify each character and to finalize word recognition depends on input quality. The images shown in the last line of Figure 1.3 require higher resolution to detect visual evidence and more attention to accept the words, compared to the images above. Also, acceptance of a recognition result depends on difficulty in the classification and similarity of the lexicon (such as *edit distance*). For instance, the last image of 'Amherst' can only be acceptable when there is no ambiguity in position of 'r' in lexicon, since the visual evidence of isolated character at position 'r' does not correspond much to the shape 'r'. If there are two lexicon entries 'Amherst' and 'Amheist', we can not accept 'Amherst' as a final decision.

In order to extend conventional recognition systems to dynamic and recursive systems to simulate human cognitive models, two concepts must be introduced: (i) altering action using feedback and environment and (ii) evaluating and regulating terminal conditions actively.

1.3 Problem Issues

The handwriting recognition task involves several subtasks such as separation of the image into meaningful units, recognition of the separated units, and decision making based on the recognition results; in addition, the task requires global system organization to maximize system performance. There are many problems that arise from these subtasks that should be solved to build an efficient and optimal system. In this section, following the above mentioned motivation, several problems related to character, word and phrase recognition processes are addressed.

Character recognition is basically related to recognition of pre-isolated character images and most research has focused on finding the best feature set and classification method in a static architecture. A method which achieves maximum separation among classes in a selected training set that is closest to the application is chosen as the optimal recognizer. A decision step follows recognition to accept recognition result within the desired performance. In this sequential and unidirectional process, feature extraction is passive and classification involves inflexible resource usage to provide the best results regardless of image quality. Also, the decision algorithm usually has a uniform criterion for acceptance. This static approach is lacking in adaptability to the input diversity and injection of dynamic operation into the processing flow is difficult.

In architectures for word recognition, several intermediate processing layers

such as feature, symbol and word level are established and methods to find the best match to a word, progressing upward from bottom to top layer, are used as the recognition process. Recognition results are provided in the form of ranked arrays of character symbols or ranked lists of lexicon entries with corresponding recognition confidences. Generally the confidence is an absolute measurement of the word image within built-in resources of the recognizer. The metric of word confidence is not adjusted to the possible character subsets nor lexicon subsets.

However, accepting a word depends not only on the absolute value of visual evidence but also on similarity of lexical entries. A regulation strategy which reflects the complexity of a lexicon is needed for optimal decision making. Contextual evidence extracted from the image other than for symbol matching (such as spacing style) can sometimes reduce computational complexity.

Research into phrase recognition involves the subtopics of word segmentation, word recognition and postprocessing. For word segmentation before word recognition, several geometric distance measuring methods are introduced. But separating a handwritten phrase into words is not always successful because of flourishes and variations in writing style. Moreover, the ensuing recognition steps do not have a way of compensating for word segmentation error. And finally issues such as incomplete lexicons and unwanted objects (such as additional writings or noise) in the input image must be tackled.

1.4 Proposed Model

Recently a general computational framework of parallel distributed processing has found use in developing theories of many aspects of neural processing, perception and cognition. The core concept is that the outcome may result from interactions

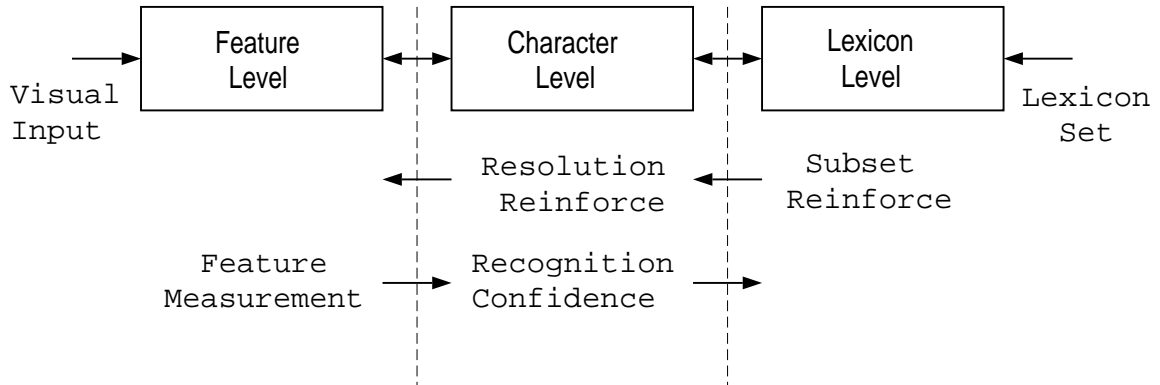


Figure 1.4: Multi-level interactive processing model

among a large network of simple processing elements. Elements become activated and the pattern of activation through the network changes by means of interaction between elements.

This computational framework has been applied to the problem of word recognition as the *time course model of word recognition* [85]. The core idea is that various facts about word recognition can be explained by considering the time course of multi-layer code activation; perception of a word results from an interactive process of nodes representing several different types of information such as feature symbol, character, and word from initial excitation of visual information.

The time course model of word recognition was implemented in a computational model of an exhaustively connected network through a small set of features and lexicon entries. The experimental results show the tracking ability of the interactive network model to the desired outputs. But several problems remain for practical implementation in automated handwritten word recognition. The major

problems are that (i) isolation of characters is not always guaranteed in handwritten words, (ii) the required feature dimension for successful handwritten character recognition is relatively large and varying to situation, and (iii) the interactive network becomes very large if the size of lexicon entries is large.

To solve these problems, we propose an interactive model as shown in Figure 1.4 which is simplified but has feedback in categorized channels from the model of [85]. A processing level hierarchy is established as in that of the time course model but interactions between levels are limited. Forward paths and reinforcement paths are established instead of exhaustive connection in order to simplify excitatory and inhibitory connections of the cognitive time course model. Feature measurements and character recognition results are generated as outputs in forward paths similarly to conventional word recognition methods. Character subset and resolution limitations are influenced from upper to lower level to excite required evidences as reinforcement paths.

The goal of decision making is to define and evaluate criteria that meet desired performance. The interactive process in a system should accompany a decision making module to reach a minimal cost operation. A theory of cognitive decision making, *epistemic utility theory* is one of the computational approaches in this area [88] [89]. It assumes that the act of making decisions involves the conflict between the desire to obtain new evidences and the desire to avoid errors. It provides a computational model of quantifying the value of avoiding an error and the value of acquiring new information.

For practical computational implementation, we propose a recursive recognition architecture as shown in Figure 1.5. The flow of processing steps is similar to that of a conventional word recognizer. But it is combined with a *loose decision maker* which adapts the concept of epistemic utility theory. The loose decision

maker rejects false hypotheses instead of accepting qualified top answer. A reinforcement step (*Update Utility* in Figure 1.5) is added as a feedback path from decision making module to the recognizer.

In order to achieve successful recognition with minimum effort, a hierarchical character recognition method is developed. The hierarchical character recognition simulates multi-resolution classification from global to detailed feature extraction. The resolution level can be controlled dynamically by the decision maker or a high level word recognizer by adjusting the usage of hierarchical feature space.

A decision metric is derived from the cognitive decision model described in epistemic utility theory of [90]. It is formulated based on two concepts of (i) evaluation of absolute confidence of projection on knowledge space and (ii) evaluation of relative credibility of information regardless of the absolute values. The first term is derived from the confidence (or distance) of character recognizer. This term provides the absolute opinion of character shape by feature measurement on built-in character knowledge. The second term is defined as *lexicon complexity* which represents the relative importance of each character in a lexicon entry against the other lexicon entries in a lexicon by means of classification difficulty.

During the feedback step, the recognition results are analyzed in term of the decision metric by the decision maker and the reinforcement to the recognizer is updated by the selected (or narrowed) hypotheses. This feedback helps the recognizer select the most efficient operation to reach terminating conditions.

This recursive model can not process *top-down* and *bottom-up* operations simultaneously like the cognitive network. But it imitates the tracking ability of the cognitive model through time delayed feedback in iterations. The difference is that delayed feedback of current recognition influences the subsequent recognition steps.

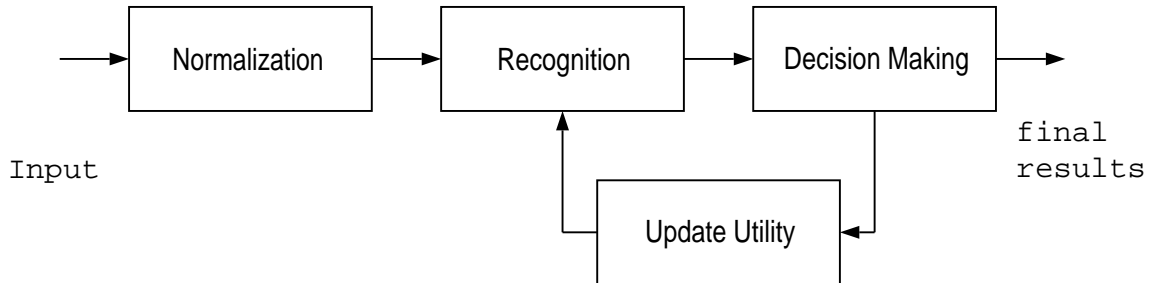


Figure 1.5: A practical model of recursive recognition system

1.5 Dissertation Organization

This dissertation is divided into two sections. The first section, comprised of Chapter 1 and Chapter 2 describes definitions, motivation, overview of proposed system and review of previous works related to the proposed model. The second section, Chapter 3 through Chapter 6, provides detailed description of the proposed model and experimental results.

In Chapter 2, previous work in character, word, and phrase recognition (including word segmentation techniques) is briefly reviewed. In addition, basic concepts of the cognitive word perception model and decision making theory are summarized.

Chapter 3 describes the basic ideas behind the computational model of hierarchical classification method. In this chapter, the mathematics for of hierarchical character recognition are based on statistical scope for an independent character recognizer. (Another hierarchical classification formulation based on a distance

metric is described in Chapter 5 for word recognition.) Various character recognition experimental results are also provided to support the validity of the hierarchical classification method.

Chapter 4 describes work related to the topic of unit segmentation. A new method of capturing writer's spacing style, *prime stroke analysis*, is introduced. Computational models of detecting gaps and generating gap confidence between characters/words are formulated.

Chapter 5 describes a basic mathematical framework for the recursive word recognizer using *epistemic utility theory* and *lexicon complexity*. The implementation details and experimental results of applying to USPS city/state name recognition are provided.

Chapter 6 describes a new method of phrase recognition by applying the segmentation method of Chapters 4. A phrase recognition scenario containing pre-screening, phrase hypothesis generation and verification is explained followed by a street name recognition application with experimental results.

Finally, a summary of the work and concluding remarks are given in Chapter 7.

Chapter 2

Related Studies

In the last four decades, handwriting recognition has been a very active area of research. Previous work on this topic can be divided into four major areas, depending on whether the recognition units are characters, words, phrases, or longer bodies of text. Some of the previous work and problems are discussed in this chapter through the division. Cognitive research on human word perception and decision making that could be applied to the development of our recursive recognition model is also reviewed.

2.1 Character Recognition

Character recognition deals with the problem of classifying pre-isolated character images within a given alphabet set. Useful reviews are found in [11] [12] [13] [14] [15] [16] [17]. Most researchers have adopted the classical pattern recognition approach in which image pre-processing is followed by feature extraction and classification. There have been many subtopics, but these may be roughly divided into the following categories: feature extraction and selection, classification method,

and combination of classifiers.

2.1.1 Feature Extraction and Selection

Feature extraction is an important step in achieving good performance for a character recognizer. Extracted features must be invariant to the distortions and variations that can be expected in a specific application. The size of the feature set is also important in order to avoid a phenomenon called the dimensionality problem [98].

The type of selected feature can determine the nature and output of the preprocessing steps; i.e., the decision to use gray-scale versus binary image, filled representation or contour, thinned skeletons versus full-stroke images depends on the nature of the features to be extracted. Two important characteristics for describing feature sets are global versus local and symbolic versus reconstructive.

Feature extraction methods using topological features can generally reconstruct the image from the feature set. Features are obtained from coefficients of various orthogonal decomposition methods by the representation properties of the image data. *Fourier descriptors* [20] [21] [22] [23], *geometric moment invariants* [25] [26] [27] [28], *Zernike moments* [24], *Wavelet descriptors* [29] [30] are the examples of reconstructive feature extraction methods. Reconstructive features generally have a multi-resolution property within the feature composition. Usually the order of the base orthogonal expansion function indicates the resolution levels. A feature vector obtained from the coefficients of the expansion base function has global as well as detailed description at the same time.

In symbolic feature extractions, the feature measurements are usually transformed into symbolic representations of geometric primitives such as line segments, convexity, concavity, convex polygons, projections etc. The base image is not perfectly reconstructible from the extracted symbolized feature set, but these methods are simple in computation and have a reasonable discrimination power. *gradient-based features* [18], *projection histograms* [19] and *gradient structural concavity* [36] are the examples.

The resolution in symbolic features depends on the area from which the feature extractions are performed. *Zoning*, i.e. extracting features from segmented grids or sub-divisions, is one approach to adjust the resolution levels of feature presentation. *Fixed grids* [18] [42] or *variable grids* [36] are adapted for zoning applications. In order to simulate multi-resolution, the feature vectors have heterogeneous composition from various levels of zoning and extraction methods.

In practice, selection of the best features for a given application is a challenging task. Several references present solutions to the feature selection problem, such as an algorithm which selects the best subset from a pre-existing feature set to maximize classification through various driving functions [31] [33], and an algorithm which automates feature generation using a random generator based on information and orthogonality measures [32].

2.1.2 Classification Method

Most recognizers have adopted classical pattern classification methods. Major approaches are *statistical based*, *structural analysis*, *template matching*, and *neural network* approaches [9] [10]. Significant progress has been made in these classification methods but more work is required to match human performance.

Recently, recognizers using multi-resolution features have been found to be more advantageous in classification than conventional methods that work with features at a single scale. One multi-resolution recognizer is the Gradient Structural Concavity (GSC) which is described in [36]. It uses symbolic multi-resolution features. The Gradient of the image contour captures the local shape of a character. The Gradient features are extended to Structural features by encoding the relationships between strokes. Concavity features capture the global shape of characters. Features at the 3 levels, **G**, **S**, and **C** are combined, and a k-nearest neighbor classification method is used to produce a multi-resolution character recognizer. But the heterogeneous composition of the feature set makes it difficult to find an appropriate common distance metric that can be used by the classifier in the multi-resolution feature space.

Generative models based on iterative computation have also been proposed as a dynamic approach to handwritten character recognition. In [35] generative models that are built from deformable B-splines with *Gaussian ink generators* spaced along the length of the spline are introduced. The splines are adjusted using a novel elastic matching procedure based on an *expectation maximization* algorithm that maximizes the likelihood of the model generating the data. A Bayesian interpretation of the fitting process is adopted to yield a practical recognition algorithm. This model has additional advantages such as description of style parameters, recognition based segmentation, small size of template, and pre-normalization compared to conventional classification algorithms. However, this approach still has problems such as computational complexity because of the burden of additional features and iterative fitting process.

2.1.3 Combination of Classifiers

Some recent research has shown improved performance using a combination of several different classification algorithms. Many examples of such multiple classifier systems can be found in handwritten digit/character recognition publications. Parallel classifier combination methods have been extensively studied, including *adaptive voting principle*, *Bayesian formalism*, *Dempster-Shafer theory*, *neural network*, *statistical regressions* and *fuzzy integrals*.

The voting principle is the simplest method of combination where the top candidate from each classifier constitutes a single vote. The final decisions can be made by certain rules such as majority, weighted sums, etc [13] [52].

The Bayesian combination approach seeks to estimate the class posterior probabilities generated by classifiers. Based on different assumptions and different techniques for probability estimation, variations of the Bayes rules have been derived [46] [51]. In [45], an evidence aggregation using the *Behavior-Knowledge Space Method* which gathers the decisions obtained from individual classifiers derives the best final decision from the statistical point of view.

Dempster-Shafer theory in which the individual evidence are combined into the belief based on probability approach, has been applied for frameworks of classifier combination. The fundamental property of Dempster-Shafer theory to aggregate evidence using uncertainty reasoning, lends itself to classifier combination tasks [48] [51].

Artificial neural networks and regression approach provide a direct estimation of the posterior probability. A comparison of neural network performance to other combination methods for handwritten digit recognition has been given in [46] [49]. Neural network combinator compare favorably to other combination methods.

A classifier combination using logistic regression has been introduced in [47]. A simple regression model based on rank values is shown to be more effective than the voting of highest rank. Fuzzy integration of classifier reliability and objective evidence are also used for combination of classifiers [50].

Recently a dynamic classifier selection approach has been presented in [44]. Estimations of each individual classifier's local accuracy in small regions of the feature space surrounding an unknown test sample are used for the dynamic classifier selection algorithm.

2.2 Word Recognition

Off-line handwritten word recognition algorithms usually take two inputs: pre-separated word image and a lexicon representing possible hypotheses for the word image. The goal is to assign a matching score to each lexicon entry or to select the best lexicon entry among the set. Various approaches for handwritten word recognition have been reported. But they can be grouped into two major approaches. The one is *analytical (model based) approach* and the other is *holistic approach*.

2.2.1 Analytical Word Recognition

The basic idea in analytical approaches is to recognize the input word image as a series of units of a predefined model set known to the unit classifier. Unit segmentation must be a part of the recognition process. Different model-based segmentation methods place different emphasis on the processes of segmentation and classification to reach the final recognition output [63] [64]. One method of character segmentation by image-features followed by character recognition (so-called

segmentation driven) and another method of recognition driven segmentation (so-called *segmentation free* in order to emphasize the recognition step which is located before determining character boundary even though segmentation is used) are examples in contrasting approaches.

Some of the most successful results have been achieved by segmentation-driven techniques combined with matching algorithms of *Dynamic Programming* and *Hidden Markov Models (HMM)*. Commonly a word image is segmented into sub-images called *primitives* without reference to the lexicon. Ideally a primitive is a character. But since perfect character segmentation is hard to achieve, in practice, an over-segmentation scheme is taken. A *character segment* can be a primitive or a union of primitives. All possible character segments are tried for matches with characters in the lexicon strings.

In [66], a discrete hidden Markov model has been used to recognize a word. The *HMM* parameters are estimated from the lexicon and the training image segments. A modified *Viterbi algorithm* is used to find the best state sequence. One *HMM* is constructed from the whole language and the optimal paths are used to perform classification. In [65], a *continuous density variable duration hidden Markov model (CDVDHMM)* is used to build character models. The character models are left-to-right *HMMs* in which it is possible to skip any number of states. The observation process is another stochastic process represented by a set of symbols which could be cluster centers of feature vectors.

Another technique using *HMM* for cursive word recognition has been developed in [67]. A series of morphological operations is used to label each pixel in the image according to its location in strokes, holes and concavities located above, within and below the core region. A separate model is trained for each individual character using word images where the character boundaries have been identified.

Word recognition is achieved by maximizing a probability function.

A recognizer using dynamic programming finds the prototype sequence which generates the best fit to the word image and yields the corresponding letter sequence as the recognized letter string or ranks the given lexicon with corresponding matching scores. Ordering primitives from left to right yields a partial ordering on the segments. A path through the partially ordered sequence of segments yields a segmentation driven by character recognition. Dynamic programming is used to find the best fit sequence of segments to match a given string. The metric of recognition confidence of the corresponding character usually drives the matching engine to generate the best fitting sequence [58] [59] [60] [61] [62].

2.2.2 Holistic Word Recognition

Holistic approaches do not attempt to label parts of the image using sets of models [70] [71]. They extract holistic features from a word image and match the features directly against the entries of a lexicon. The process of constructing a holistic feature representation for all lexicon entries should be accomplished before matching to the extracted features. [68] [69]

Holistic methods described in the literature have used a variety of holistic features but they commonly include a variety of structural descriptions. Structural features are extracted from the image and describe the geometric and topological characteristics : *local shapes* such as edges, curvatures and end points, *perceptual features* such as loops, ascenders, descenders and directional components and *global features* such as polygonal approximations and topological contour description.

These features are represented more robustly by a graph or a string of symbol codes, each code referring to a different feature or a combination of features. A

location-coded string representation captures the locations in the image of each feature. The feature symbols and spatial relationships between features are represented by a graph representation that can show two dimensional relationships.

Dynamic programming, edit-distance metric or stochastic similarity measures are adopted as matching methods. Syntactic recognition methods attempt to derive the string representing the feature description of word classes using production rules. Some holistic classifiers have been developed for use in reduction of large lexicons and verification of recognition results obtained from other classifiers [73] [74].

2.2.3 Word Recognition Combination

To optimize system performance, multiple word recognizer combination [73] and control strategies [72] have been proposed. In [72], multiple recognizers are connected in a serial path. A decision making strategy is applied at each intermediate connection on this path. At each stage, a recognition result can be produced if the recognition results are of a high enough quality. More recognition effort is paid to poor quality input to reach a final output. This control strategy achieves a better performance than that obtained by using a single recognizer only.

2.3 Phrase Recognition

Some character recognition and word recognition studies have been extended to text and phrase recognition in applications of postal address interpretation [1] , bank check processing [3] and tax form processing [5]. Previous research efforts related to handwritten phrase/text recognition problems most often assumes that

words are already isolated or can be perfectly isolated.

The study of handwritten phrase recognition algorithms has been lacking compared to recognition of isolated units. There are a limited number of research reports that discuss word separation for unconstrained handwritings. Generally the scheme of phrase recognition systems follows the processing flow consisting of these steps: word segmentation, word recognition and postprocessing. Separated word images are sent to a word recognizer and contextual information and linguistic constraints are used in postprocessing to complete the phrase selection process, assuming that word images are perfectly isolated. As another approach, to compensate for the imperfections of word segmentation, one or more grouped segments are sent to be identified by the word recognizer and the best matching path is chosen as the recognized phrase, using dynamic matching through the entire phrase image.

2.3.1 Word Segmentation

When a text is handwritten in unconstrained environment separating the text into words is challenging, since handwritten text lacks the uniform spacing which is normally found in machine-printed text or constrained handwriting where words are written in boxes whose locations are known.

Word segmentation precedent to recognition depends primarily on geometric relationship between components. The commonly adopted computational approach to solve such issues is to extract the distance between pairs of adjacent components. However, because of different writing style and flourishes, a simple geometric distance between components is not enough to determine word break

points. Several distance measuring methods have been investigated, such as *run-length Euclidean heuristic distance* [76], *convex hull* [75] etc. These methods analyze relations between adjacent connected components and find gap metrics to cope with various spacing styles.

Intuitively the horizontal distance between the bounding boxes of two neighboring components can be used as a metric to measure word gaps, where the bounding box is defined as the smallest rectangle enclosing a component. But horizontal distance between bounding box pairs is not sufficient to deal with separation of components where the bounding boxes are overlapped.

The run-length and Euclidean heuristic distance (*RLE*) [76] uses the minimum run-length method for connected components that overlap vertically by more than a threshold; otherwise, the minimum Euclidean distance is used. Several heuristics are added to the *RLE* method and this method achieves about 90% accuracy in detection of word gaps.

The convex hull method has been designed to avoid problems associated with the bounding box method [75]. First, it is not constrained to rectangular approximations of the shape of each component. Rather it uses the convex hull, which is the smallest convex polygon enclosing the component. Second, it does not take the horizontal distance between the rightmost and leftmost extremities of the neighboring convex hulls. Rather it finds the line segment joining the centers of gravity at the convex hulls, determines the points at which this line intersects each hull and computes the Euclidean distance between these points. This metric performs better than *RLE* in correctly detecting word gaps.

As an effort of encoding writing styles, a model using consecutive stroke transition as basic raw material to detect word gaps, was proposed in a neural network

[78]. The style is captured by characterizing the variation of spacing between adjacent characters as a function of the corresponding characters themselves. Two consecutive intervals, component heights, a flag for inter-component transition, a flag for the last segment, and averages of intervals and heights are used as input parameters to a neural network with one hidden layer. Experiments showed good performance in isolating words from text.

2.3.2 Phrase Recognition Models

A phrase recognition model generally involves (i) word separation from the given phrase image, (ii) phrase hypotheses generation combining image segments and possible *lexemes*. (iii) interpretation with constraints or language models. The *lexemes* are defined by the exclusive atomic word list which constitute a lexicon. Several recognition strategies have been introduced to solve various application problems.

Systems finding the best match between word image segments and lexicon entries have been developed to overcome imperfections of word segmentation. One or more grouped image segments are sent to a word recognizer to be matched against entries in a lexicon pool. The best matching path through the entire image is chosen as the recognized phrase using dynamic matching in applications for bank check amount recognition [3] [4].

In order to overcome incompleteness of lexicon entries as well as word segmentation failure, a recognition-driven partial matching strategy has been developed for street name recognition of postal addresses. A method of recognition-driven

segmentation into constituent words by encoding the author's writing style is presented in [78]. A lexicon is generated so that it contains only the keywords. Matching between the lexicon entries and the input image is performed to interpret the street name.

A lexicon expansion method which enlarges vocabularies with abbreviations, and numeric word replacement is also used to overcome incompleteness of a lexicon [81] [82]. For efficient processing, a method to minimize computational complexity has been sought by maximizing the advantage of the lexicon driven word recognizer. Segmentation statistics which determine the transition probability between segments of a character is used for the early rejection of lexicon entries precedent to actual recognition.

An interpretation method for phrase lines which have several different fields of city-name, state-name and ZIP code in address block has been developed [79]. The interpretation mechanism involves the generation of multiple hypotheses in the parsing sub-system and the usage of very large lexicon recognition such as city-names (38405) and ZIP codes (43030). Word separation within multi-line text, separating the digit or characters within a word and efficient word recognizer driving techniques have been implemented.

2.4 Cognitive Models

2.4.1 Word Perception Model

In cognitive research areas, several word perception models have been presented. One of them, a *time course model of word recognition* has been developed in [85].

This model depends upon three basic assumptions; (i) Perceptual processing

takes place within a system in which there are several levels of processing. For visual word perception, a visual feature level, a letter level, and a word level are assumed as processing levels. (ii) Visual perception involves two types of parallel processing. Information covering a region in space at least large enough to contain a four-letter word is processed simultaneously. Also visual processing occurs at several levels at the same time. (iii) Perception is fundamentally an interactive process. *Top-down* or *conceptually driven* processing works simultaneously and in conjunction with *bottom-up* or *data driven* processing to provide multiplicity of constraints that jointly determine what we perceive.

The model has been implemented with a discrete relaxation network composed of nodes and connections between them. Nodes are organized into a hierarchical network structure of three levels mentioned above. The feature level includes a node that detects every feature at each of four letter positions. Only four-letter words are considered because of the second assumption of spatial parallel process. The letter level includes a node for every letter at each letter position. The word level includes a node for each allowable four-letter word.

Nodes sharing common information are linked by excitatory pathways; nodes that do not share information are linked by inhibitory pathways. Recognizing a particular type of information increments the activation level of its node and activation spreads by means of the excitatory and inhibitory links through the network. The entire network operates in discrete and recursive fashion. The input for the network is derived from the feature detectors that output activation values proportional to the probability that their features exist in the image. These activation values are weighted by constants associated with each arc. The results are added to the excitatory connection and subtracted from the inhibitory connections

of each feature detector. Letter nodes that are so-activated send weighted excitatory and inhibitory activation values to word nodes and the word nodes then send activation values to the letter level. A word or letter is recognized when its activation level exceeds a threshold amount.

Though this model does not fully represent the human word perception process because it lacks several types of relevant information such as phonology or orthography, this model still provides a computational approach of the cognitive interactive process based on image contents within lexicons.

2.4.2 Epistemic Utility

A theory of cognitive decision making, *epistemic utility theory* provides a means of quantifying the value of avoiding error and the value of acquiring new information [89] [88]. An epistemic system is described in terms of *knowledge corpus*, *information valuation*, and *truth valuation* and designed to reflect the preferences that a system ought to have for the various possible correct and incorrect answers available to it when a system is called upon to choose among competing propositions.

Epistemic utility theory employs two independent utilities within the knowledge corpus. One is designed to characterize the truth support of the propositions being evaluated as subjective truth valuation, and the other is designed to characterize the informational value of rejecting them.

The fundamental contents of an epistemic utility-based decision making approach is that, by endowing the two utilities with the mathematical structure of probabilities, the attributes of propositions are quantified and compared. Those propositions whose truth-support does not outweigh their informational-value-of-rejection should be rejected. All other propositions should be retained as serious

possibilities. This procedure retains all propositions that are considered *good* without focusing exclusively on a search for one that is deemed the *best*. The first principle, *epistemic utility*, regulates evaluation of information regardless of its truth-value, and the second principle, *credal probability*, regulates the credibility of information regardless of its worth. These two principles lead to a dynamic decision making.

A theory of *satisficing decision* which accepts one that is not optimal but meets minimum requirements has been introduced as a method for obtaining sub-optimal solutions [93] [94] [95]. The epistemic utility theory has been adapted to the satisficing decision in [90] [91] [92], as a computational model which evaluates evidence toward a goal. The notion of truth and information is re-interpreted in a practical setting. The notions of credal probability and epistemic utility are extended to *accuracy support*, conformity to a given standard and *liability exposure*, susceptibility or exposure to something undesirable.

Chapter 3

Hierarchical Character Recognition

Feature extraction and classification methods are usually involved in a character recognizer. Most character recognition techniques described in the literature use a *one model fits all* approach, *i.e.*, a set of features and a classification method are developed and every test pattern is subjected to the same process, irrespective of the constraints present in the problem domain. We argue that such methods are passive. Most recognition methods that adhere to this paradigm have pre-determined (hence static) parameters such as: (i) dimension of feature vector, (ii) the distance or confidence metric considered as threshold for recognition acceptance, (iii) the set of classes contending to match with a test pattern and (iv) equal computational resources available to every test pattern.

In the optimality paradigm, if the correspondence between the pre-selected training set and the actual input stream in an operating environment is reasonably large (such as machine-printed text with the same font set), the minimum number of selected feature sets which satisfy desired classification performance is one of the optimal solutions. However, empirically it is realized that, in handwriting recognition, the similarity between a training set and an actual input stream is

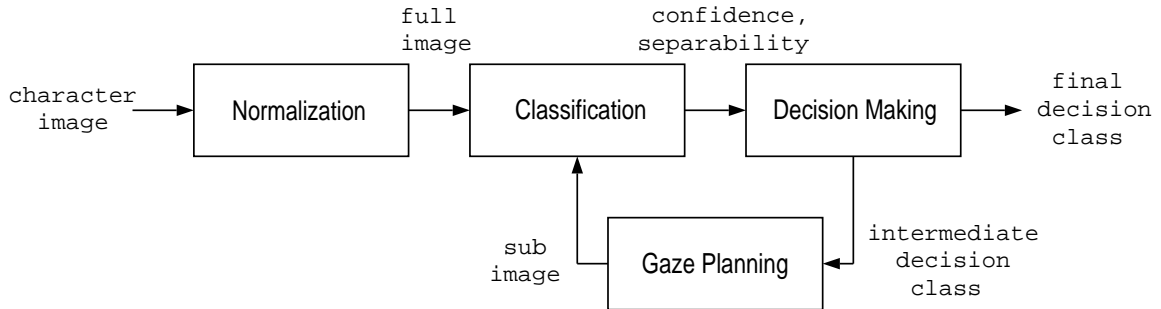


Figure 3.1: Block diagram of hierarchical classifier

relatively weak and even in a training set, the homogeneity in shape within a class is not so strong.

Using a high dimensional feature space is one of the solutions to deal with the poor homogeneity in order to add discriminatory power to a classifier. Intuitively speaking, a high dimensional feature set which is selected in order to maximize recognition performance usually generates excessive separation in a ‘good enough’ input. In a cost-optimal approach, selection of features usually generates a trade-off problem between the computational and storage cost and the recognition performance.

Recently several successful approaches have been introduced by selecting features from different scales [17]. These approaches are combined with classification algorithms and improve the recognition performance dramatically [17] [36]. However, approaches that mix features of different scales in high dimension have several disadvantages as well. Correlations among features extracted in different scales but in similar types are often unavoidable. A distance metric that captures features at all the scales and at the same time represented by a single scalar value

is hard to find because the heterogeneous composition of the feature set makes it difficult to find an appropriate common distance metric. The computational complexity of a multi-resolution recognizer such as the GSC described in [36] is known to be high because of the burden of additional features needed to accommodate the multi-resolution aspect of the feature space.

In this chapter we present a dynamic and multi-resolution character recognizer (Figure 3.1) which uses a hierarchical feature space that preserves the benefits of the multi-resolution models while circumventing the disadvantages cited above. The recognizer begins with features extracted in a coarse resolution and focuses on smaller sub-images of the character on each recursive pass, keeping the same feature space (dimension). In other words, successive passes render the same features at a finer resolution. Thus the recognizer is effectively working with a finer resolution of a sub-image in each time till the classification meets acceptance criteria.

The sub-images are selected by an approach called *gaze planning*. It guides the successive passes through a series of sub-images where it believes the features of interest are present. This is done in an adaptive and dynamic manner. By beginning with a global view, and progressively increasing the resolution by sub-images of the test pattern we are able to accomplish an adaptive classification with optimal resource usage in terms of storage space and computational cycles of recursion. Depending on the classes competing for the top spot different sub-images might be selected at any given stage. For example, if '3' and '5' are the top contenders at a particular recursive stage, the features from the upper-left zone of the test pattern holds greater discriminatory power however, the lower-left zone does not as shown in Figure 3.2. If '7' and '5' are the top contenders at a particular recursive stage, the features from the lower-left zone has larger discriminatory power

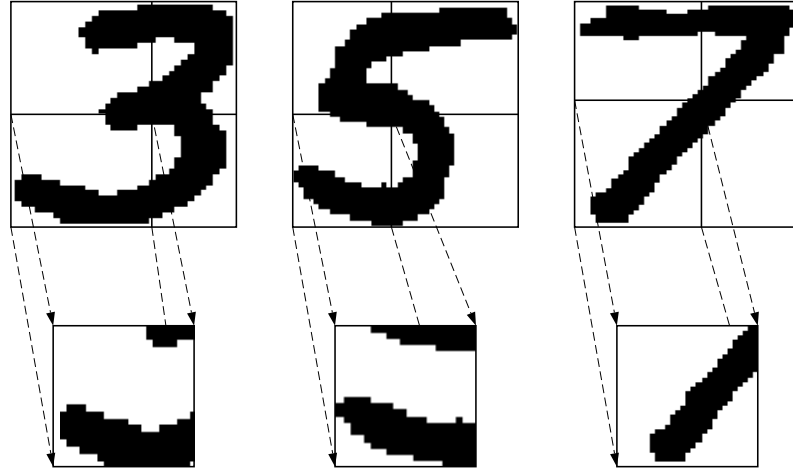


Figure 3.2: An example of using sub-images

compared to the case of '3' and '5'.

3.1 Hierarchical Classification

Pattern recognition techniques typically extract and evaluate features of a pattern image. Each feature is represented by its values in a feature vector of one or more dimensions. In statistical pattern recognition techniques, a probability value along each feature dimension is computed based on the strength of the features extracted.

Let \mathbf{x} be a feature vector with dimension N_x and C be a set of mutually exclusive classes of size N_c ,

$$\mathbf{x} = [x_0, x_1, \dots, x_{N_x-1}], \quad (3.1)$$

$$C = \{c_0, c_1, \dots, c_{N_c-1}\} \quad (3.2)$$

A *posteriori* probability of class c_i with feature vector \mathbf{x} of an input image is given by Bayes' Rule as,

$$P(c_i|\mathbf{x}) = \frac{P(\mathbf{x}|c_i)P(c_i)}{P(\mathbf{x})} \quad (3.3)$$

where the probability of a feature vector \mathbf{x} is

$$P(\mathbf{x}) = \sum_j P(\mathbf{x}|c_j)P(c_j). \quad (3.4)$$

If all the features are statistically independent, the class probability of an input pattern can be readily computed. If features are somewhat correlated to each other, it is hard to find the joint probability from individual probability of each feature value. This problem is serious in symbolic feature space compared to coefficient based feature space such as Fourier descriptor or Zernike moments, since the coefficient features are extracted from the mathematical orthogonal decomposition. However well selected and properly quantized combination of features which are combined with nearest neighbor classification such as GSC can achieve high accuracy in classification, even if it does not provide any direct solutions to the correlation problem.

3.1.1 Probability in vector quantized hyper feature space

Since an actual joint probability of feature vectors is hard to obtain if the features are correlated, a direct sampling of the probability in the multi-dimensional feature space can help to overcome the correlation problem in symbolic features (Figure 3.3). By dividing multi-dimensional feature space into non-overlapping subsections, a quantized hyper feature space is established where the size of the hyper space is the product of all the feature dimension and the quantization resolution.

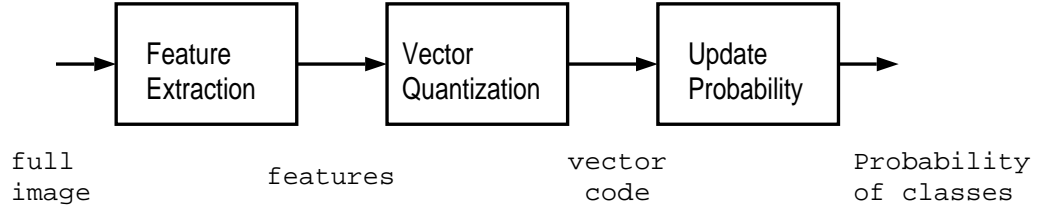


Figure 3.3: Processing steps of class probability sampling

As one of the simplest way, given a threshold vector

$$\mathbf{x}_T = [x_{T0}, x_{T1}, \dots, x_{T(N_x-1)}] \quad (3.5)$$

for the extracted feature vector \mathbf{x} , a binary feature vector can be obtained.

$$\mathbf{b} = [b_0, b_1, \dots, b_{N_x-1}] \quad (3.6)$$

\mathbf{b} is a $N_x \times 1$ matrix and 2^{N_x} such matrices are possible. Each component is binarized as

$$b_i = \begin{cases} 0 & \text{if } x_i \leq x_{Ti}, \\ 1 & \text{otherwise} \end{cases} \quad (3.7)$$

where $i = 0, 1, \dots, N_x - 1$.

By defining a vector code space,

$$V = \{0, 1, 2, \dots, 2^{N_x} - 1\} \quad (3.8)$$

every binary feature vector can be mapped onto the vector code space. In other words, every N_x dimensional hyper feature vector can be then mapped to a unique code as shown Figure 3.4.

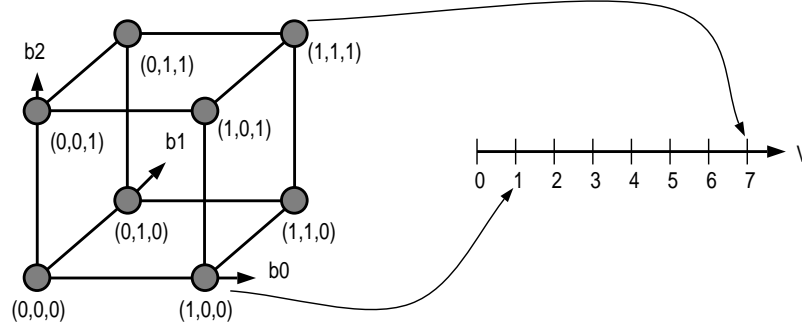


Figure 3.4: Mapping from binary feature vectors to vector codes when $N_x = 3$

A *posteriori* probability $P(c_i|v)$ of class c_i with a vector code $v \in V$ is then given by

$$P(c_i|v) = \frac{P(v|c_i)P(c_i)}{P(v)} \quad (3.9)$$

A *priori* probability $P(c_i)$, measurement-related information $P(v|c_i)$ and vector code probability $P(v)$ are obtained during training on a pre-classified image set and the $P(c_i|v)$ is stored in a look-up table of size $|V| \times N_c$ which makes for a very efficient program. We will denote the *a posteriori* probability $P(c_i|v)$ by the notation $P_i[v]$. The storage needed would be multiplication of 2^{N_x} , N_c and number of bytes per entry. For example, if $N_x = 9$ features and there are 4 bytes assigned per entry, the table size for holding $P_i[v]$ would be $2^{N_x=9} \times 4 \text{ bytes} = 2048 \text{ bytes per class}$.

We obtained 54% accuracy of top choice on a set of numerals using 9 features ($N_x = 9$), which is not acceptable in practical applications. The number of features, and hence the feature space dimensionality must increase for accurate classification of handwritten characters, especially in practical applications where the data tends to be noisy. A common approach to addressing this issue, as suggested by

pattern recognition literature, would be to add new features. For example, the multi-resolution method of GSC described in [36] uses a 512 dimensional binary feature vector, the table for which would be impossible to create or store since $|V| = 2^{512} \times N_c$. Bailey and Srinath use 30 real coefficients to obtain satisfactory recognition performance [24]. However, binarization process loses the classification accuracy obtained by using the real coefficients.

3.1.2 Simulation of multi-resolution

An alternative approach to the tradeoff dilemma (between accuracy and storage) described above is to keep the number of features small to begin with and focus on ‘important’ sub-images of the original image. Sub-images are obtained at any stage by using a division rule such as quad tree or quin tree (Figure 3.5) [104]. Features extracted in earlier passes correspond to coarse global features while those in later passes provide local features. The multi-resolution of features is thus achieved, although the actual feature extraction process and the type of the features remains the same.

This adds to some additional storage space. For example, in a quin tree of maximum depth 4, there are 156 sub-images and the additional storage size is $2^{N_x=9} \times 4 \times 156 = 312Kbytes$ per class. Still it is very reasonable compared to $2^{512} \times 4$ per class warranted by the GSC recognizer.

Let $P_i[v(d, l)]$ be a *a posteriori* probability of vector code $v(d, l)$

$$P_i[v(d, l)] = \frac{P(v(d, l)|c_i)P(c_i)}{P(v(d, l))} \quad (3.10)$$

where

$$P(v(d, l)) = \sum_j P(v(d, l)|c_j)P(c_j) \quad (3.11)$$

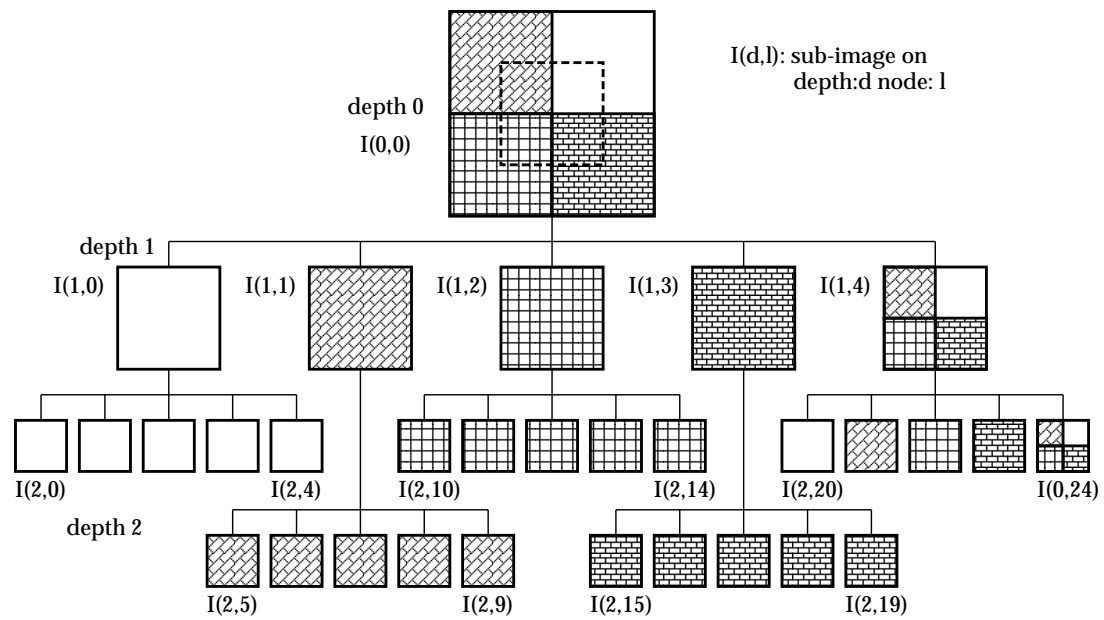


Figure 3.5: Sub-images generated by a quin tree division rule

and $v(d, l)$ is a vector code which is extracted from the sub-image $I(d, l)$ of l^{th} node at tree depth d for class c_i . The tree depth and node are limited as,

$$d \in \{0, 1, \dots, N_D - 1\}, \quad l \in \{0, 1, \dots, N_B^d - 1\}. \quad (3.12)$$

N_D is the maximum depth and N_B is the number of branches per node.

Let us assume that features of sub-images in an input pattern image are strongly correlated and the degree of correlation is uniquely characterized by the class. We define a confidence function $Q(c_i|d)$, in a short form Q_i^d , of class c_i of an entire layer at depth d by the product of *a posteriori* probabilities of vector codes for the same class within all sub-images as,

$$Q_i^d = Q(c_i|d) = \prod_l P_i[v(d, l)] \quad (3.13)$$

Based on similar assumptions, the confidence $Q(c_i)$, in a short form Q_i of class c_i considering all the layers of a tree is defined as the product of all the sub-nodes.

$$Q_i = Q(c_i) = \prod_d Q_i^d = \prod_d \prod_l P_i[v(d, l)] \quad (3.14)$$

$\sum(N_B^d - 1)$ multiplications are necessary to obtain the answer from the confidence function of a class, using equation (3.14). Applying the logarithm function to (3.14), we can convert it into the equivalent addition operations.

3.1.3 Efficient processing of the recursive tree

Based on the assumptions described in the previous section, a confidence of a class can be updated recursively using a sub(partial)-tree selected from a whole tree starting from the root extending to the sub-nodes successively. An intermediate confidence of a classes which is the confidence in a sub-tree at a certain time (recursion) can be obtained by the products of the probabilities of the sub-nodes processed up to the time. If the sub-tree is optimally and adaptively selected for the

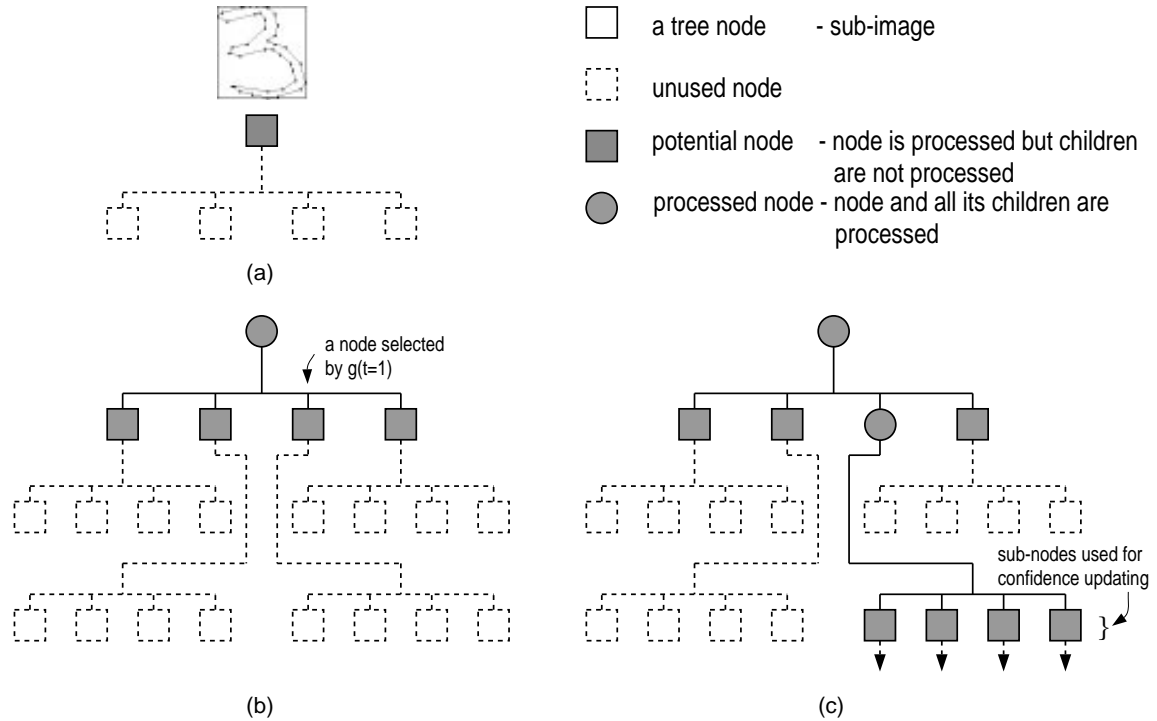


Figure 3.6: An example of confidence update recursion; (a) starting from root, (b) after 1 recursion pass, (c) after 2 recursion passes

input image pattern, an adaptive multi-resolution recognition can be simulated. Also if the recursive recognition is performed until satisfying the acceptance criteria, a cost optimal recognition which adapts resource usage to the input quality can be achieved.

Gaze planning is a means of expanding only some of the nodes in a tree in order to select the optimal sub-tree. We fix the number of nodes whose probabilities are used to update confidences on each pass and we limit the selection of new nodes to a subset of *potential nodes*. A potential node is defined as a node whose

probabilities are used in the previous pass but whose sub-nodes are not used. The shaded squares in Figure 3.6 show the potential nodes.

The intermediate confidences can be computed recursively at time t ,

$$Q_i(t) = Q_i(t-1) \cdot Q_i^{sub}(t-1) \quad (3.15)$$

where $Q_i^{sub}(t)$ is the product of a *posteriori* probabilities of feature vector codes for class c_i obtained from the images which are the sub-nodes of a node that is selected at time t (Figure 3.6). Since the number of sub-nodes is finite, the time index t is only effective in the range of

$$t \in \{0, 1, \dots, \sum_{d=1}^{N_D-1} N_B^{d-1} - 1\} \quad (3.16)$$

Let a function $g(t)$ determine whether a node is to be examined further in the next iteration. We can write $Q_i^{sub}(t)$

$$Q_i^{sub}(t) = \prod_{k \in sub(d,l) | (d,l)=g(t)} P_i[v(d_k, l_k)] \quad (3.17)$$

where $sub(d, l)$ is a function which gives the indices (d_k, l_k) of sub-nodes on $(d+1)^{st}$ layer which belong to the node (d, l) . If we choose $g(t)$ as a gaze function that finds the most informative node within the pool of potential nodes for updating class confidence, a confidence value of the true class will rapidly approach its steady state.

Usually a classification result of a recognizer is accepted based on the confidence of top choice and its separability from the second choice [36] [18]. Using the same criterion, a normalized confidence $\kappa_i(t)$ of a class c_i at recursion time t is defined by the ratio of the current (intermediate) confidence and the maximum possible confidence of the class regardless of the specific feature vector,

$$\kappa_i(t) = \frac{Q_i(t)}{R_i(t)} \quad (3.18)$$

$$R_i(t) = \prod_{k=0}^{t-1} R_i^{sub}(k) \quad (3.19)$$

where $R_i(t)$ is the product of possible maximum probability of any vector code in class c_i through the selected node sequence. The node sequence is accumulated outputs of a gazing function g till pass $t - 1$.

In a recursive representation,

$$R_i(t) = R_i(t - 1) \cdot R_i^{sub}(t - 1), \quad (3.20)$$

where

$$R_i^{sub}(t) = \prod_{k \in sub(d,l) | (d,l)=g(t)} \max(P_i[v(d_k, l_k)] | \forall v(d_k, l_k)) \quad (3.21)$$

The recursive equation of the normalized confidence for class c_i is given as

$$\kappa_i(t) = \kappa_i(t - 1) \cdot \kappa_i^{sub}(t - 1) \quad (3.22)$$

where $\kappa_i^{sub}(t)$ is the normalized confidence of newly added branches of node (d, l) ,

$$\kappa_i^{sub}(t) = \kappa_i^{g(t-1)} = \frac{Q_i^{sub}(t)}{R_i^{sub}(t)} \quad (3.23)$$

From (3.17) and (3.22), the intermediate confidence and normalized confidence of each class is obtained at every pass. After ordering of classes using normalized confidences. separability, $\sigma(t)$, of the top two choices is given by

$$\sigma(t) = \frac{Q_{top}(t)}{Q_{second}(t)} \quad (3.24)$$

when a class c_{top} has the highest confidence at time t .

If the result does not meet acceptance criterion built by normalized confidence and separability, features within the same space are computed at a different scale in the next pass.

For driving the gazing function, local separability which presents the support of current top ranked class at time t is extracted in each potential node. If a class c_{top} is the top choice and a class c_{second} is the second choice in the intermediate confidence result of time t , the local separability $\sigma^{(d,l)}(t)$ of a potential node (d, l) is given by

$$\sigma^{(d,l)}(t) = \frac{P_{top}[v(d, l), t]}{P_{second}[v(d, l), t]} \quad (3.25)$$

By comparing local separability of each potential node, the next gazing node which is expected to be the most supportive of the current recognition results can be determine. If the top and second rated classes change by updating, the local separability should be renewed at all the potential nodes in which the updating has ended. Figure 3.7 summarizes the recognition algorithm. To reduce computation complexity, if we use a logarithm function, the equations of products for the parameters are transformed to summations.

3.2 Implementation

In this section, a detailed implementation procedure of the proposed model is described. An efficient contour representation, normalization, feature extraction and vector code generation and methods for dividing sub-images are explained.

3.2.1 Image Representation

A character image is scanned and binarized. The binarized image is represented by *piecewise linear contours* for basic image representation in our implementation. In order to generate piecewise linear contours, the pixel based contours are first

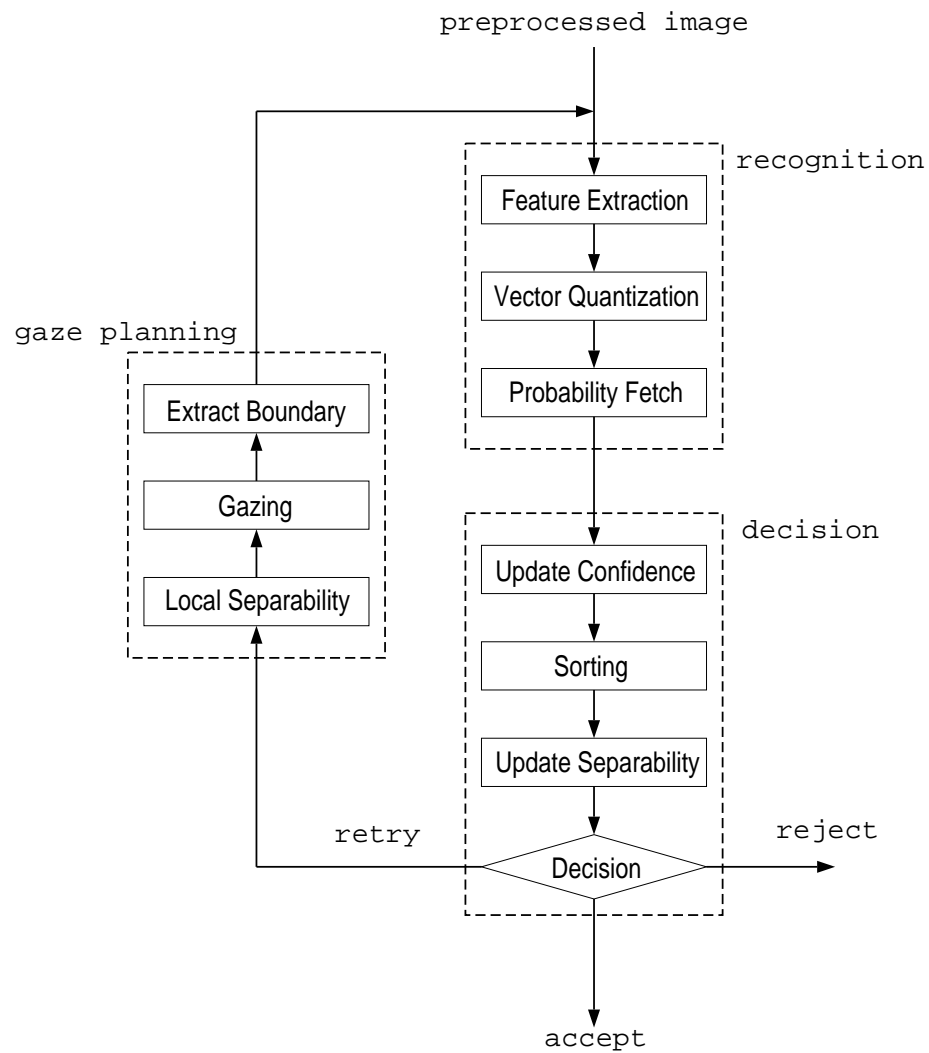


Figure 3.7: Methodology of simulating a multi-resolution feature space by recursive processing of *selected* sub-images

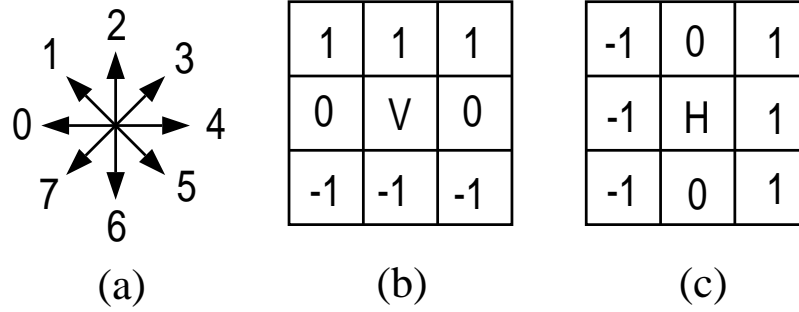


Figure 3.8: Contour representation; (a) contour slope division, (b) vertical weight map and (c) horizontal weight map

obtained from an input image using the chain code method described in [99] [100] which keeps the coordinate and 8 directional inletting slope information from the previous contour pixel as shown in Figure 3.8-(a). An equally weighted window is applied to a contour component and its neighbors to obtain the curvature at that point. The window size N_w is determined by the height of the bounding box and the image density. The bounding box is defined as the smallest rectangle enclosing a character and the image density is obtained from the ratio of the number of black(inked) pixels to the size of the bounding box.

$$N_w = 2 \cdot k_w \cdot height \cdot density + 1 \quad (3.26)$$

k_w is a control parameter.

At each contour point, horizontal and vertical slopes are found using the masks of Figure 3.8-(b,c) and binarized using a threshold (0.65 in our implementation). Contour points which have high curvature (derivative of the binarized slope) are identified as critical points. Contour segments between adjacent critical points are approximated to be linear. This representation adds to the overall processing speed

by reducing the total data points and also by eliminating the need for smoothing and noise removal. Figures 3.9-(b) and (c) illustrate the procedure.

3.2.2 Normalization

The major sources of variation in handwriting are size(height), slant, skew(slope), stroke width and rotation [58] [64]. Minimization of these variations will help to boost the recognition performance. Since the input is a pre-isolated character image, we assume that shape variations caused by rotation and skew are very small. In our implementation, since the feature extraction is primarily based only on the contour component (gradient and moment) and since the extraction grids are variable to the bounding box, the variation generated by stroke width and size do not generate a significant effect.

For slant normalization, continuous vertical line components (spanning several critical points) whose length is greater than a threshold are identified. Slant angle is estimated by the average of the slant angles (from the horizontal) of all vertical line components, where the angles are weighted in proportion to their length in pixels. Using the estimated slant angle, the critical points are mapped for slant correction. Figure 3.9 - (d) shows the mapped points after slant correction.

3.2.3 Sub-images

A variable size rectangular grid is used to define sub-images for the quad and quin trees. The center of mass of the contour is first computed. The bounding box of a character image is divided into four rectangular regions using the center point as shown in Figure 3.10-(c). A vertical and horizontal line through the center of mass determines the four regions for the quad tree. The quin tree structure

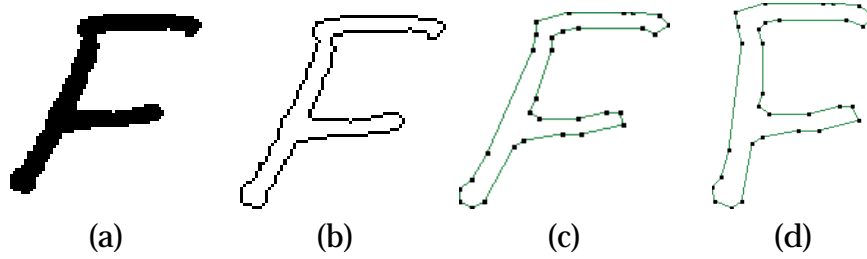


Figure 3.9: Contour representation and preprocessing; (a) input image, (b) pixel based contour representation, (c) piece-wise linearized contour and (d) piece-wise linear contour representation after normalization

is similar with an additional fifth sub-region which is located in the central area formed by joining the centers of the other four sub-regions. Subsequent layers are successively constructed using the same method.

Quad and quin trees with a maximum depth of 4 were implemented. Figure 3.10-(c)s, (d)s and (e)s show sub-regions for each layer. Figure 3.10-(f)s show the reconstructed contour images using center of the finest sub-regions at the deepest (4^{th}) level.

3.2.4 Features

Two types of contour histogram in gradient and moment based projections are used as feature extraction methods. The size of the feature vector is kept small for reasons described earlier. An eight bit binary feature vector for each sub-region in the quad tree representation and a nine bit feature vector for each sub-region in the quin tree representation are generated after binarization of feature measurement.

Contour slope angles are quantized in the pairs of $\pi/4$ slots that are opposite

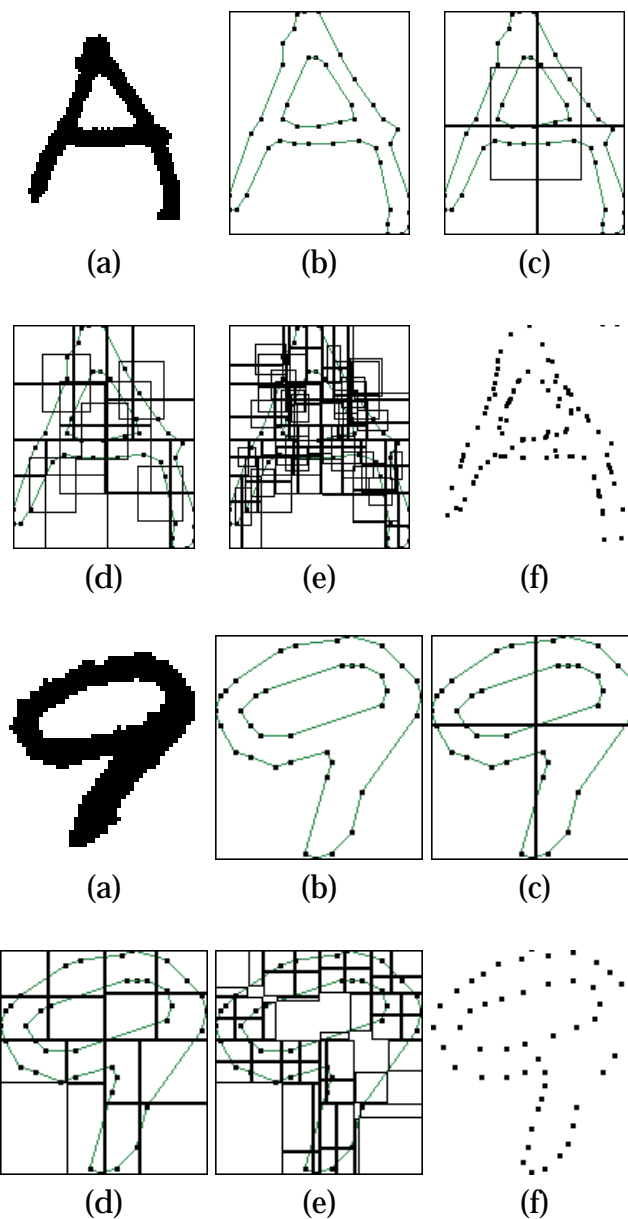


Figure 3.10: Character image and sub images of quin tree (class A) and quad tree (class 9) : (a) original image, (b) piecewise linearized contour, sub-regions at (c) layer 1, (d) layer 2, (e) layer 3, and (f) pseudo reconstructed contour image from center points of final layer

each other labeled 0 to 3 in Figure 3.11 - (a). The angle between the horizontal line and the linearized contour piece between two adjacent critical points is assumed to be the slope angle of the contour segment and the length of a contour piece is taken as the Euclidean distance between adjacent critical points as shown in Figure 3.11-(b).

Gradient feature measurement x_i^{grad} is determined by the ratio of the total contour length r_i^{grad} corresponding to the contour segments whose slope angle falls in division i and the total contour perimeter r_{total} in that sub-image.

$$x_i^{grad} = \frac{r_i^{grad}}{r_{total}} \quad (3.27)$$

Each gradient binary feature is generated by thresholding using its corresponding x_T^{grad} .

$$b_i^{grad} = \begin{cases} 0 & \text{if } x_i^{grad} \leq x_{Ti}^{grad}, \\ 1 & \text{otherwise} \end{cases} \quad (3.28)$$

Moment features are used to obtain the other four (or five in the case of quin tree) bits. Moment measurements are obtained from the distribution of the contour within the sub-images. A moment measurement x_i^{moment} is defined by the ratio of the length r_i^{moment} of the sum of contour segments which are present within a sub-image i to the total contour perimeter.

$$x_i^{moment} = \frac{r_i^{moment}}{r_{total}} \quad (3.29)$$

Each moment binary feature is determined by thresholding using its corresponding x_T^{moment} .

$$b_i^{moment} = \begin{cases} 0 & \text{if } x_i^{moment} \leq x_{Ti}^{moment}, \\ 1 & \text{otherwise} \end{cases} \quad (3.30)$$

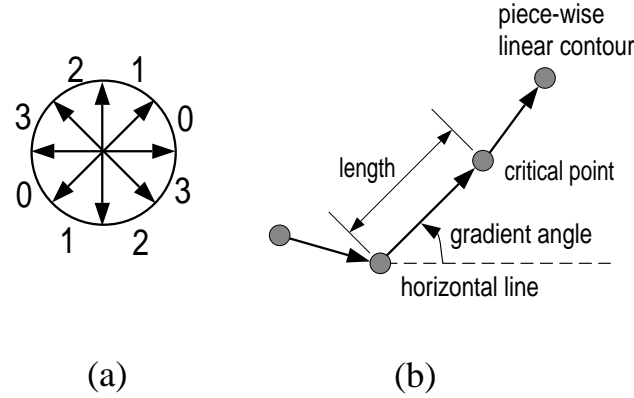


Figure 3.11: Gradient feature generation; (a) gradient division and (b) gradient angle measure

When the binary feature vectors are generated, the thresholds values are uniformly set at 0.25 for both gradient and moment feature measurements except for the moment feature measurement of the fifth sub-images of the quin tree. Since the fifth sub-image of the quin tree is a region overlapping the other four sub-images and the character contours are relatively concentrated around the center of mass, for this measurement, the threshold is set at 0.5.

3.2.5 Training

A eight (or nine for quin tree) bit size feature vector code is generated on each sub-region by the combination of the binary feature vector. All probabilities, $P(c_i)$ s and $P(v(d, l)|c_i)$ s are independently obtained from a training image set. Using each $P(v(d, l)|c_i)$ and $P(c_i)$, all the $P_i[v(d, l)]$ in every sub-image are calculated and stored in a look-up table.

As mentioned in section 3.1.2, a logarithm function is applied to the confidence function in order to reduce computation time. It can replace all the multiplications of equation 3.14 into additions which are much faster in conventional computers. Following this replacement, all $P_i[v(d, l)]$ values are stored in logarithm scale.

Since all values of probabilities of every sub-image are stored in a table, large storage size is required even though the feature vector dimension is kept small. In our implementation, 256 step resolution (1 byte size) for logarithm scale is used in order to reduce the storage size. If we use a maximum depth of 4 for the quad tree (as used in the experiments), the total storage size of this implementation for digit recognition becomes $10 \text{ classes} \times (1 + 4 + 16 + 64) \text{ sub-images} \times 2^{N_x=8} \text{ feature vector codes} \times 1 \text{ byte resolution} = 212.5 \text{ Kbyte}$

3.3 Experiment

Experiments have been performed using isolated digit and character images present in the NIST database 3 (SD 3) [102], and CEDAR database collected from handwritten postal addresses [101]. The images were divided into mutually exclusive training and testing sets for various experiments shown in Table 3.2. A 782 MIPS *SUN ULTRA 2* is the experimental platform.

In order to test validity of the proposed classification method, three types of experiments have been performed; (i) recognition performance for various digit/character recognition applications using correct recognition rate without rejection, (ii) recognition performance comparison using two different sub-image

Experiment type	No. of classes	Training images	Testing images	Image source(s)
EXP 1	10	25656	12242	Postal+NIST
EXP 2	10	159228	53301	NIST digit
EXP 3	26	34961	11481	NIST alpha upper-case
EXP 4	26	29059	9623	NIST alpha lower-case
EXP 5	52	64020	21104	NIST alpha full-class
EXP 6	52	64020	21104	NIST alpha folding

Table 3.1: Experiment types and size of training and testing sets

generation methods (quin and quad tree), and (iii) comparison of different classification methods using the same feature vector extracted from hierarchical feature space. A neural network and a k-nearest neighbor classifier have been implemented as other classification methods for comparison.

Being different from other classification methods, the proposed algorithm has a recursive recognition process which can be terminated after satisfying a terminal condition for recognition in the decision stage. This recursive loop can generate dynamic characteristics by changing the termination condition, since the usage of resources and recursion depth depend on the input image quality. An experiment has been performed to see dynamic behavior between recognition rate and resource usage using various degrees of separability between top two choices as a terminal condition.

Experiment type	In training set		In testing set	
	top 1 (%)	top 2 (%)	top 1 (%)	top 2 (%)
EXP 1	99.5	99.9	96.2	98.8
EXP 2	98.6	99.7	98.2	99.5
EXP 3	99.2	99.9	96.4	98.8
EXP 4	98.7	99.8	90.4	97.1
EXP 5	93.4	98.7	78.2	94.7
EXP 6	97.5	99.5	92.4	97.1

Table 3.2: Recognition results using a quin tree in maximum depth of 4

3.3.1 Recognition Performance

Table 3.1 shows the types of experiments and settings. For EXP 1, the image sets for training and testing were obtained from selected sets of CEDAR and NIST data. Selected sets were exclusively divided into two parts, training and testing groups, respectively. For the other experiments, NIST sets were used. The first three sub sets were used as training sets and the last one set was used as the testing set. We assumed that the class distributions between the training and testing sets were close enough. EXP 6 used the same sets as EXP 5 and 52 full classes were used in the classification process the same as in EXP 5 but the upper and lower cases were folded when the answers were checked to avoid the errors generated by upper and lower cases having the same shape as in classes such as (C,c) (S,s) etc.

Table 3.2 shows the recognition performance using the proposed method with the quin tree in maximum depth of 4. For an input image, $1+5+25+125$ sub-images were generated and a maximum of 1404 binary feature vectors were used in units of 9. The same probability table which had been generated from the training set

Experiment type	EXP 1		EXP 3		EXP 6	
Tree type	Quad	Quin	Quad	Quin	Quad	Quin
Top 1 correct rate (%)	95.8	96.2	95.7	96.4	92.0	92.4
Top 2 correct rate (%)	98.7	98.8	98.6	98.8	97.0	97.1
Recognition time (<i>msec/char</i>)	1.60	2.80	1.94	3.24	2.49	4.06
Table size (<i>Kbyte</i>)	213	780	553	2028	1105	4056

Table 3.3: Performance comparing by changing sub-image generation rule using quad and quin trees.

was used for the recognition of both training and testing sets. Column *Top 1* indicates the success rate of the top ranked candidate being the truth and the column *top 2* indicates the rate when the true answer can be found within the top two candidates. If the top two candidates were the same class after folding in EXP 6, the third candidate was assumed to be second candidate.

Table 3.3 shows a comparison using quad and quin tree sub-image division methods. For both methods, a maximum depth of 4 was used. In quad tree, 1+4+16+64 sub-images were generated and a maximum of 680 dimensional binary feature vector was used in units of 8. The recognition time represents the average processing time per character image for the recognition including preprocessing. Table size shows the required storage size only for the probability table. Methods using quin tree achieve superior recognition performance but require additional processing time and storage space.

Table 3.4 shows a comparison of different classification algorithms on setting of EXP 3. To keep the comparison meaningful, the same features have been used for all the 3 classification methods: proposed method, a neural network method, and

Classification Method	Proposed Method	Neural Net		k -Nearest Neighbor		
		168	680	1-NN	3-NN	6-NN
Top 1 correct rate (%)	95.7	95.0	96.4	95.1	95.7	95.7
Template size (Kbyte)	612	74	976	3,777		
Recognition Time (msec/char)	1.45	1.83	11.5	384		
Training Time (hour)	1	2	24	1		

Table 3.4: Performance in various classification methods using setting of EXP 3; The proposed method operated on a dynamic mode of $\log(\kappa) = 40$ and feature vectors for the neural network of 168 sized input layer were generated using a quad tree of depth 3 (1+4+16)

a nearest neighbor method, even though this feature size is not the optimal size for the other classification methods. A quad tree of depth 4 was used for feature generation and 680 dimensional binary features were generated for a given image. Proposed method operated on a dynamic mode with separability $\log(\kappa) = 40$.

The neural network has fully connected three layers, an input layer of 680 units, a hidden layer of 353 units and an output layer of 26 units. Also a neural network with a 168 input feature vector has been implemented using features generated by a quad tree of depth 3 (1+4+16) to check the performance difference of multi-resolution in static classification. The neural network of 168 size feature vector has fully connected three layers of 168-97-26 units.

The k-nearest neighbor method described in [9] has been implemented using Hamming distance metric for binary feature codes. The performance numbers using 1, 3 and 6 nearest neighbors are shown in the table.

Storage requirement for the proposed method refers to the size of the lookup table. In the case of the neural network, it refers to the storage size of total number

of weights, and in the case of k - nearest neighbor classifier, it refers to the size of prototypes in the training set. Training time indicates wall-clock time required for generating probability table, templates or training of the networks.

The relative merits of each of the three methods are highlighted in the table. The tables shows that proposed method has some advantages in size and time while maintaining the same or showing better performance.

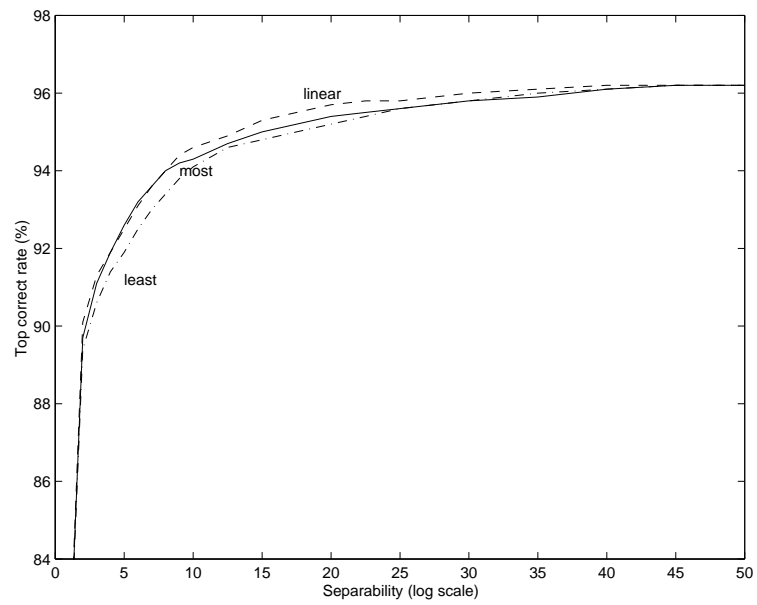
3.3.2 Dynamic Characteristics

A simple terminal criterion based on the separability was applied in order to check dynamic behavior. If the separability value of an intermediate classification result reaches the given threshold, a recognition output is generated at the point, otherwise the confidences are updated continuously until the recursion exhausts the hierarchy.

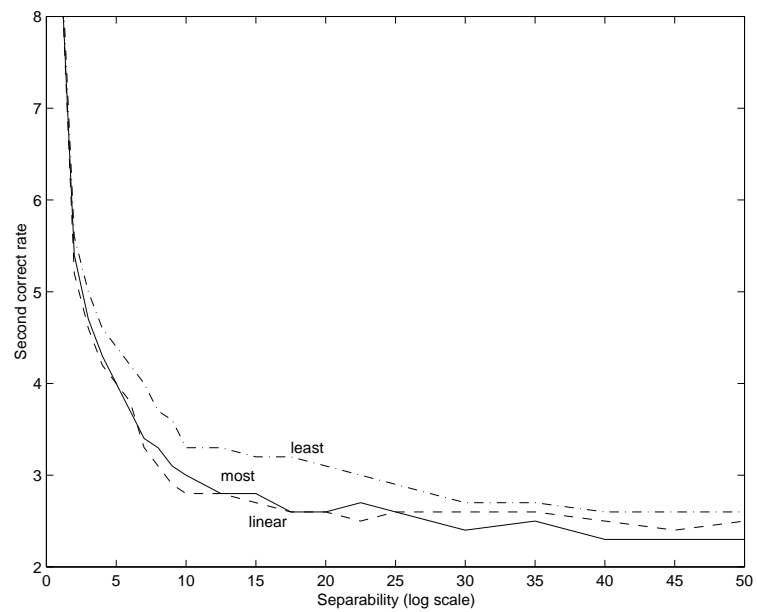
In order to see the effect of the gazing function, three different kinds of gazing algorithms have been examined as choosing a sub-image (i) which has the highest local separability, denoted as *most*, (ii) which has the lowest local separability, denoted as *least*, and (iii) which allows linear selection by index number order, denoted as *linear*. on the data set of EXP 1 with exactly the same settings using a quin tree of depth 4.

Figure 3.12 shows recognition performance in terms of top choice and second choice correct rate. Figure 3.13 shows resource usage in terms of the average of used nodes (sub-images) to the separability threshold in the three gazing strategies, and relation between the resource usage and the correct recognition rate of *most* gazing strategy.

As shown in figure 3.12 -(a), all the curves of top correct rate reach their steady

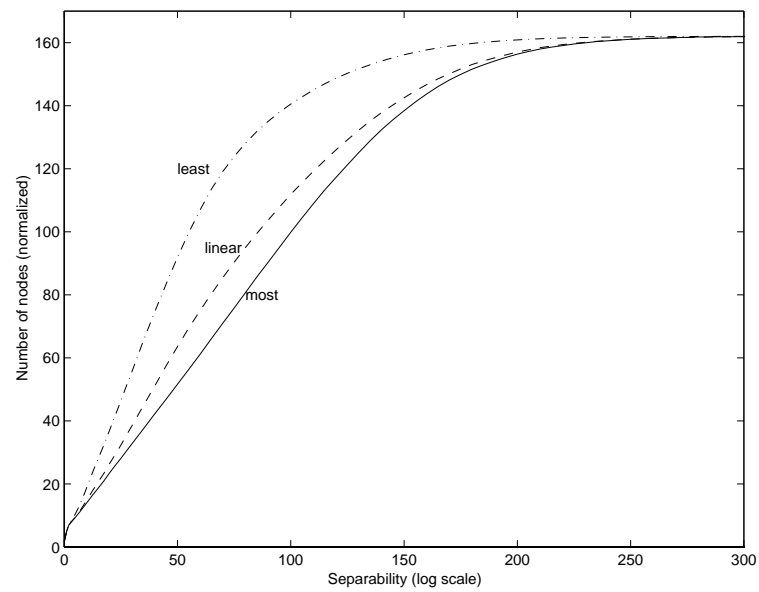


(a)

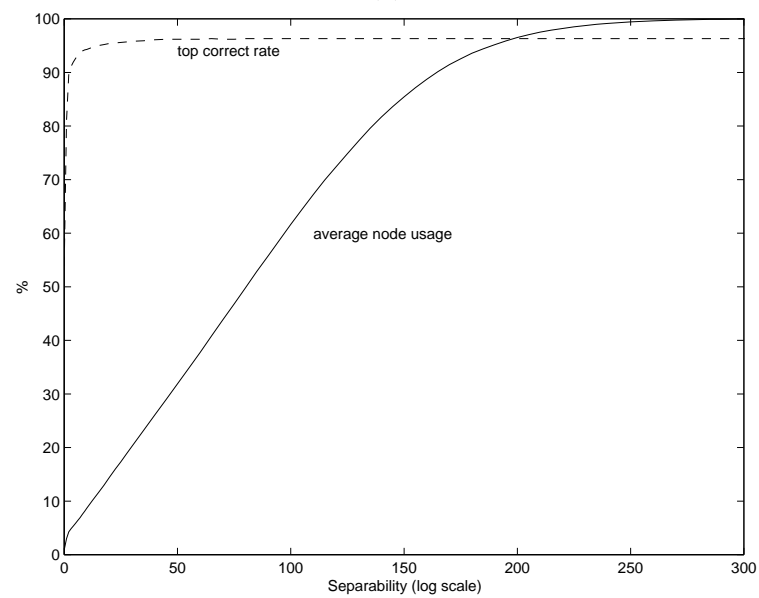


(b)

Figure 3.12: Dynamic recognition performance. (a) top choice correct rate and (b) second choice correct rate

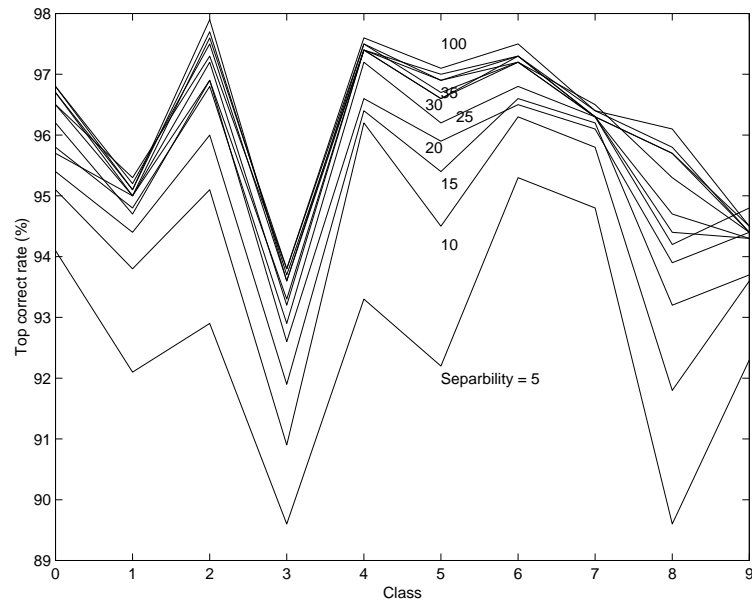


(a)

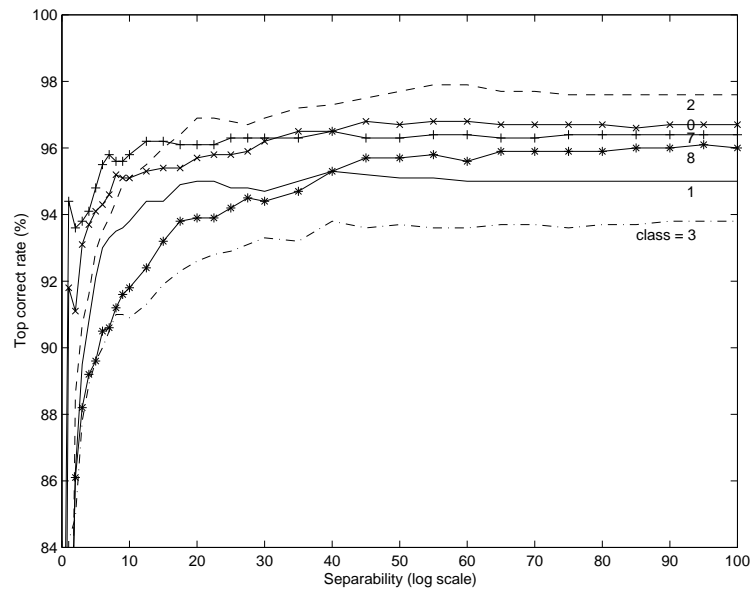


(b)

Figure 3.13: Average resource usage. (a) node usage in each searching mode and (b) normalized node usage vs. performance



(a)



(b)

Figure 3.14: Performance by class

state beyond 45 of separability threshold. However the minimum required resource usage to reach the steady state is different for the three gazing methods as shown in figure 3.13 -(a). The gazing strategy *most* uses minimum nodes to reach the same performance. For comparison in Figure 3.13-(a), the number of used nodes are normalized to the top correct rate respectively.

If we assume that a steady state is reached after separability threshold 45 in this experiment as shown in 3.12 -(b), we can obtain the top candidate correct rate using 25 % of node usage in average. Figure 3.14 shows classification performance by each class using the *most* gazing strategy.

3.4 Discussion

We have defined the confidence of a class based on the assumption that features in sub-images of an input pattern image are strongly correlated and the degree of correlation is uniquely characterized by the class. In other words, we expect the features to be highly correlated from the viewpoint of the true class and correlated to a lesser degree or not at all when considered from the viewpoint of the non-true classes.

For example, let us consider a set of feature vector codes extracted from sub-images of an input image. First, let us assume that the feature vectors are perfectly correlated to each other, that is, if we know one feature vector code, we can guess all the other vector codes without actual feature extraction. If confidences of classes are updated by multiplications of probabilities of vector codes extracted from other sub-images, confidence value for the true class will remain the same but those for the other classes will be decreased since the correlation degree is weak. Ideally, if the used sub-images become larger, the stronger separability between the true class

and the others can be obtained. But in actual application, perfect correlation can not be obtained. If vector codes which are extracted from sub-nodes are strongly correlated for the true class compared to other classes, we can obtain better separation between the true class and other classes by multiplying the probabilities of feature vector codes in sub-nodes, even though the confidence value of the true class might decrease.

In section 3.1.3 and equation (3.15) we keep updating the confidence values $Q_i(t)$ from successive and recursive cycles. This would indeed lower the updated confidence of a class. However, based on the correlation assumption made, the confidences of the classes that are not true should decrease at a more rapid rate. Hence the separability between the true class (top choice) and the non-true class (second choice) will increase. This justifies the subsequent recursive cycles and using the product of confidences in Equation 3.14. The assumption of correlation is based on empirical observations. It is supported by the recognition accuracy of the method which is comparable to those described in [44] [36] [39] that have trained and tested on similar, if not the same data sets.

Equations 3.24 and 3.25 do suggest that if an incorrect class is picked as the top choice in an early recursive cycle the error is going to be propagated to subsequent recursive cycles. Dashed lines in Figure 3.15 illustrate that a if mis-classification is made in an early recursive cycle confidence value of the top choice (not the true class) rapidly drops and the separability becomes unstable and remains near zero even though the number of sub-images used by the new recursion increases. On the other hand, when the top choice is the true class, the separability parameter holds steady (solid lines in Figure 3.15).

Since in our experiment, a simple decision strategy using only separability was

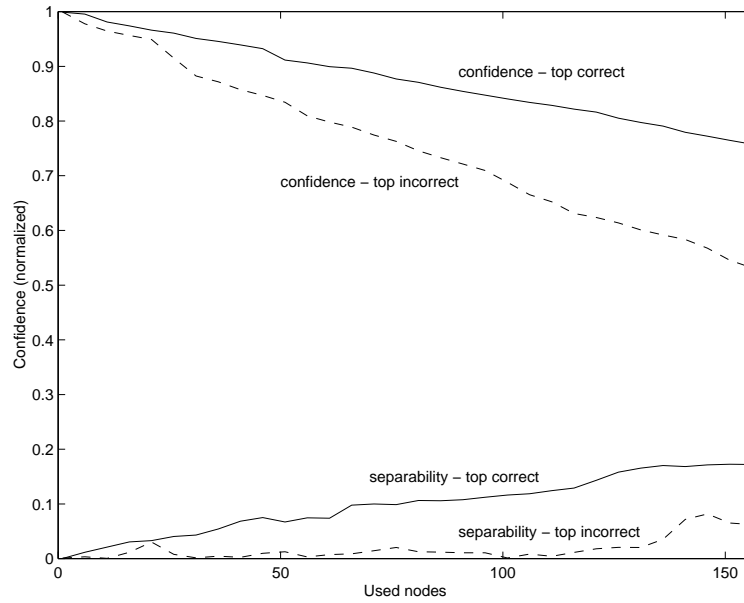


Figure 3.15: Confidence and separability changes - in EXP 1 with quin tree. (solid line shows a case of successful and dashed line shows a case of failure that true answer is not top ranked)

taken to see the behavior of proposed method, the type of error-generation mentioned above could not be prevented, even if the separability remained near zero. A more intelligent terminal criterion which considers classification results chronologically as well as horizontally throughout classes can prevent that type of error and can discard them at an early stage before reaching the exhaustive usage stage.

Chapter 4

Stroke Analysis

A *prime stroke analysis*, which detects global spacing style through an entire phrase, is presented in this chapter. Rule based stroke segmentation and prime stroke classification are used. The global spacing style is extracted as the dominant stroke frequency characteristic of a writing sample. This analysis is then applied to word segmentation and estimation of number of characters as a source of basic parameters. Finally, the gaps are identified as inter-word and inter-character gaps by simple thresholding or ranking.

Following a spacing model described in the next section, a prime stroke analysis method is presented involving the following steps: *normalization* which concerns slant and skew correction, *stroke segmentation* which separates word images into primitives, and *word break detection* which classifies primitives and estimates dominant prime frequency. Our analysis assumes that words are ordered from left to right in a series of horizontal text lines. However, unlike previous work in this area [75] [76], we assume that a single connected component can belong to more than one word.

4.1 Segmentation Model

The cue of spatial separation is most commonly used to perform word separation, but humans use several other cues such as presence of punctuation marks, presence of uppercase characters following a string of lower case characters, transition between numerals and alpha characters, and actual recognition of words prior to word separation [78].

In handwriting, variations in spacing are generated from difference in writing style and flourishes (Figure 4.1). Since handwritten text lacks uniform spacing, earlier methods compute the gaps between adjacent connected components with a given metric and sort the gaps in order of the metric [75] [76].

Previous approaches are simple, efficient and reasonably accurate in estimating word gaps but still leave several problems such as detecting global spacing style. Even though various inter-character and inter-word spacing styles exist in handwritten text due to different writing styles, the assumption that the space between two adjacent words is larger than that of characters is acceptable in general. The underlying problem of unconstrained word segmentation is not just the different way of measuring distance between components but the fact that accurate inter-character gaps are hard to detect.

Our primary goal in stroke segmentation is to detect possible character and word gaps to drive a word recognizer with a proper subset of lexicon entries. Segmentation is supposed to be precedent to recognition and detecting global spacing style in terms of the prime stroke period is used as the primary parameter. So the last three cues mentioned in the preceding paragraph can not be used in our proposed model. Only spatial and limited punctuation separation cues (, , -) are considered in our application.

ROANOKE PIKWY #507

(a)

Schooner Bay Lane

(b)

Union Street

(c)

Cheyenne Drive

(d)

Laurence Ave

(e)

Harmony Lane

(f)

Figure 4.1: Examples of various spacing style: (a)-(c) character based spacing, (d)-(f) stroke based spacing.

A writer's spacing style is generally ruled by the strokes or character bodies which are located within the primary horizontal writing band (imaginary) established throughout lines of handwritten text. The others, such as descenders, ascenders and decorative strokes are not ruled by the inter-character spacing style.

Figure 4.1 shows two types of spacing style. In a phrase which is written in hand printed or mixed style as the examples in Figure 4.1 (a)-(c), the immanent interval is ruled by characters in general. In cursive handwriting, the immanent interval is usually referenced by the distance between vertical strokes rather than by the interval of characters (Figure 4.1(d)-(f)). Generally word spacing is significantly larger than the space which is formed by the immanent interval between adjacent characters or vertical strokes.

If we define a *prime stroke* as a stroke which has a significant vertical component within the primary writing band (exclude flourish and ligature), we can generalize spacing style in terms of the prime strokes. The dominant immanent interval can be judged in term of *prime spatial periods* which is the most typical distance between prime strokes. Further, if the dominant prime stroke frequency is known, word break points can be detected based on the spatial period of stroke frequency assuming that word spacing is larger than character spacing.

Based on the above assumptions, an analytic handwriting spacing model can be considered as one where

- A character consists of one or more pieces of strokes called *primitives*.
- The primitives are classified into two categories, *prime*: a basic stroke of characters and *auxiliary*: such as ligature, hat or punctuation mark.
- A character has at least one prime primitive.

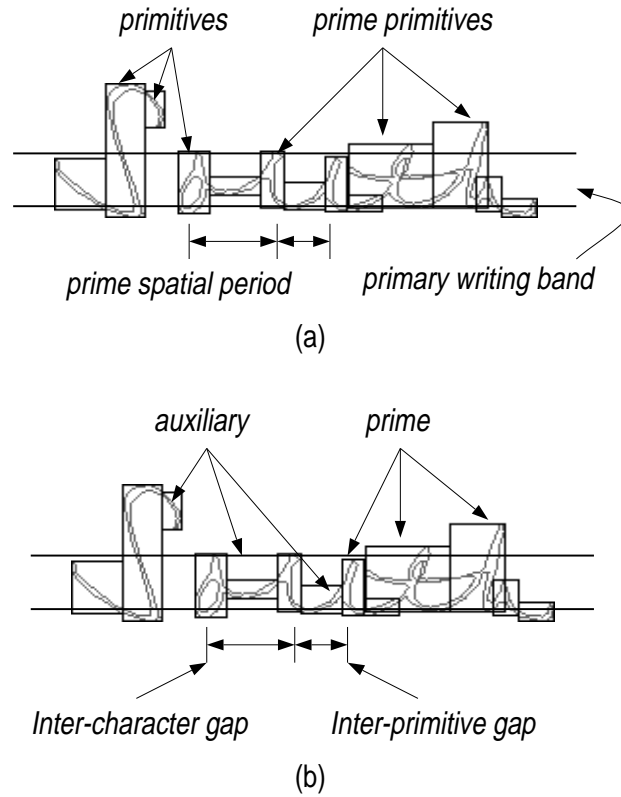


Figure 4.2: Prime spatial period (a) and types of stroke segments (b)

- The distance between two adjacent prime primitives of a character is equal or less than the dominant prime primitive frequency.
- Word gaps are larger than the prime spatial period.

As an image representation, the conventional contour method is used [99]. However for efficient resource usage, we use piecewise linearized contour model instead of pixel based contour (Figure 4.4 (b)). A phrase or word image is scanned and binarized. Binarized images are represented by piecewise linear contours.

Since the image scanner which is used has a fixed resolution regardless of the size of the target phrase image, the window size for generating piecewise linearized contour is adjusted by image density. The density is obtained from the ratio of black pixels to the minimum rectangle bounding box of whole phrase image.

4.2 Normalization

Since input images are assumed to be unconstrained handwritten phrase images (street name lines in address blocks of US mail pieces in experiments), there are several type of variations. The major sources of the variation are size(height), slant, skew(slope), stroke width, and rotation [58] [63] [64].

- *Size*: The size of letters will vary by writer's writing style even for the same task and for a given writer for different tasks.
- *Slant*: The slant is the deviation of vertical strokes from the true vertical, and it can vary between words and between writers.
- *Skew*: This is the angle of the baseline of a word if it is not written horizontally
- *Stroke width*: This depends on such factors as the writing instrument used, the pressure applied, and the angle of the writing instrument, as well as the paper type.
- *Rotation* : Usually this factor is caused by misalignment of the scanner and scanning target.

Minimization of these variations will help to boost the recognition performance. Assuming that rotation is very small or removed by scanning equipment

and size normalization can be incorporated into the character recognizer, slant and skew are the major targets of normalization. Stroke width is not a target of normalization but it is used as a reference parameter to segment strokes. Normalization processes are shown in Figure 4.3.

Slant Detection Vertical and near vertical lines are extracted by identifying *critical points* whose connection to the next critical point forms an angle with the horizontal axis which is larger than a threshold. The *critical points* are the endpoints of line segments in our piece-wise linear contour representation. In Figure 4.4 (c), such lines are shown as thick contour lines. The slant angle in an given phrase image is estimated by the average of all the angles formed by the vertical line segments, weighted by their length.

Skew Detection For skew angle detection, the baseline is estimated by extracting local minimum points on contours. Baseline estimation consists of the following steps:

1. Local minimum points which are on the *lower half contours* between leftmost and rightmost extremities of a component and located below of the center of mass of the whole image.
2. The line of best fit for the extracted local minimum points is estimated.
3. Average and variance of the distance between the fitting line and the local minimum points are obtained.
4. The local minimum points whose distance from the best fitting line exceeds a threshold expressed as a function of variance are ignored.

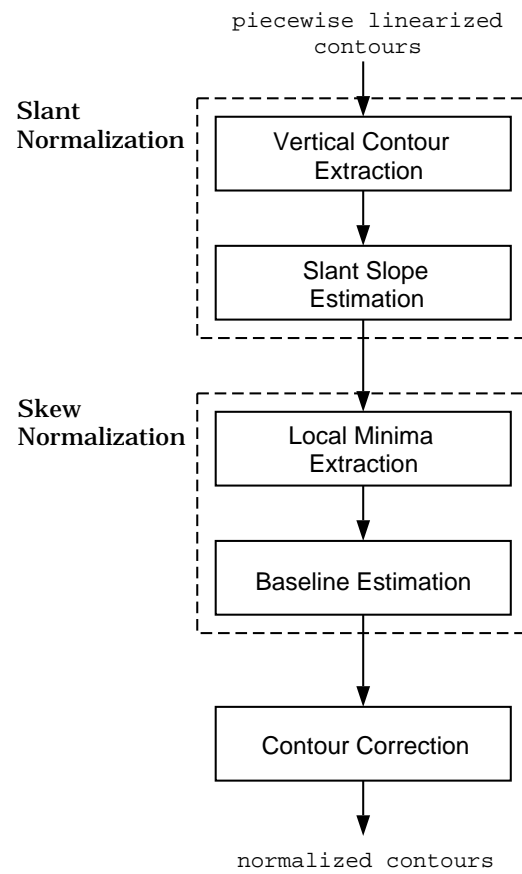


Figure 4.3: Normalization steps

5. The best fitting line is re-estimated by the remaining local minimum points.

The angle between baseline and horizontal line is assumed as the skew angle. Figure 4.4 (c) shows extracted vertical lines and local minimum points along linearized contours. Figure (d) shows contours after normalization process.

4.3 Stroke Segmentation

Stroke segmentation is performed in order to generate an array of primitives. From the the primitive arrays, prime frequency is estimated as a quantified reference of the writer's spacing style.

Connected components are split into more detailed symbolic primitives than what is done in character segmentation procedure. Stroke segmentation points are determined using features of convexity, concavity and ligature by tracing contours. Considering geometric relations among convexities, concavities, and ligatures, primary segmentation points are determined by a heuristic rule based algorithm. To extract more accurate geometric relations among strokes, the primary segmentation points are expanded along ligature strokes. These expansions can separate flourish and ligature strokes from basic primitives. Figure 4.5 shows the steps of stroke segmentation.

As shown in Figure 4.6 (b), using a predetermined horizontal sensing line, average stroke width is estimated by averaging of cross distance captured by contours. Convexities in the lower contour path and concavities in the upper contour path are identified as in (c). Using the vertical distance between the upper and lower contours and thresholding with the estimated average stroke width, ligature strokes are identified in (d).

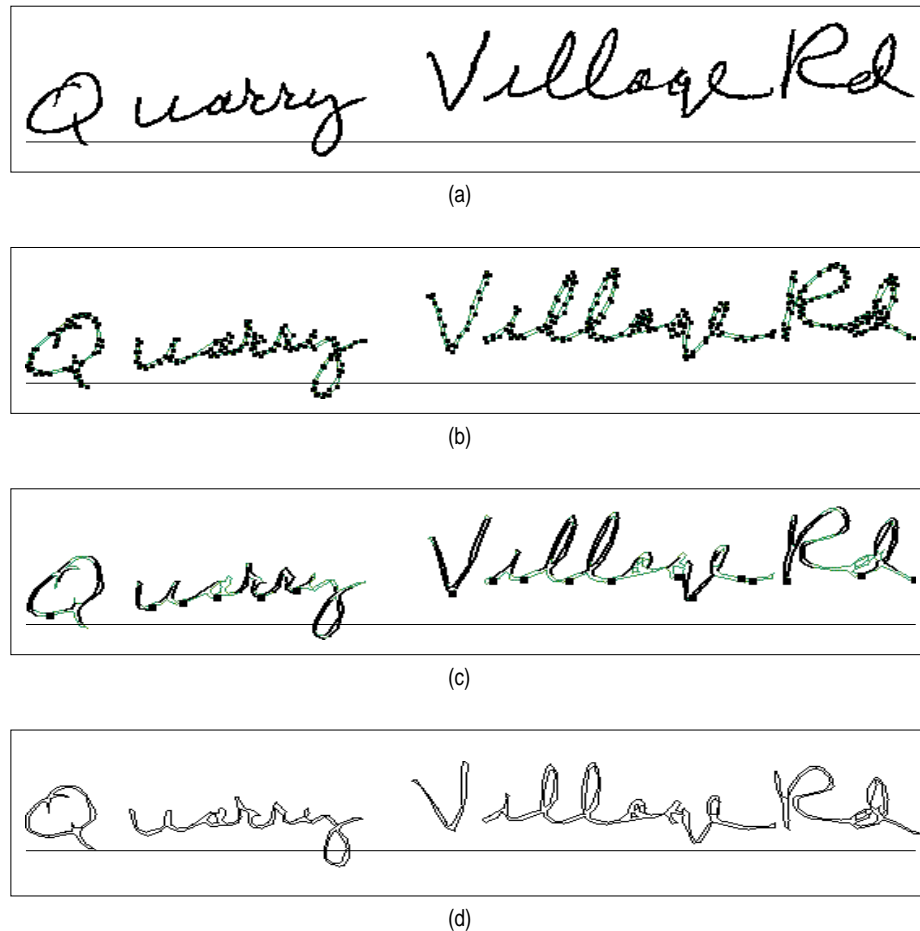


Figure 4.4: Normalization processing example, (a) a street line image, (b) critical points and piecewise linearized contours, (c) vertical contour components (thick lines) and filtered local minimum points (square dots), (d) slant and skew normalized contours

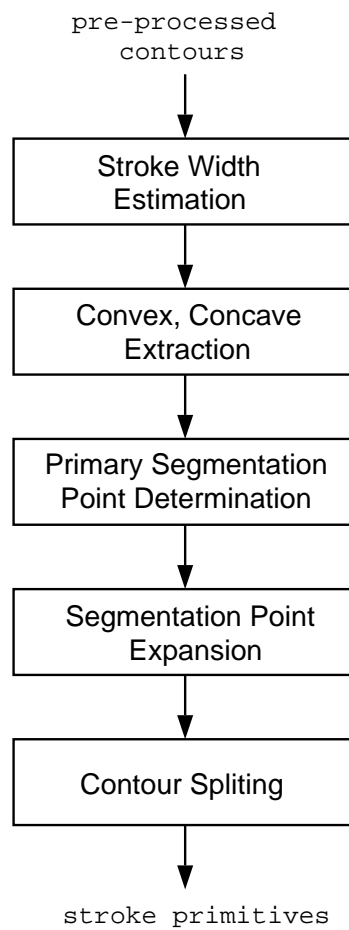


Figure 4.5: Stroke segmentation steps

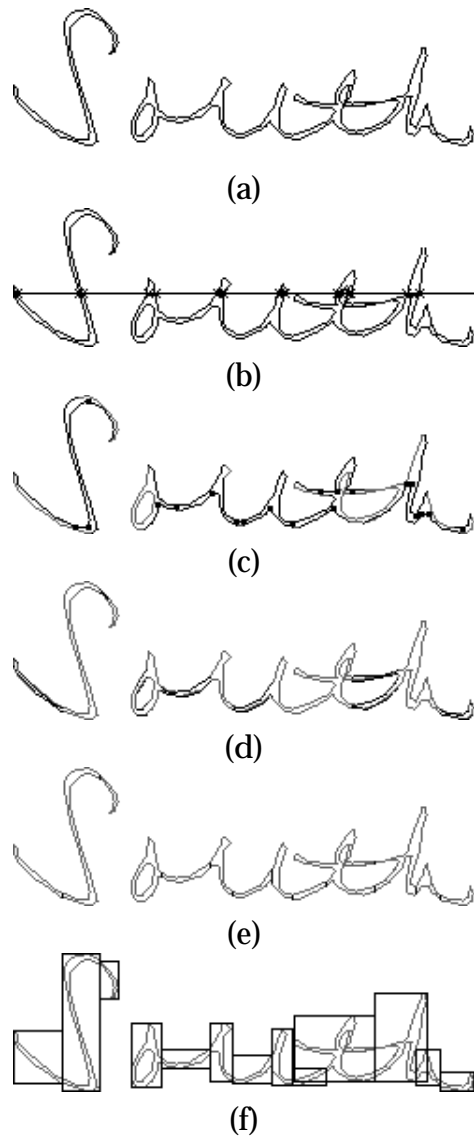


Figure 4.6: Stroke segmentation processing example: (a) normalized contours, (b) stroke width sensing, (c) extracted convexities and concavities (d) extracted ligature zone (e) final segmentation points (f) segmented primitives.

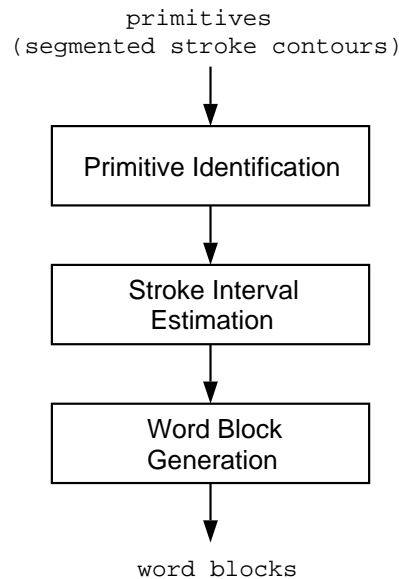


Figure 4.7: Word segmentation steps

4.4 Word Break Detection

As mentioned in the previous section, the writer's spacing style is generally dominated by the immanent interval of vertical strokes which stretch up and down along the angle of slant and a word space is usually larger than the immanent interval. If a dominant vertical stroke frequency can be extracted, spatial word break points can be detected based on the period of stroke frequency (called as the prime spatial period).

If we extend this assumption into the cases of hand-printed or mixed text, using meaningful stroke segments called prime primitives instead of vertical strokes, prime primitive frequency is used as a basic parameter to determine word break point following our assumption. Figure 4.7 shows the word segmentation process.

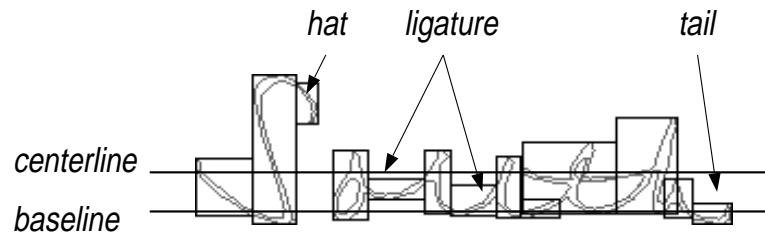


Figure 4.8: Primitive classification

4.4.1 Primitive Classification

After stroke segmentation, each split stroke primitive is classified into *Prime*, *Hat*, *Tail* and *Ligature* in order to separate auxiliary strokes from prime strokes as shown in Figure 4.8.

- *Prime*: a basic primitive comprising a character - most of the primitive body is laid on or crossed through the primary writing band.
- *Hat*: a primitive located above the basic writing band such as the bar in character 'T' or dot in character 'i'.
- *Tail*: a primitive located below the basic writing band such as the descenders in characters of 'y' or 'z' or comma.
- *Ligature*: a connecting stroke between two characters or inter horizontal stroke in cursive characters as 'w' or 'u'

Two reference lines, *halfline* and *baseline* are used to establish the horizontal primary writing band and *centerline* is used to estimate the two lines. The *centerline*

is a horizontal line passing through the global mass center of all contours. The *baseline* is the best fitting line of local minimum points of concavities of exterior contours below the center line (extracted from normalization step). The *halfline* is a line parallel to the baseline, the best fitting to local maximum points of convexities of exterior contours above the center line. These lines are extracted by angular histogram method after normalization and are assumed to be horizontal.

The *merging ratio* which is the ratio of contour component length within primary writing band to the whole length of contour component is used as the reference to separate both *Hat* and *Tail* from prime primitives. The ligatures are generated from horizontal pen flow of finishing one character to starting of next character. In general, components of ligature contours are horizontal runs and the structure is relatively simple as a single concavity over which the vertical distance between upper and lower contours stays within the average stroke width. However, since ligatures are usually located in the primary writing band, spatial reference is not useful. *Vertical component run ratio*, *ligature stroke ratio* and number of convexities and concavities throughout entire trace of primitive contours are used as features to classify ligature primitives.

The vertical component run ratio is a ratio of the run length of vertical contour components to the the whole run length of contour component. The vertical contour components are extracted in similar way of generating vertical and near vertical lines in slant detection of section 4.2. *Ligature zone* is established as the horizontal zone whose vertical distance between the lowest point of the upper contour path and the highest point of the lower contour path is less or equal than the average stroke width. The ligature stroke ratio is the ratio of the ligature zone to the whole horizontal span of contour component.

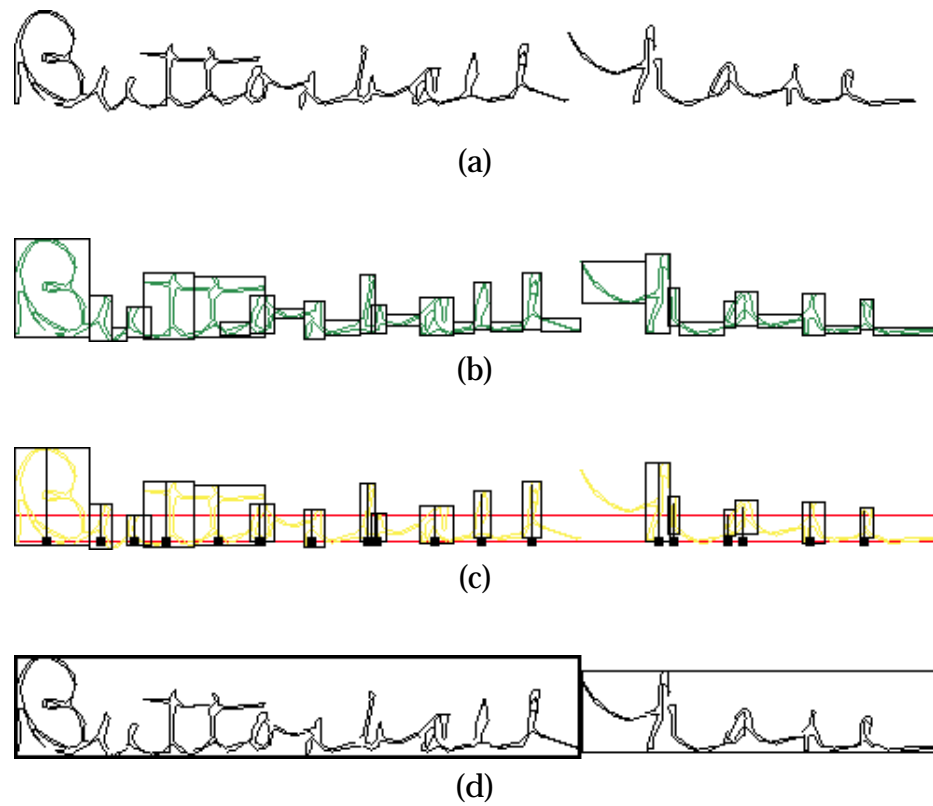


Figure 4.9: Word segmentation processing example: (a) slant and skew normalized contours, (b) segmented primitives, (c) stroke interval estimation, (d) generated word segments

4.4.2 Estimating Prime Frequency

Assume N_s stroke primitives are generated in a given word image and N_p primitives are classified as being of prime type. Let i be the index of a prime primitives as $i \in 0, 1, \dots, N_s - 2$ and (Z_{x_i}, Z_{y_i}) be the coordinate of the centroid of i^{th} prime primitive.

The interval of the i^{th} prime primitive p_i is

$$p_i = Z_{x_{i+1}} - Z_{x_i}. \quad (4.1)$$

Let the average and variance of intervals be \bar{p} and v_p^2 respectively. Since \bar{p} is the average of a mixture of inter-stroke and inter-character intervals, \bar{p} can not be directly taken as the period of prime frequency.

Let $T(p)$ be a filtering function as

$$T(p) = \begin{cases} p & \bar{p} - \gamma v_p \leq p \leq \bar{p} + \gamma v_p \\ 0 & \text{otherwise.} \end{cases} \quad (4.2)$$

where γ is a control parameter.

The prime period \tilde{p} is defined by the average value of filtered prime intervals as

$$\tilde{p} = \frac{1}{N_p} \sum_{i=0}^{N_s-2} T(p_i) \quad (4.3)$$

where N_p is the number of filtered prime intervals whose values are not zero. Since handwriting does not guarantee uniform spacing among strokes or characters, error rate \tilde{v}_p is defined by variance of filtered prime intervals to support the confidence of the extracted prime period.

$$\tilde{v}_p^2 = \frac{1}{N_p} \sum_{i=0, T(p_i) \neq 0}^{N_s-2} (T(p_i) - \tilde{p})^2 \quad (4.4)$$

Finally, a metric of word gap confidence q_i in i^{th} prime primitive is defined in terms of the *Maharanobis Distance* as

$$q_i = \begin{cases} 0 & p_i \leq \tilde{p} \\ \frac{p_i - \tilde{p}}{\tilde{v}_p} & \text{otherwise.} \end{cases} \quad (4.5)$$

After computation of all word gap confidences, the possible word break points are detected by simple thresholding. Figure 4.9 shows word segmentation examples after word break point detection.

4.5 Experiment

Street name images were generated from live mail pieces using *Handwritten Address Interpretation System* described in [7]. And the images which had been successfully isolated from the address block and the primary number (street number) field were selected manually as the experiment set. A total 1342 image set was used as the testing set and the first 150 images of the set were used in parameter tuning. The proposed word segmentation algorithm has missed actual word segmentation point in about 2% of all images, maintaining perfect word segmentation in 48% and over segmentation points in 50% of image set (Table 4.1).

Number of Images	1342	100%
Perfect Segmentation	643	47.9%
1 Over Segmentation	418	31.2%
2 Over Segmentation	181	13.5%
3 or More segmentation	70	5.2%
Fail(under-segmentation)	30	2.2%

Table 4.1: Segmentation Performance

Chapter 5

Recursive Word Recognition

One of the most common architectures of word recognition system is the serially connected two processing steps of recognition and decision making. In this architecture, a recognition module provides multiple choices of recognition results with associated confidence values and a decision making modules filters the raw recognition results to meet the required system performance. Since the flow of process is monotonous, the word recognition engines are always operated in using all the resource built in the recognizer regardless of recognition difficulty in order to generate the best recognition results. If several choices of raw recognition results are made available by a word recognizer, a decision maker chooses a recognition result to represent as the true answer. The goal of decision making is to meet the performance requirements that are imposed by the user of the recognition system.

The most cost-effective solution is hard to obtain without interaction between a recognition engine and a decision making module, since exhaustive processing

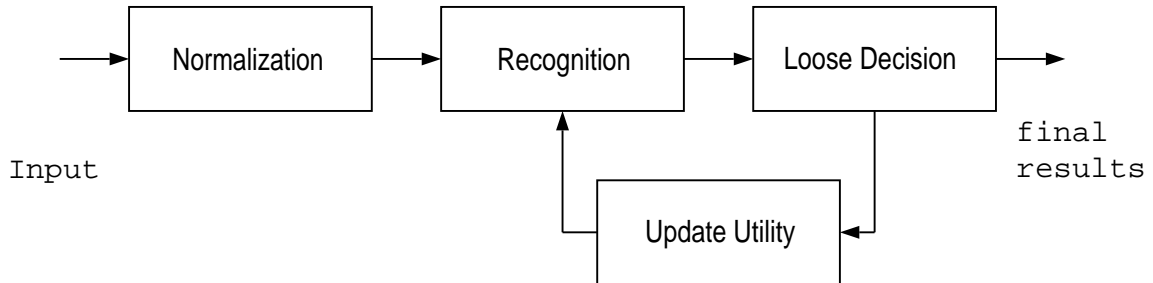


Figure 5.1: Recursive recognition

actions are seldom required. In order to achieve successful recognition with minimum required processing, an active combination of recognition and decision making is necessary since the relation between recognition and decision making is similar to the relation between an inquiry and seeking the answer. One of the simple solutions to reach a cost optimal status is to have a system in which the recognition and the decision making are connected in a closed loop as shown in Figure 5.1 which operates recursively with successive upgrades of recognition accuracy. The recursion can eventually reach a satisfactory terminal condition or a state of exhaustive use of all the resources.

In this chapter such a recursive recognition model is described which adopts two cognitive models of *time course word recognition* [85] and *epistemic utility theory* [88]. The described model basically follows the multi-processing levels of the time course word recognition model but parallel processing is modified to simplify recursive updating. The basic *satisficing decision* [91] concept which adapts *epistemic utility theory* is implemented in handwritten word recognition by using the confidence values generated by a character recognizer.

5.1 Applying Utility Theory

Let Ψ denote the full truth set and H a pre-selected subset of Ψ including the null set ϕ . The H consists of N_H possible true hypotheses.

$$H = \{h_0, h_1, \dots, h_{N_H-1}\} \quad (5.1)$$

In handwriting recognition, Ψ represents all possible combinations of characters and H is a meaningful lexicon set of N_H elements. H is derived from Ψ independently for the word recognizer. We assume that all lexicon entries are unique and are not subsets of other lexicon entries. And a lexicon entry h is represented by a length N_h array of symbols which belong to a possible character set,

$$h = [h_0, h_1, \dots, h_{N_h-1}] \quad (5.2)$$

where $h \subset C'$ and C' is an extended class set of C including the null set ϕ ,

$$C' = \{\phi, c_1, c_2, \dots, c_{N_c}\} \quad (5.3)$$

Let us denote by $A(h)$ an accuracy support function which generates an absolute distance measure associated with h . $L(h|H)$ be a function which characterizes liability exposure of h in given lexical entries H . The epistemic utility function $U(h|H)$ of each lexicon entry h in a lexicon H is defined by a convex combination of accuracy support A and liability exposure L [90], such as

$$U(h|H) = \alpha A(h) + (1 - \alpha)L(h|H) \quad (5.4)$$

where $\alpha \in [1, 0]$ and

$$U(h|H) = [u_0, u_1, \dots, u_{N_h-1}] \quad (5.5)$$

$$u_i = U(h_i|H) \quad (5.6)$$

After positive linear transformation the function is

$$U(\mathbf{h}|H) = A(\mathbf{h}) - \beta L(\mathbf{h}|H) \quad (5.7)$$

The parameter β is referred to the *index of rejectivity*. Large values of β imply an increased willingness to reject hypotheses. $\beta \cong 0$ implies minimal concern for liability, $\beta \cong 1$ implies equal concern for accuracy and liability.

If multiple symbols are assigned to an h , for instance, allowing folding of lower and upper characters, each epistemic utility value is defined by

$$u_i = \min_{\forall c \in h_i} (U(c|H)) \quad (5.8)$$

If we set $|U|$ to be a function which generates an overall epistemic utility value of a hypothesis \mathbf{h} as the average value of each epistemic utility, then $|U|$ becomes

$$|U(\mathbf{h}|H)| = \frac{1}{N_h} \sum_{i=0}^{N_h-1} U(h_i|H) \quad (5.9)$$

The minimum distance of each test pattern from the prototypes can be used as the value of accuracy support $A(h_{ij})$ of character hypothesis h_{ij} where i is a lexicon index and j is a index of the symbol sequence. In the above equation, if we set β as zero, $|U|$ becomes a conventional method of finding the average distance between a lexicon entry and a given image using a character recognizer which provides the distance in feature space.

The liability exposure represents the information value of rejecting the hypothesis in a pool of hypotheses. As a liability exposure function, $L(h_{ij}|H)$ for hand-written word recognition applications, *lexicon complexity* is defined by the degree of relative importance of a character h_{ij} in a lexicon entry \mathbf{h}_i against the lexicon set H .

The importance of each character in a lexicon entry depends not only on the lexicon but also on the separation characteristics of a character recognizer. For example, in two lexicon entries if all characters are the same symbols and in the same sequence except only one symbol, the utility of that character recognition result is very high for the purpose of separating those two lexicon entries. But the utility should be weighted by inter-class confusion of a character recognizer since inter-class confusion reflects imperfections of the character recognizer. If the recognizer being employed has high inter-class confusion for that particular character pair of previous example we can not succeed in separating the two elements of the lexicon.

5.1.1 Revision of Utility

The essence of a decision-making strategy lies in rejecting only bad hypotheses and leaving other hypotheses for further consideration. Our priority is avoiding actual errors rather than minimizing the probability of error. If we recursively reject only bad hypotheses, we will eventually reach the conditions of termination in which we can choose several type of results such as satisfying one answer, group of possible answers, or reject all hypotheses.

There are two ways in which a recognizer may want to modify its decision boundary without adopting a new process schedule from Equation 5.7. First, a word recognizer can renew the accuracy support of hypotheses using a more accurate character recognizer. Second, the recognizer can update the liability exposure by discarding bad hypotheses which are no longer believed to be true.

Since accuracy support basically depends on the distance measures returned by a character recognizer, it can be renewable by applying the character recognizer

in higher accuracy or a new character recognizer which is more accurate.

Liability exposure is renewable by considering the lexicon complexity of a reduced lexicon set. If a decision maker discards bad hypotheses from a lexicon as false, the lexicon complexity changes and the liability exposure values can be updated. If the decision maker no longer discards bad hypotheses, evidence should be upgraded to expand the decision boundary. After applying more accurate degree of classification, the liability exposure is also renewable using inter-character confusion of the character recognizer in enhanced accuracy.

For this control scheme, we can rewrite the epistemic utility function in Equation 5.7 as a time dependent equation,

$$U(\mathbf{h}|H(t)) = A(\mathbf{h}, t) - \beta L(\mathbf{h}|H(t)) \quad (5.10)$$

where t is a recursion time index. If the time index increases, a character recognizer which is used in the function $A(\mathbf{h}, t)$ applies more accurate degree of classification by increasing the feature extraction resolution. $H(t)$ denotes a narrowed hypotheses set at time t as determined by a decision-maker after rejecting bad hypotheses. The lexicon complexity $L(\mathbf{h}|H(t))$ is updated by the narrowed the lexicon subset $H(t)$ and enhanced degree of inter-character confusion.

5.2 Simulating Multi-Resolution in A Distance Metric

Recognizers using features extracted in multi-resolution have been found to be more accurate than conventional methods using features at a single scale [17] [36]. Features extracted in finer resolution have the advantage of compensating for poor discriminatory power of confused local shapes.

However, features in higher resolution is not always necessary to reach final

acceptance even though such features add more separation between hypotheses. Sometimes additional features in detailed level are redundant for the required classification depth and excessive features generate dimensionality problems. Moreover, they have disadvantage of increasing computational complexity and unregulated metric problems, because of the burden of additional features needed to accommodate the multi-resolution aspect of the feature space.

In Chapter 3, a hierarchical character recognizer is formulated using statistical tools, which preserves the benefits of a multi-resolution model while circumventing the disadvantages mentioned above. However, for a word recognition application, the statistical model of characters in a word does not usually correspond to the probability of general training characters, since words have their own contextual composition probability among characters such as *digram* and *trigram*.

Thus, we extend the hierarchical classification to a nearest neighbor method using a distance metric. For computational efficiency, templates obtained from training data are clustered and the primary metric is changed to minimum distance between templates and extracted feature vector instead of probability.

5.2.1 Multi-resolution Feature Space

Let C be a set of mutually exclusive classes of size N_c , and \mathbf{x} be an extracted feature vector from a supposed character image as in Chapter 3,

$$C = \{c_0, c_1, \dots, c_{N_c-1}\}, \quad (5.11)$$

$$\mathbf{x} = [x_0, x_1, \dots, x_{N_x-1}] \quad (5.12)$$

Also, consider a tree structure as Figure 3.5 where N_D is the maximum depth and N_B is the number of branches per node.

Let $\mathbf{x}(d, l)$ be an extracted feature vector from sub-image $I(d, l)$ from the tree where,

$$d \in \{0, 1, \dots, N_D - 1\}, \quad l \in \{0, 1, \dots, N_B^d - 1\}. \quad (5.13)$$

Features extracted in a larger region provide the global level features while those in a smaller sub-region provide the localized features in finer resolution. The multi-resolution of features is thereby achieved, although the actual feature extraction process and the type of the features remains the same.

Let $\bar{\mathbf{x}}_i$ be a template vector of class c_i which has the same dimension N_x and which has been obtained from a training set. The distance $\nu_i(\mathbf{x})$ between a feature vector \mathbf{x} and templates in a class c_i is defined by the minimum distance to any template vector of the class c_i ,

$$\nu_i(\mathbf{x}) = \min \delta_f(\mathbf{x}, \bar{\mathbf{x}}_i) \quad (5.14)$$

where $\delta_f(\mathbf{x}_i, \mathbf{x}_j)$ is a distance measure function between two feature vectors which satisfies the triangular inequality.

The distance of class c_i in a layer of a tree, denoted as $\dot{\nu}_i$, is defined by a weighted sum of the local distance measurements of class c_i generated from the sub-images of the layer as,

$$\dot{\nu}_i(d) = \sum_{l=0}^{N_B^d-1} \omega_L(d, l) \nu_i(\mathbf{x}(d, l)) \quad (5.15)$$

where $\nu_i(\mathbf{x}(d, l))$ is the local minimum distance of an extracted feature vector from sub-image $I(d, l)$ to the templates of class c_i , and $\omega_L(d, l)$ is normalized weight of l^{th} sub-image in layer d of a tree, which satisfies

$$\sum_{l=0}^{N_B^d-1} \omega_L(d, l) = 1 \quad (5.16)$$

Similarly, the distance of class c_i in a tree, denoted as $\ddot{\nu}_i$, is defined by a weighted sum of all the $\dot{\nu}_i$ in the tree as,

$$\ddot{\nu}_i = \sum_{d=0}^{N_D-1} \omega_T(d) \dot{\nu}_i(d) \quad (5.17)$$

where $\omega_T(d)$ is normalized weight of layer d of the tree, which satisfies

$$\sum_{d=0}^{N_D-1} \omega_T(d) = 1 \quad (5.18)$$

5.2.2 Distance in A Recursive Tree

Usually, localized features extracted in finer resolution are useful for enhancing separability of a confused subset in which global features can not provide adequate discrimination. But, empirically, it is realized that localized features are not always useful for better separation, because of the poor homogeneity in patterns of the same class.

Consider a hierarchical image tree and a sub-tree which is a partial tree of the hierarchical image tree. Let us assume that a sub-tree is obtained by successive expansion of a parent image to child images as sub-images by a quad or quin division, and feature vectors extracted from the sub-tree provide the minimum discrimination which satisfy an accepting terminal criterion. We limit a sub-tree can always be selected within a tree, starting from the root of the tree extending to the sub-nodes successively. The distance of a class c_i in a sub-tree G , denoted as $\tilde{\nu}_i(G)$, is defined by a weighted sum of local distance of class c_i generated from the sub-nodes which belong to G as,

$$\tilde{\nu}_i(G) = \frac{\sum_{\forall(d,l) \in G} \omega_T(d) \omega_L(d,l) \nu_i(\mathbf{x}(d,l))}{\omega(G)} \quad (5.19)$$

where

$$\omega(G) = \sum_{\forall (d,l) \in G} \omega_t(d) \omega_L(d, l) \quad (5.20)$$

Consider a hierarchical feature tree which is extracted from a hierarchical image tree of a given image. Let us assume the feature tree matches perfectly one of the template feature trees of the true class which are built in a classifier, and that there is no identical tree pair in the template trees. Also let us assume that a classifier is caught in a situation that the top two class separability of the distance using global features of the root node does not satisfy the acceptance criterion. Then, the updated distance using sub-nodes will remain zero for the true class but the updated distance to the other classes will be increased, since the generated local distance of sub-nodes will be added up to the weight of the sub-nodes for the non-true classes. In an ideal case, the more nodes which are used, the bigger the separability which can be obtained.

But in actual applications, a tree which is identical to one of template trees is seldom obtained. Thus, if we assume that an extracted feature tree is close enough to one of template trees of the true class, but relatively not close to those of non-true classes, we can obtain better separation between classes by using sub-node expansion, even though the distance to the true class might increase.

Let a function $g_J(t)$ determine which node is to be examined in the next recursion based on a potential generating function J . The $g_J(t)$ is defined to return a sub-node index (d, l) as an output at time t from the pool of expandable nodes. If we choose $g_J(t)$ as a gaze function that finds the most useful node within the pool of *extreme nodes* for updating class distances similar to the gaze function described in Section 3.1.3, the separation will rapidly approach to the terminal criterion. An extreme node is defined by a child node of examined nodes but unused until t .

If we assume that a sub-tree is recursively expanded by only one node per recursion time, we can rewrite $\tilde{\nu}$ of Equation 5.19 as a recursive equation,

$$\check{\nu}_i(t) = \frac{\check{\omega}(t-1)\check{\nu}_i(t-1) + \omega_T(d)\omega_L(d,l)\nu_i(\mathbf{x}(d,l))|_{(d,l)=g_J(t)}}{\check{\omega}(t)} \quad (5.21)$$

where

$$\check{\nu}_i(t) = \tilde{\nu}_i(G_t)|_{G_t=G(t)} \quad \text{and} \quad \check{\omega}(t) = \omega(G_t)|_{G_t=G(t)} \quad (5.22)$$

$G(t)$ denote a sub-tree expanded until time t and $\check{\nu}_i(t)$ denotes a recursive distance function of class c_i at time t . Since the sub-node expansion is finite, the recursion index t is only effective in the range of

$$t \in \{0, 1, \dots, \sum_{d=0}^{N_D-1} N_B^d - 1\} \quad (5.23)$$

The recursive equations of the weight function and the recursive sub-tree are

$$\check{\omega}(t) = \check{\omega}(t-1) + \omega_T(d)\omega_L(d,l)|_{(d,l)=g_J(t)} \quad (5.24)$$

and

$$G(t) = G(t-1) \uplus I(d,l)|_{(d,l)=g_J(t)} \quad (5.25)$$

where \uplus denotes the grafting operation of a sub-node I to a tree.

We define a potential measurement of an extreme node, denoted as $J(d,l)$, by weighted separability between the global top two choices of the parent node as,

$$J(d,l) = \omega(d)\omega(d,l)(\nu_{second}(\mathbf{x}(d_p, l_p)) - \nu_{top}(\mathbf{x}(d_p, l_p))) \quad (5.26)$$

where c_{top} and c_{second} is the global top two choices at the time t , and (d_p, l_p) is the index of parent node of (d, l) . After finding all potential measurements of extreme nodes, an extreme node which has the largest potential measurement is selected as an output of $g_J(t)$ to be examined in the next recursion. If global top two choices are changed, potential measurements in all extreme nodes are updated, otherwise remain the same.

5.2.3 Word Recognition Accuracy

A time dependent accuracy support of a lexicon entry \mathbf{h} is defined by an array of minimum distance of the individual class subset at time t as,

$$A(\mathbf{h}, t) = [A(h_0, t), A(h_1, t), \dots, A(h_{N_h-1}, t)] \quad (5.27)$$

where

$$A(h, t) = \min_{\forall c_i \in h} \check{\nu}_i(t) \quad (5.28)$$

The conventional recognition distance of a lexicon entry is given by the average of character recognition distance in the lexicon entry as,

$$|A(\mathbf{h}, t)| = \frac{1}{N_h} \sum_{i=0}^{N_h-1} A(h_i, t) \quad (5.29)$$

From the above equations, recursion depth t limits the usage of a hierarchical feature tree for individual character recognition. Since a recursive sub-tree is a partial tree of a hierarchical feature tree which starts from the root and expands to sub nodes successively, at $t = 0$, classification results are obtained in the most global view (coarse resolution). By increasing recursion depth, more built-in feature templates are used by comparing to the extracted features in finer resolution. From the recursive equation of 5.21, an active classification which has multi-resolution and minimal usage of resources can be simulated.

In segmentation driven word recognition in which character matching is limited on a segmented stroke frame, each independent sub-tree is built based on the image generated from each possible combination of stroke segments for character matching of a lexicon entry. But all characters which belongs to any of the character subset of all lexicon entries associated with the character image boundary of a word image share the same sub-tree.


5.3 Lexicon Complexity

A lexicon driven word recognizer finds the best matches between symbol strings of a given lexicon and an array of segmented strokes in a given images. In [59], a two step dynamic matching scheme has been implemented. Comparisons are made between feature vectors extracted from the possible combinations of segments and reference prototypes of codewords in order to find the best matching path. In the second phase, a global optimum path is obtained by Dynamic Programming using the scores obtained from the previous step.

However, the recognition using only visual evidence, *i.e.* distance between feature vector extracted from a given image and reference prototypes built in the recognizer, cannot always be successful due to the limitations of information and resources available during training. Some of these problems can be overcome by using contextual knowledge which can be obtained from given selectable choices, *i.e.* lexicon entries.

Let us consider an input image and three different lexicons as shown in Figure 5.2, and consider three different situations in which input images are the same one but the lexicons are given by Lexicon 1, 2, and 3 respectively. For explanatory convenience, a phrase image and phrase lexicons are shown and the situations are denoted by the name of the lexicon sets.

Initially, all the three lexicon entries of each lexicon and their constituent characters in the three situations are of equal importance. Let us assume that the lexicon entry ‘Sunset Avenue’ is commonly rejected at the first recognition recursion in all the three situations because of the word length. After rejecting ‘Sunset Avenue’ in *lexicon 1*, the first two characters ‘We’ and ‘Ea’ are more informative (useful) than the other characters of the two lexicon entries. In the situation with *lexicon 2*, the



Lexicon 1: West Central Street
 East Central Street
 Sunset Avenue

Lexicon 2: West Central Street
 West Main Street
 Sunset Avenue

Lexicon 3: West Central Street
 West Central Avenue
 Sunset Avenue

Figure 5.2: An example of lexicon complexity

substring ‘Central’ is more informative than the rest of the remaining contents of lexical entries for choosing ‘West Central Street’ as the truth. In the situation with *lexicon 3*, the suffix ‘Street’ is the most informative than any other contents.

If the recognizer does not provide strong separation between the substrings ‘We’ and ‘Ea’ in *lexicon 1*, it can not be confident in accepting ‘West Central Street’ as the final decision. In the same way, if the recognition engine provides a strong separation between ‘Street’ and ‘Avenue’ in *lexicon 3*, ‘West Central Street’ can be accepted as a final result even if the overall confidence scores are close in the two lexicon entries.

As shown in the example, considering relation among lexicon entries is sometimes useful for decision making of recognition results. A liability of each character's recognition results which describes importance of characters in a lexicon must be considered to maximize the performance. However the liability depends not only on the lexicon entries also on the recognition quality of a recognizer. since the basic lexical entropy is given by a lexicon but the verifying results come from a recognition engine built in the word recognizer. Lexical entropy in view of separation difficult in a lexicon and reliability of a character recognizer serve to capture the dynamism of the environment. A parameter, named as *lexicon complexity* is defined by measuring these two factors.

The distance between two finite character arrays can be measured by *Edit Distance* [97]. Given two strings, the edit distance is defined as the minimum cost of transforming one string to the other through a sequence of weighted edit operations. The edit operations are usually defined in terms of *insertion*, *deletion*, and *substitution* of symbols. Pure lexical entropy of a lexicon set can be evaluated considering edit distance between lexicon entries of the lexicon in a character transform cost model. However in word recognition, especially in segmentation driven word recognition, since the word recognition is driven by a character recognizer in a segmentation framework, the simple cost model between characters is not accurate enough to express the actual cost of character transformation. We extend the edit distance in order to let it work in a segmentation frame by adding three edit operations; *splitting*, *merging* and *converting*.

5.3.1 Edit Cost in A Segment Frame

Consider two different lexicon hypotheses, \mathbf{h}_i and \mathbf{h}_j over a finite alphabet, of length N_{h_i} and N_{h_j} respectively, given as

$$\mathbf{h}_i = [h_{i,0}, h_{i,1}, \dots, h_{i,N_{h_i}-1}] \quad (5.30)$$

$$\mathbf{h}_j = [h_{j,0}, h_{j,1}, \dots, h_{j,N_{h_j}-1}] \quad (5.31)$$

Let us assume that a stroke segmentation algorithm is applied to a given image and a stroke segment array \mathbf{s} of length N_s is generated as

$$\mathbf{s} = [s_0, s_1, \dots, s_{N_s-1}] \quad (5.32)$$

Also assume that a Dynamic Programming based word recognizer similar to [59] finds the best matching path of a lexicon entries within a segment array conforming to rules such as (i) at least one character matches one stroke segment, (ii) one character can only match within a window, and (iii) lexicon entries with a length smaller than N_s or the maximum allowable number of characters of the stroke array are discarded.

We can define the edit operations between the best matching paths of two different lexicon entries without altering the framework of matching stroke segments as follows

- **nullifying:** $r(h_{i,m} \rightarrow \phi | s_{a \sim b})$
- **substituting:** $r(h_{i,m} \rightarrow h_{j,n} | s_{a \sim b})$
- **splitting:** $r(h_{i,m} \rightarrow [h_{j,n_1}, \dots, h_{j,n_2}] | s_{a \sim b})$
- **merging:** $r([h_{i,m_1}, \dots, h_{i,m_2}] \rightarrow h_{j,n} | s_{a \sim b})$
- **converting:** $r([h_{i,m_1}, \dots, h_{i,m_2}] \rightarrow [h_{j,n_1}, \dots, h_{j,n_2}] | s_{a \sim b})$

where $r(h_m \rightarrow h_n | s_{a \sim b})$ is a transform function from a symbol subset h_m to h_n within stroke segments array of s_a to s_b . The m and n denote indices of characters within a lexicon entry, i.e. the $h_{i,m}$ denotes the m^{th} character in lexicon entry i where $h_{i,m} \in C'$. The a and b are indices of the stroke segment array, and $s_{a \sim b}$ denotes any character combination of stroke segment s_a through s_b where $0 \leq a \leq b \leq N_s - 1$.

Nullifying is a transform operation for a character match to certain boundary of stroke segments to be null. *Substituting* is a replacement of a character match to another within exact same matching bound. *Splitting* is a split of a character match into several character matches with the segment frame and *merging* is an inverse operation of splitting. *Converting* is transforming a string of multiple character matches into another string matches when matching bounds of characters in the two strings are cross each other within the segment frame.

The overall matching distance of a lexical entry to an array of stroke segments depends primarily on the distance measurement between a extracted feature vector of one or more combined stroke segments and corresponding prototypes. Let $\varepsilon(c_m \rightarrow c_n | s_{a \sim b})$ be a elementary transform cost function from c_m to c_n and $\delta_h(c | s_{a \sim b})$ be a character recognition distance function of c within stroke segments s_a and s_b .

The transform cost functions are defined as

$$\varepsilon(h_{i,m} \rightarrow \phi | s_{a \sim b}) = \delta_h(h_{i,m} | s_{a \sim b}) \quad (5.33)$$

$$\varepsilon(h_{i,m} \rightarrow h_{j,n} | s_{a \sim b}) = \delta_h(h_{i,m} | s_{a \sim b}) - \delta_h(h_{j,n} | s_{a \sim b}) \quad (5.34)$$

$$\begin{aligned} \varepsilon(h_{i,m} \rightarrow [h_{j,n_1}, \dots, h_{j,n_2}] | s_{a \sim b}) &= \delta_h(h_{i,m} | s_{a \sim b}) - \min(\delta_h(h_{j,n_1} | s_{a \sim a_1}) + \\ &\delta_h(h_{j,n_1+1} | s_{(a_1+1) \sim a_2}) + \dots + \delta_h(h_{j,n_2} | s_{a_{n_2-n_1-1}+1 \sim b})) \end{aligned} \quad (5.35)$$

$$\begin{aligned} \varepsilon([h_{i,m_1}, \dots, h_{i,m_2}] \rightarrow h_{j,n} | s_{a \sim b}) = \min(\delta_h(h_{j,m_1} | s_{a \sim a_1}) + \\ \delta_h(h_{j,m_1+1} | s_{(a_1+1) \sim a_2}) + \dots + \delta_h(h_{j,m_2} | s_{a_{m_2-m_1-1}+1 \sim b})) - \delta_h(h_{j,n} | s_{a \sim b}) \end{aligned} \quad (5.36)$$

$$\begin{aligned} \varepsilon([h_{i,m_1}, \dots, h_{i,m_2}] \rightarrow [h_{j,n_1}, \dots, h_{j,n_2}] | s_{a \sim b}) = \min(\delta_h(h_{j,m_1} | s_{a \sim a_1}) + \\ \delta_h(h_{j,m_1+1} | s_{(a_1+1) \sim a_2}) + \dots + \delta_h(h_{j,m_2} | s_{a_{m_2-m_1-1}+1 \sim b})) - \min(\delta_h(h_{j,n_1} | s_{a \sim a_1}) + \\ \delta_h(h_{j,n_1+1} | s_{(a_1+1) \sim a_2}) + \dots + \delta_h(h_{j,n_2} | s_{a_{n_2-n_1-1}+1 \sim b})) \end{aligned} \quad (5.37)$$

where $a_1, a_2, \dots, a_{n_2-n_1-1}$ are indices of stroke segments within $s_a \sim s_b$ satisfying $a \leq a_1 \leq \dots \leq a_{n_2-n_1-1} \leq b$.

5.3.2 Matching Distance

Given a transforming path from \mathbf{h}_i to \mathbf{h}_j , the matching distance has been defined directly in terms of transform cost of the paths (or traces)

$$r(\mathbf{h}_i \rightarrow \mathbf{h}_j) = \mathbf{r}_{i,j} = [r_{i,j,0}, r_{i,j,1}, \dots, r_{i,j,N_{h_i}-1}] \quad (5.38)$$

where $\mathbf{r}_{i,j}$ is a path of transform operations of symbols from \mathbf{h}_i to \mathbf{h}_j within a segment framework of \mathbf{s} .

The associated transform cost for a path is

$$\lambda(\mathbf{r}_{i,j}) = \sum_{k=0}^{N_{h_i}-1} \varepsilon(r_{i,j,k}) \quad (5.39)$$

where k is a index of matching bound array of the path,

$$k \in \{(s_0 \sim s_{b_0}), (s_{a_1} \sim s_{b_1}), \dots, (s_{a_{N_{h_i}-1}} \sim s_{N_s})\}, \quad a_{i+1} - b_i = 1 \quad (5.40)$$

Let $\bar{\mathbf{r}}_{i,j}$ be the transform path of minimum cost from \mathbf{h}_i to \mathbf{h}_j as

$$\bar{\mathbf{r}}_{i,j} = [\bar{r}_{i,j,0}, \bar{r}_{i,j,1}, \dots, \bar{r}_{i,j,N_{h_i}-1}] \quad (5.41)$$

Then the matching distance from \mathbf{h}_i and \mathbf{h}_j is defined as

$$\lambda(\bar{\mathbf{r}}_{i,j}) = \min(\lambda(\mathbf{r}_{i,j})) \quad (5.42)$$

A overall matching distance of a character of a lexicon entry can be evaluated by the average of minimum transform costs of the character to all the other lexicon entries in the lexicon. The overall matching distance of a true character matched in a portion of a segment array can be minimized in negative value, if corresponding characters of other lexicon entries matched to the same segment portion are all different and a character recognition engine generates maximum separation between them. Otherwise, if most of corresponding characters of other lexicon entries in the same portion are the same true class or the classifier does not generate good separation between them, the overall matching distance increase. In non true class it may become a positive value. Thus, the overall matching distance reflects the importance of a character in a lexicon as well as the separation degree provided by a recognition engine.

Lexicon complexity of a lexicon hypothesis is defined by an overall transform distance array of the character string of the lexicon hypothesis to all the other lexicon hypotheses. The liability exposure $L(\mathbf{h}_i|H)$ can be evaluated by the lexicon complexity.

$$L(\mathbf{h}_i|H) = -\frac{1}{N_H - 1} \sum_{j=0, j \neq i}^{j=N_H-1} \varepsilon(\bar{\mathbf{r}}_{i,j}) \quad (5.43)$$

The absolute recognition distance of a lexicon entry in Equation 5.29 is the minimum distance measurement between feature vector and templates of the best path within the stroke segment array. This lexicon complexity of a lexicon entry reflects the lexical (or character) entropy of a lexicon entry projected into a lexicon in terms of matching distance.

Time dependent equation is defined as

$$L(\mathbf{h}_i|H(t)) = -\frac{1}{N_{H(t)} - 1} \sum_{j=0, j \neq i}^{j=N_{H(t)}-1} \varepsilon(\bar{\mathbf{r}}_{i,j}) \quad (5.44)$$

where $H(t)$ is a narrowed lexicon subset at time t .

The overall lexicon complexity is

$$|L(\mathbf{h}_i|H(t))| = \frac{1}{N_{H(t)} - 1} \sum_{j=0, j \neq i}^{j=N_{H(t)}-1} \lambda(\bar{\mathbf{r}}_{i,j}) \quad (5.45)$$

5.4 Implementation

A segmentation based word recognizer is chosen as the base implementation model of our approach. In one of the simplest approach of segmentation based word recognition, after preprocessing, a character segmentation algorithm is applied to a given word image and a character recognizer is applied to each segment assuming that perfect character separation is achieved. However since perfect character segmentation is not always guaranteed, an over-segmentation strategy (called as *stroke segmentation*) is usually adopted to avoid character segmentation errors which feed forward to character recognizer.

After segmentation procedure, possible combinations of segmented strokes are sent to a character recognizer to find matching scores of characters, and this scores are used as a metric in the next word matching. Various word matching schemes have been implemented according to the metric used in a character classification and the method of handling a lexicon.

In our implementation, a matching scheme using Dynamic Programming and a distance metric are used to find the best match between a symbol string and segmented strokes. Our implementation follows the flow of segmentation and

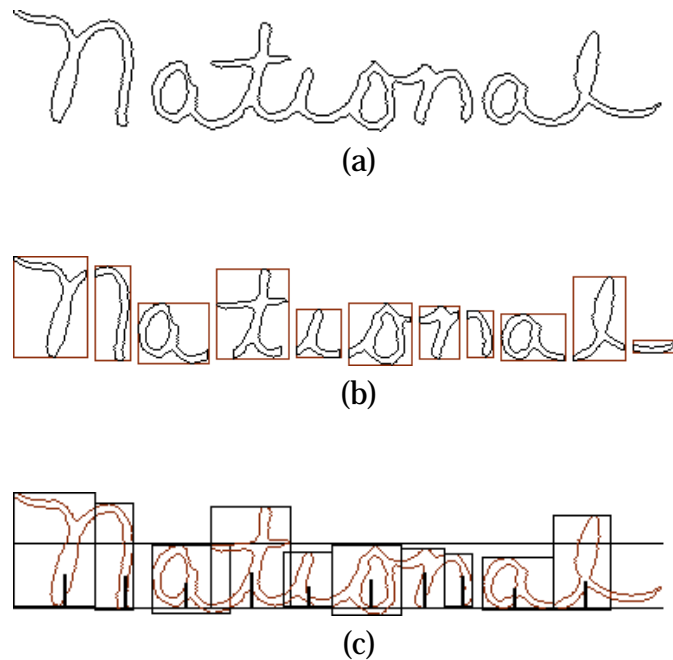


Figure 5.3: Stroke segmentation: (a) normalized contours, (b) segmented primitive array and (c) prime stroke period estimation

recognition after preprocessing but it has a recursive loop between recognition and decision making through an updating step.

5.4.1 Stroke segmentation

In the segmentation stage, word components are split into more detailed symbolic primitives using the method of Chapter 4. An example of segmentation results is shown in Figure 5.4.1. Then possible combinations of stroke primitives are sent to a character recognizer.

A *prime stroke analysis*, which have been introduced in Chapter 4 is used as a

basic parameter provider. A dominant stroke interval period, *prime spatial period* is estimated in a given word image and used as a reference for choosing a character matching window size. For computational efficiency, a permissible window size for character matching is generated by a function of the number of characters in a lexicon, the number of prime segments in a given image, and extracted prime stroke period from the image.

5.4.2 Features

Features are directly generated from a pixel based contour representation of a character [99]. A given image is divided into N_f by N_f grids in equal area where N_f is the size of divisions. A N_g slots histogram of gradient of contour components are obtained at each grid and an array of the histogram measurements of all grids are used as a feature vector. Two global features, aspect ratio of bounding box and a vertical/horizontal component ratio on overall contours are included in the feature vector in order to prevent mis-classification of non-character components. A $(N_f \cdot N_f \cdot N_g + 2)$ dimensional feature vector is extracted for a given image or a sub-image.

Sub-nodes of a hierarchical feature tree are generated by quin or quad division shown in Figure 3.5. Sub-images are obtained by dividing the parent image by the vertical and horizontal lines which cross at the centroid of the contours as shown in Figure 5.4.2 in order to minimize the effect of diversity of handwritten characters such as character size, aspect ratio, etc. The feature extraction procedures on a sub-image are exactly the same as those applied on the whole image. But the feature extraction for sub-images is only applied to a sub-image selected by gaze function, if it is deemed necessary for updating character recognition.

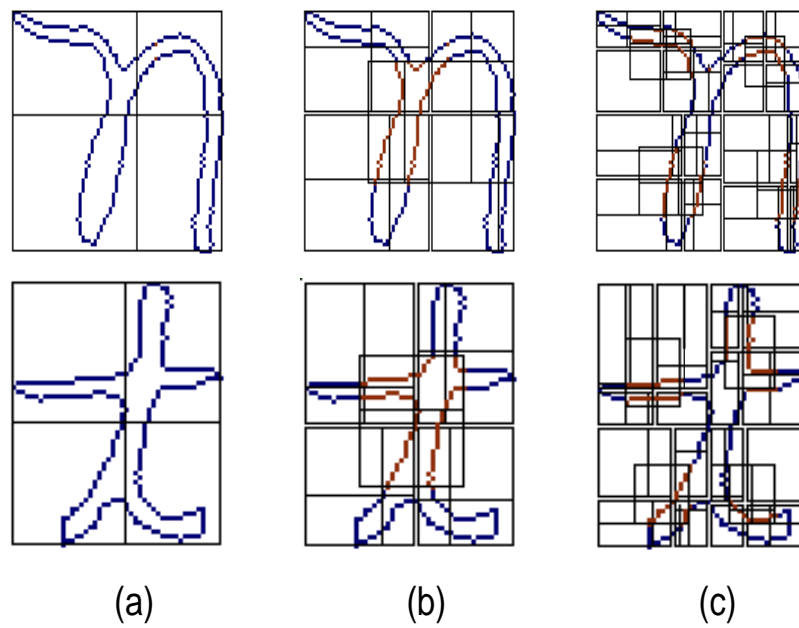


Figure 5.4: Sub-image division in a grouped character segment

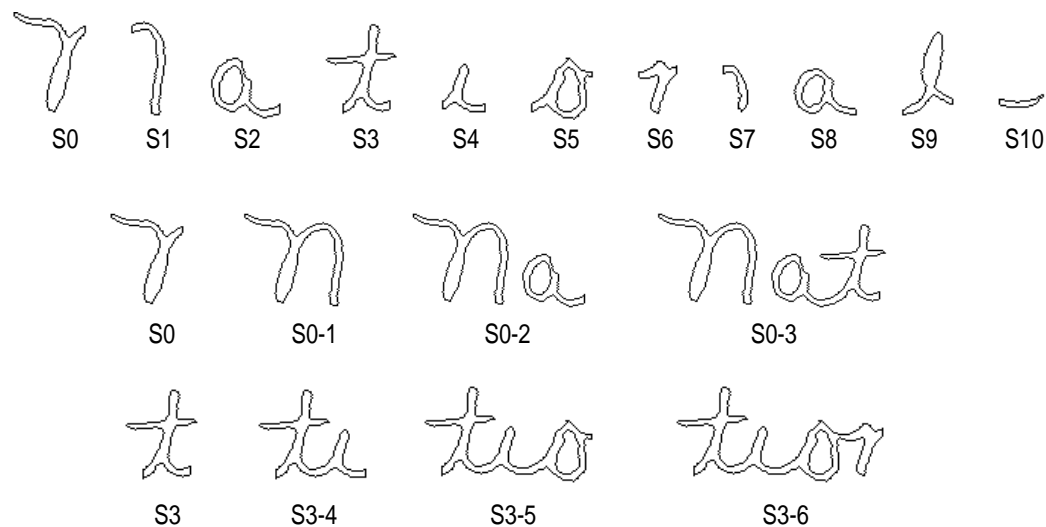


Figure 5.5: Character matching

In order to generate prototype templates, feature vectors of the hierarchical feature tree are extracted using a training character set and the same feature vector extraction method. And the extracted feature vectors which belong to a class are clustered independently in each sub-node, using a *K-means* clustering algorithm. A clustering procedure is completed if the mean square error of clusters to the centers is less than a predetermined threshold or the number of centers reaches the maximum allowable number of clusters.

The training character images are generated from pre-isolated character images which are manually segmented and collected from word images taken from address blocks of mail pieces in USPS main stream.

5.4.3 Character Matching

The major goal of character matching following stroke segmentation is to provide evidence in a given metric in order to find the best matches of given lexicon entries in a next word matching stage. All possible combinations of strokes within a permissible window are sent to a character recognizer. And the character recognizer provides corresponding confidence projected in its knowledge base. In our implementation, distance between feature vector extracted from the given combination of strokes and the built in prototypes of associated character subset, is used instead of confidence.

Following the stroke spacing model described in Chapter 4, it is assumed that a character has at least one prime strokes. The number of characters in a lexicon entry and the number of prime primitive in a given image are used to decide an actual size of matching window of each character. For more meaningful alignment of a sub-array of prime primitives to a character, filtering using a weighted matching window which is a function of a prime stroke period is applied to the combination of segmented strokes for a character.

A basic permissible window size for a i_{th} character h_i of a lexicon h on a stroke primitive array s is determined as

$$s_{min} = \min(i + 1, N_s - \eta_s \cdot (N_h - i)) \quad (5.46)$$

$$s_{max} = \max(\eta_s \cdot (N_h - i), N_s - (N_h - i)) \quad (5.47)$$

where s_{min} and s_{max} are the minimum and maximum of a window, N_h is the length of the lexicon h , and η_s is the maximum allowable number of primitives for a character matching.

A weight function for a matchable combination is defined as a step function

referenced by a prime stroke period as,

$$\tau(a, b) = \begin{cases} 1 & \frac{(z_{x_b} - z_{x_a})}{\bar{p} + \bar{v}_p} \leq \eta_w \\ \infty & \text{otherwise.} \end{cases} \quad (5.48)$$

where η_w is a threshold and a and b are indices of bound stroke segments associated.

Since the segmentation process is not the part of the recursion loop, once the effective combinations of stroke segments is determined, they remain the same until the end of recognition process. However, since the number of lexicon entries in lexicon $H(t)$ are dependent on the recursion time, the eligible subset of characters to its associate window is changed by the recursion depth.

A hierarchical character recognizer in a distance metric which described in the previous section is adopted for the character matching engine of this implementation. The classification accuracy of the hierarchical character recognizer is controlled by the number of nodes that can be used and the tree depth that can be accessible at recursion time t . The actual distance measurement of a subset of symbol h in a lexicon entry h for a combination of all strokes between s_a and s_b is given by a product of the window function and the recursive distance function of a hierarchical character recognizer as

$$A^{a,b}(h, t) = \tau(a, b) \cdot \min_{\forall c_i \in h} \check{\nu}_i(t) \quad (5.49)$$

The liability exposure of Equation 5.44 is obtained by direct replacement of $\delta_h(c_i | s_a \sim b)$ into $\check{\nu}_i(t)$ generated from image tree of $s_a \sim s_b$

5.4.4 Word Matching

Dynamic programming technique is used to find the best fit between the possible symbol strings of a lexicon entry and the segmented stroke frame of a given image. The matching distance found by a character recognizer is used as a metric to find word matching distance. One of the matching paths which generates minimum distance measure throughout a possible symbol string of a lexicon entry is chosen as the best match of the lexicon entry, as

$$A(\mathbf{h}, t) = [A^{a_0, b_0}(h_0, t), A^{a_1, b_1}(h_1, t), \dots, A^{a_{N_h-1}, b_{N_h-1}}(h_{N_h-1}, t)] \quad (5.50)$$

where $a_0 = 0$, $b_{N_h-1} = N_s - 1$ and $a_{i+1} - b_i = 1$.

The individual matching distance of character recognition are combined to compute the accumulated distance over the word string. The accumulated distance is normalized to the length of a lexicon entry and generated the overall distance of the lexicon word.

$$|A(\mathbf{h}, t)| = \frac{1}{N_h} \sum_{i=0, (a_{i+1}-b_i)=1}^{N_h-1} A_{h_i}^{a_i, b_i}(t) \quad (5.51)$$

Characters in a handwritten word often have flourish and ligature strokes which generally do not appear in independently written characters. These extra strokes are hard to be generalize in common prototypes, since they are usually unique to the writing style. The training character set can not cover all the variations of these strokes even if pre-isolated characters from word images are used for training because of the uniqueness of these extra strokes.

The character matching performance can be lowered because of these extra strokes. These strokes can be hidden behind of series of unaligned matches near the extra strokes. These mismatches cause degradation of global matching score

and may result in final decision of *rejection* in the following decision step even if the truth is ranked as the top choice.

The assignment of null to non-prime primitive strokes upon the global measurements of all the symbols can minimize the loss caused by extra strokes mentioned above. But in practical system, this *relaxed matching process* needs large computational effort due to the large number of combinations. To reduce this large computational effort, a matching method using cascaded classifiers is used. Figure 5.6 shows an example of relaxed matching.

A pre-classifier (called *stroke classifier*) is located before the major character matching and limits possible null assignment. The pre-classifier can only classify possible null assign primitive strokes as flourish and ligature stokes (the stroke segments of 1,4,7,9 and 11 in Figure 5.6). These null-allowable primitive strokes can be assigned to a null stroke before the first, after the last of a symbol string and between two adjacent symbols. In Figure 5.6, the flourish strokes (segments 1,11) and ligature strokes (segments 7,9) are matched to a null assignment. By such a relaxed matching, we can maximize the matching confidences of a lexicon entry.

Let us assume that a prime stroke primitive array \dot{s} of length $N_{\dot{s}}$ where $N_{\dot{s}} \leq N_s$ is generated by a stroke classifier from a stroke primitive array of s ,

$$\dot{s} = [\dot{s}_0, \dot{s}_1, \dots, \dot{s}_{N_{\dot{s}}-1}] \quad (5.52)$$

A basic permissible window size for a i_{th} character h_i of a lexicon h is determined as

$$\dot{s}_{min} = \min(i + 1, N_{\dot{s}} - \eta_s \cdot (N_h - i)) \quad (5.53)$$

$$\dot{s}_{max} = \max(\eta_s \cdot (N_h - i), N_{\dot{s}} - (N_h - i)) \quad (5.54)$$

where \dot{s}_{min} and \dot{s}_{max} are the minimum and maximum of a window, N_h is the length of the lexicon h , and η_s is the maximum allowable number of prime primitives

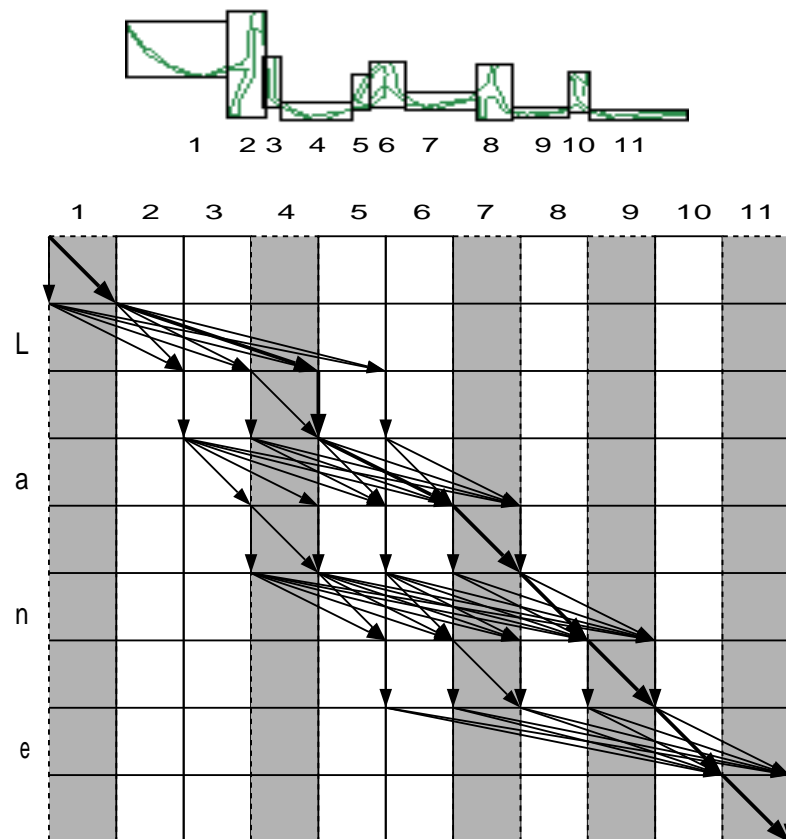


Figure 5.6: Word matching based on dynamic programming with nullity

for a character matching. Any continuous sub-array including non-prime stroke primitives within a window boundary can be a matchable combination.

Equation (5.50) can be replaced by a prime stroke sequence \dot{a}_i, \dot{b}_i instead of a primitive sequence a_i, b_i

$$A(\mathbf{h}, t) = [A^{\dot{a}_0, \dot{b}_0}(h_0, t), A^{\dot{a}_1, \dot{b}_1}(h_1, t), \dots, A^{\dot{a}_{N_h-1}, \dot{b}_{N_h-1}}(h_{N_h-1}, t)] \quad (5.55)$$

where $\dot{a}_0 = 0$, $\dot{b}_{N_h-1} = N_s - 1$ and $\dot{a}_{i+1} - \dot{b}_i = 1$. And the overall distance is

$$|A(\mathbf{h}, t)| = \frac{1}{N_h} \sum_{i=0, (\dot{a}_{i+1} - \dot{b}_i)=1}^{N_h-1} A^{\dot{a}_i, \dot{b}_i}(t) \quad (5.56)$$

As a pre-classifier, the primitive classification methods which are used to generated prime primitive sequence in Section 4.4.1 are directly applied to limit null assignment.

5.5 Experiment

The proposed recursive word recognition scenario is implemented in a lexicon driven word recognizer for experiments. For a character recognition engine, a hierarchical character recognizer which has a quad tree structure of depth of 4 has been implemented. Thus, 85 (1+4+16+64) sub-nodes have been built. From setting of $N_f = 2$ and $N_g = 8$, total 34 feature measurements are extracted for a feature vector. Using manually isolated character images from word images which are exclusive to testing word images, a maximum of 30 templates per class have been generated for each sub-node by a $K - means$ clustering algorithm with a mean square error threshold of 0.1. A stroke classifier has been used before the word matching step to generate a prime stroke sequence. And Dynamic Programming

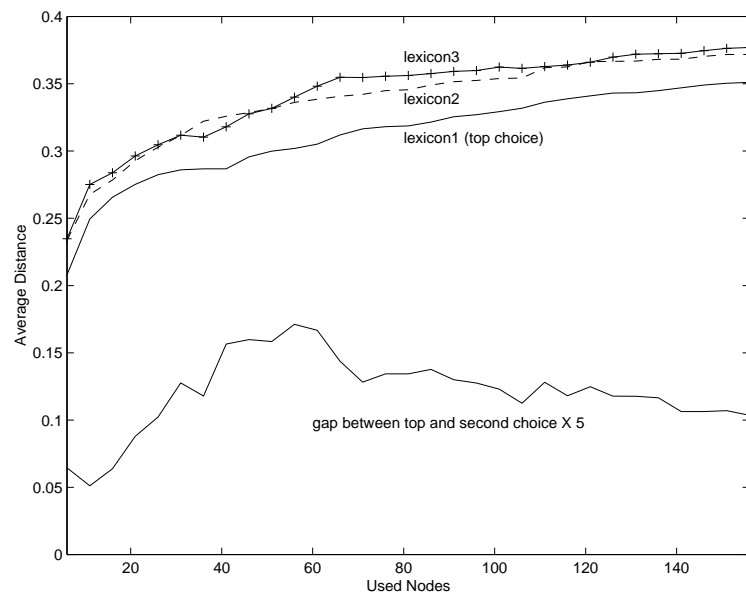
technique is used to find the best matching path on the prime stroke array with nullity option of non prime strokes.

3000 word images (digitized at 212 dot-per-inch) including city names, state names, and street names collected from address block of mail pieces in USPS mail stream were used as testing images. The lexicons were generated as subsets which has randomly selected lexicon entries in fixed size (10, 100 or 1000) from the pool of all the possible city, state, and street names. The true word was always present in the lexicon. The word recognizer have achieved maximum 96.3% of top choice correct rate with size 10 lexicons.

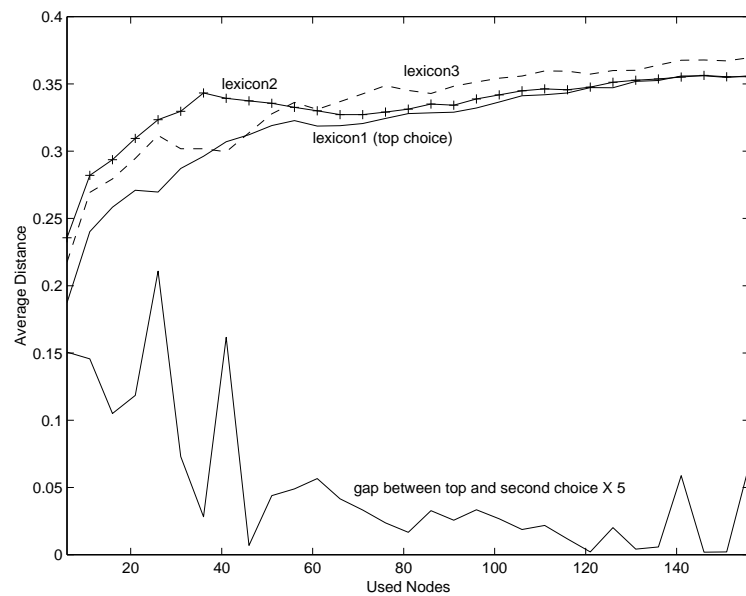
Figure 5.7 shows dynamic characteristics of the implemented recursive word recognizer without considering the complexity of the lexicon. As more resources are used, the average minimum distance grows and the separability between distance of the lexicon entries become also larger. If the usage of sub-nodes increases continuously the gap between distance of the lexicon entries stays in a steady state or starts to decrease at a certain point (because of dimensionality problem). However as shown in Figure 5.7-(b), the behavior of distance gaps between lexicon entries is unstable, if the recognizer fails to rank the truth as a top choice. The processing time is linearly proportional to the number of used sub-images.

Figure 5.8 shows dynamic characteristics of adjusted distance by the liability exposure in term of lexicon complexity. The index of rejectivity, β , is 0.7, the maximum number of usable nodes is fixed at 96, and the lexicon size is 100. If the rejected lexicon entries become larger, the distance gap between the first and second choices becomes large.

Figure 5.9 shows the overall effect of using lexicon complexity throughout the testing set, with fixed usage of knowledge built in. The usable nodes were fixed as the root node and 4 child sub-images of the root without an active selection of

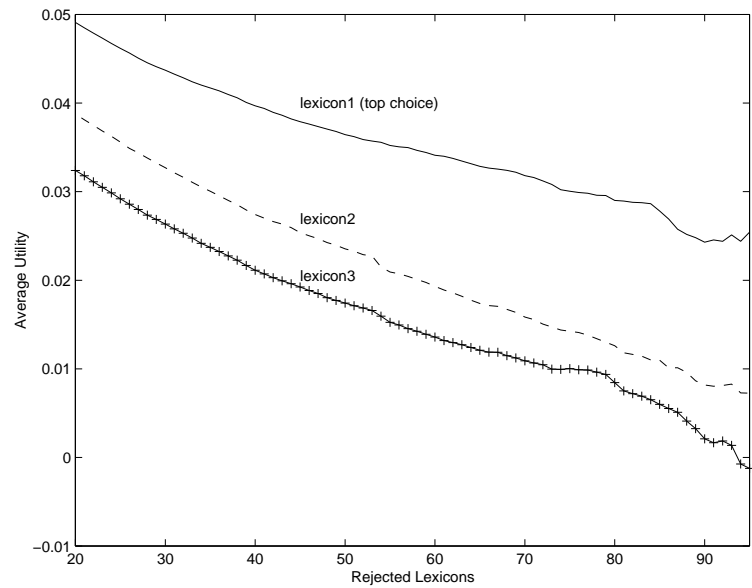


(a)

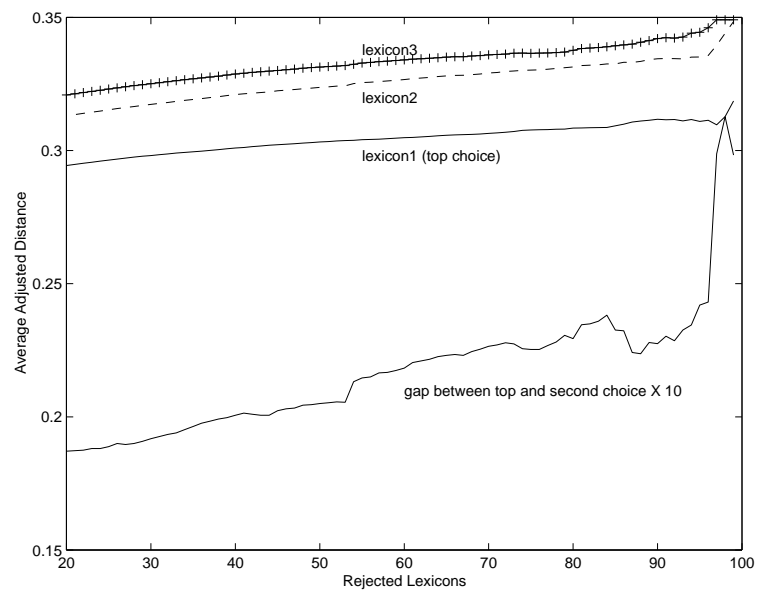


(b)

Figure 5.7: Accuracy support changes: (a) success case and (b) failure case



(a)



(b)

Figure 5.8: Adjusted distance changes: (a) utility values and (b) adjusted distance

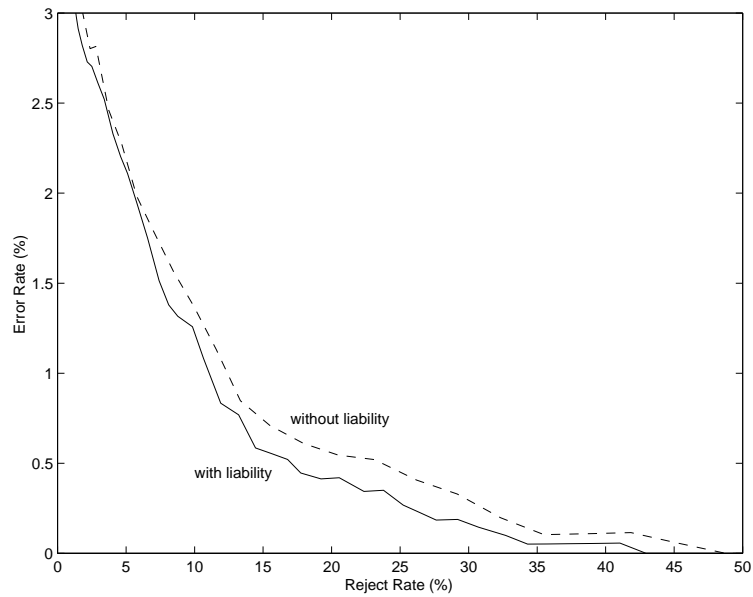


Figure 5.9: Recognition performance with and without liability

them. The lexicons of size 10 were used. The index of rejection β was set at 0.7 when liability (lexicon complexity) was applied. The recognizer using both character recognition distance and lexicon complexity (liability) has better performance compared to the recognizer using only the distance.

Figure 5.10 shows word recognition performance. The top choice correct rate is changed by adjusting the usage of nodes. The performance with feedback loop of active recursion controlled by decision maker which reacts to the image quality is superior than that of non-recursive with fixed usage of resources. With same degree of usage nodes, the recursive system achieves better top choice correct rate, even though the both steady state recognition performance are same.

Figure 5.11 shows the dynamic character recognition characteristics when the

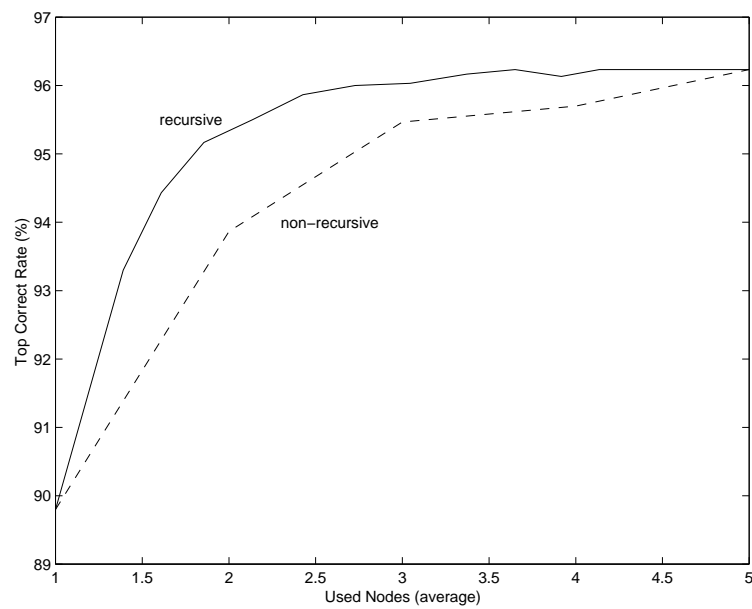


Figure 5.10: Dynamic recognition characteristics - comparison between recursive system with active searching and non-recursive system with passive adding of sub-nodes

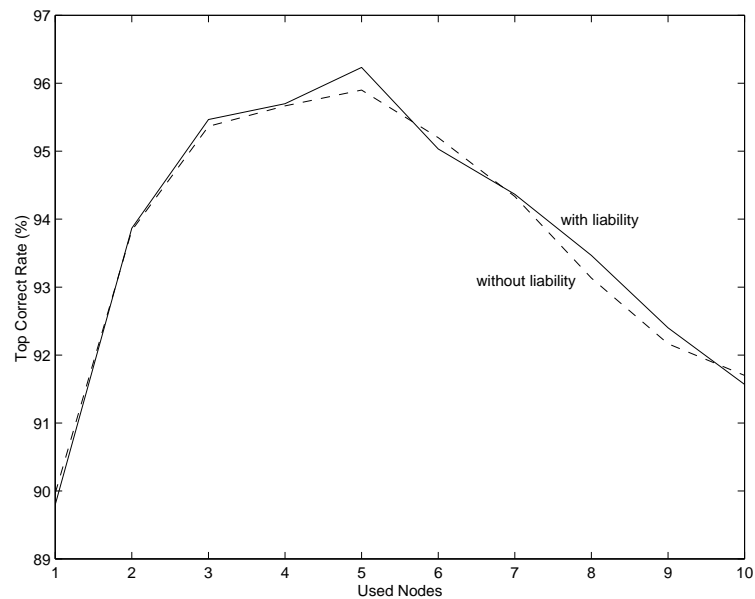


Figure 5.11: Dynamic recognition characteristics - recognition performance changes by various acceptance criterion

acceptance degree of decision making is changed. The x-axis displays the average usage of nodes. If the accepting terminal condition is more restrictive, the usage of features is increased because more discrimination power between classes is necessary. When liability of character confidence within lexicon set is used, the performance can be maximized.

Chapter 6

Dynamic Phrase Recognition

Comparison between feature vectors of several possible combinations of stroke segments and prototypes of characters of a lexicon entry are made to find the best match of the word in analytic word recognition which uses stroke segmentation before character recognition. Similarly, in dynamic phrase recognition which uses word segmentation before word recognition, several possible combinations of word image segments are submitted to a word recognizer to find the best match between an array of word image segments and a phrase lexicon entries.

If the number of word segments in a given phrase image and the lexicon size are small, an exhaustive testing of all the combinations of the word segments and the lexicon entries can maximize the success rate. But it is not acceptable in practice due to the limitation of resources and the burden of processing time.

A dynamic phrase recognition method shown in Figure 6.1 is proposed to achieve high recognition performance while keeping a system cost-efficient. A prime stroke analysis described in Chapter 4 is used to detect the dominant stroke

frequency. This analysis is also extended to word gap classification which provides numeric confidence, and is applied to a method which provides the estimated number of characters within a segmented word image. The word gap confidence and the estimated number of characters of a word segment are used as basic parameters to generate a subset of lexemes. Phrase hypotheses are generated by an array of combinations of one or grouped word image segments and eligible lexemes and tested by a lexicon driven word recognizer.

6.1 Prescreening

By assumptions 3 and 4 of the spacing model in Section 4.1, the number of characters within an image segment can be estimated using prime spatial period, which is the period of estimated prime frequency. By counting the number of exclusive prime primitives pairs whose horizontal distance between them is larger than the prime spatial period, the number of character contents in a word image segment can be obtained roughly.

Let K be the estimated number of characters in a word image segment obtained by the method described above. Because of flourish strokes, different character size, and spacing style, the accuracy of K is not high enough to use as a stand-alone parameter. So, in order to compensate of the inaccuracy, an error range, e is defined to support accuracy of the estimated number of characters as,

$$e = \frac{K \cdot \tilde{v}_p}{\tilde{p}} \quad (6.1)$$

where \tilde{p} is an extracted prime spatial period and \tilde{v}_p is the variance of the prime spatial period which is defined in Section 4.4.2.

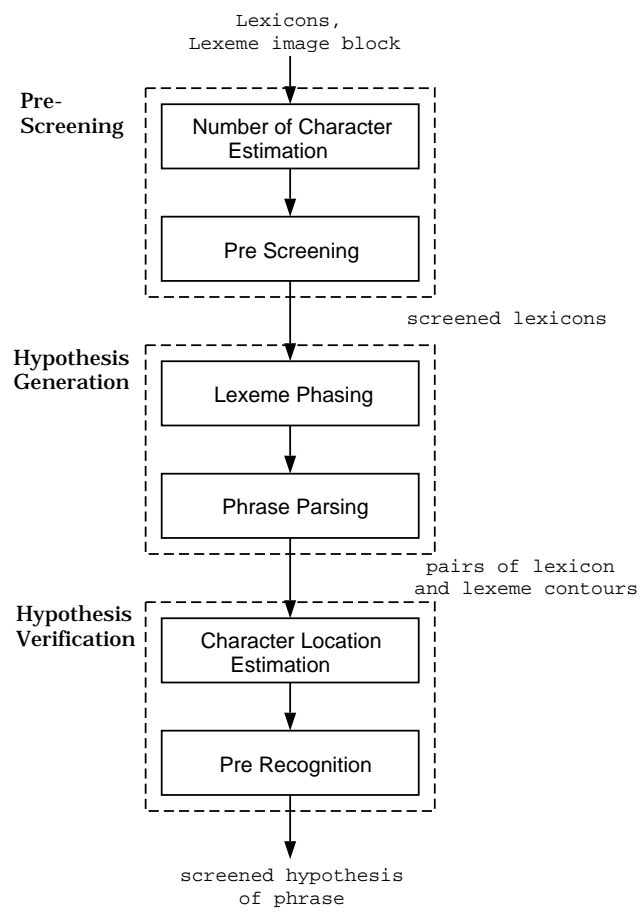


Figure 6.1: Dynamic matching phrase recognition

A boundary of estimated number of characters is defined as

$$K_{max} = \begin{cases} K + \mu_{max} \cdot e & \text{if } (K + \mu_{max} \cdot e) < N_s \\ N_s & \text{otherwise.} \end{cases} \quad (6.2)$$

and

$$K_{min} = \begin{cases} K - \mu_{min} \cdot e & \text{if } (K - \mu_{min} \cdot e) > 1 \\ 1 & \text{otherwise.} \end{cases} \quad (6.3)$$

where N_s is the total number of prime primitives in a word image segment. μ_{max} and μ_{min} are control parameters.

Using the estimated number of characters and the number of word segments, the size of a given lexicon can be reduced after rejecting lexicon entries which does not meet the below conditions.

1. The number of words in a lexicon is smaller than the number of word segments.
2. Total number of characters in a lexicon is equal or less than the sum of K_{max} .
3. Total number of characters in a lexicon is equal or greater than the sum of K_{min} .
4. If the number of words is same as the number of image segments, the number characters of each word is within bound of K_{max} and K_{min} in the corresponding image segment.

In an phrase recognition application, if a lexicon contains only subset of the true string, the third condition mentioned above can not be applied as a condition of selecting a lexicon subset, because provided lexicon entries should be allowed

to be matched with the part of the phrase images. If an input phrase image has additional string components, it is the same situation as the above case.

6.2 Hypotheses Generation

The over-segmentation of a phrase image can lower the error rate of missing word break points. However, because of false detections of word breaks which are not actual word breaks, word recognition must be extended to consider grouped image segments as a possible single word image. A strategy using over-segmentation followed by a dynamic matching process between one or a group of word image segments and a single lexeme (called *lexeme matching*) provides better recognition performance but it increases the burden of word recognition.

Dynamic matching process between an image segment to multiple lexeme should be considered in order to recover the error cases of under-segmentation. However, because the chance of under-segmentation is very small as shown in Table 4.1, a matching type of an image segment to multiple lexeme is not considered in this work.

In order to reduce the burden of lexeme matching, the globally unsuitable lexicon entries are discarded from a lexicon by the prescreening process described in the previous section. However, the number of hypotheses is still going to be large, since estimating the number of characters in a word segment is difficult. Because the boundary of estimation is large, the number of hypotheses is going to be very large.

Meanwhile, checking possible combinations of image segments to each lexeme can reduce the burden of word recognition. If the gap for a detected word break is significantly larger than the prime stroke period, it reasonably believed that the

gap is a valid word break. If anchors are selected among candidate of word breaks, the matching complexity can be lower.

Word break candidates are categorized into two groups; *hard word point* and *soft word point* and grouping of image segments crossing a hard word point is not allowed as a permissible segment combination. The gap confidence which is defined in Section 4.4.2 is used for the classification of word break types. In Figure 6.2, thick lined boxes show hard word blocks and the thin lined boxes show soft word blocks.

To explain the generation of hypotheses, let us assume that our application is street name recognition and assume that inputs are the image and the raw lexicon shown in Figure 6.2 and Table 6.2 respectively.

As shown in Figure 4.1, there are several variations of writing the prefix and suffix of a street name. For example, the full word, ‘Drive’, instead of a abbreviation, ‘Dr’, can appear in a street name image. Moreover, additional writing objects associated with the street name can be found such as ‘2nd Fl’, ‘Apt’, ‘Suite’, etc. Certainly, the lexicon entries do not anticipate their presence.

To overcome these kinds of variations, and missing or incorrect prefix and suffix, each lexicon entry is expanded in all conceivable ways as Table 6.3. We assume all additional writing objects can only be attached at the end of a street name.

All possible combinations of image segments and lexemes are generated as lexeme hypotheses. Type of word break, boundary of estimated number of characters and order of a word in a lexicon entry are used as restriction. In the example, the lexeme ‘Steamboat’ can not match image *segment 0*, because the number of characters of ‘Steamboat’ is out of the estimated boundary of number of characters in *segment 0*. Also since a hard word point is laid between image *segment 0* and *segment 1*, both segments can not be grouped to match ‘Steamboat’. So, the

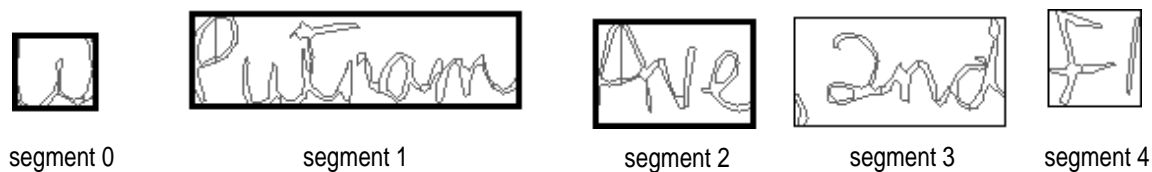


Figure 6.2: Word image segments

image segment index	0	1	2	3	4
no. of primitives m	4	12	5	6	3
esti. no. of char. n	3	9	5	5	2
error bound r	1	3	2	2	1
n_{max}	4	12	5	6	3
n_{min}	1	2	1	1	1

Table 6.1: Estimated number of characters

phrase lexicon entries which contain ‘Steamboat’ at first lexeme are rejected. Table 6.1 shows the boundary of estimated results and Table 6.4 lists all the generated lexeme hypotheses.

After generating lexeme hypotheses, phrase hypotheses are generated as an array of lexeme hypothesis indices. A phrase can have multiple phrase hypotheses sharing common lexeme hypotheses. (Table 6.5).

Pre	Name	Suffix
	LAKE	AVE
	STEAMBOAT	DR
W	PUTNAM	AVE

Table 6.2: Raw lexicon

Pre	Name	Suffix
	Lake	Ave
	Lake	Avenue
	Lake	
	Steamboat	Dr
	Steamboat	Drive
	Steamboat	
W	Putnam	Ave
W	Putnam	Avenue
W	Putnam	
West	Putnam	Ave
West	Putnam	Avenue
West	Putnam	
	Putnam	Ave
	Putnam	Avenue
	Putnam	

Table 6.3: Expanded lexicon

index	lexeme	start	dur
0	Lake	0	1
1	Ave	1	1
2	#	2	2
3	Avenue	1	1
4	#	1	4
5	W	0	1
6	Putnam	1	1
7	Ave	2	3
8	Ave	2	1
9	#	3	2
10	Ave	2	2
11	#	4	1
12	Avenue	2	3
13	Avenue	2	2
14	West	0	1

Table 6.4: Lexeme hypotheses - the symbol # indicates remaining image segments corresponding to additional informations

Pre	Name	Suffix	Add
	Lake (0)	Ave (1)	# (2)
	Lake (0)	Avenue (3)	# (2)
	Lake (0)		# (4)
W (5)	Putnam (6)	Ave (7)	
W (5)	Putnam (6)	Ave (8)	# (9)
W (5)	Putnam (6)	Ave (10)	# (11)
W (5)	Putnam (6)	Avenue (12)	
W (5)	Putnam (6)	Avenue (13)	# (11)
W (5)	Putnam (6)		# (2)
West (14)	Putnam (6)	Ave (7)	
West (14)	Putnam (6)	Ave (8)	# (9)
West (14)	Putnam (6)	Ave (10)	# (11)
West (14)	Putnam (6)	Avenue (12)	
West (14)	Putnam (6)	Avenue (13)	# (11)
West (14)	Putnam (6)		# (2)

Table 6.5: Phrase hypotheses - the number within parenthesis indicates index of lexeme hypothesis

6.3 Hypotheses Verification

A lexicon driven word recognizer is preferred for hypothesis verification rather than a character based word recognizer, because the hypotheses consist of word segments and a subset of possible lexemes. Dynamic programming based matching scheme through an array of stroke segments and the character strings of lexicon entries is a well known technique to find the best match between a word image and a lexicon.

Since the normalized distance passing through a matching path or average of individual matching scores of characters is usually used as metric to find the best matching path, a mismatch of a character in a lexicon entry, which is usually caused by an error of stroke segmentation, is compensated by the other good matches, if the lexicon entries are the true.

The stroke segmentation method used in our experiment, has been implemented in a heuristic rule based algorithm. Since this segmentation method is not perfect to separate all the touching characters, a longer string has a better chance to be in successful matches. Thus, recognition using an entire array of image segments consisting a phrase has better performance than individual lexeme based recognition in hypothesis verification. In our implementation, grouped image segments for a phrase hypothesis is submitted to a word recognizer with a eligible subset of phrase lexicon strings. The word recognition scores are retained as the hypotheses confidence.

Total Images	1342
Images with additional unwanted segments	116 (8.64%)
Raw Lexicons (average)	22.0
Expanded Lexicons (average)	77.2

Table 6.6: Information of experiment image set

6.4 Experiment

The proposed phrase recognition method has been applied to street name recognition. In this experiment, the goal is to select the best match between a street name line image extracted from a USPS mail piece and provided lexicon entries in full or partial matching. The recognition performance of the proposed phrase recognition method is compared to that of word mode recognition. The word mode recognition is a simplistic approach where a street name image is treated as a single word by ignoring inter-word spacing.

The same street name images which had been used in the experiment of Chapter 4 were used. Each raw lexicon has been obtained from postal databases using live queries of a ZIP code and a primary number which were shown on the same mail piece from which the street image had been extracted. The ZIP codes and primary numbers were manually truthed.

Since the testing set came from live mail pieces, the size of lexicon of each images was not fixed. The number of raw lexicon entries was ranged from 1 to 100. Among our test set, about 9% of the images have additional unwanted writing objects such as apartment numbers (Table 6.6) in the input image. Table 6.7 shows the

Number of Image	1342
Raw Lexicons	22.0
Expanded Lexicons	77.2
Accepted Lexicons	62.0
Lexeme Hypotheses	59.2
Phrase Hypotheses	98.5

Table 6.7: Experiment environment - all numbers are average values

environment before hypotheses are submitted to a word recognizer. Each value is averaged to the number of images.

Figure 6.3 shows top choice correct rate of phrase mode and word mode recognition method by adjusting a threshold of confidence. Figure 6.4 shows the reject-error curves. The phrase mode recognition has better performance.

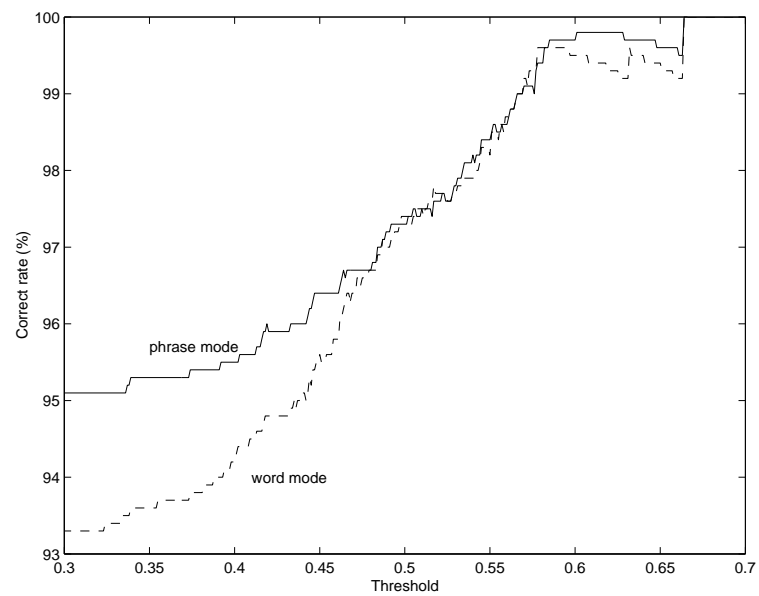


Figure 6.3: Performance of phrase recognition algorithm; Top choice correct rate

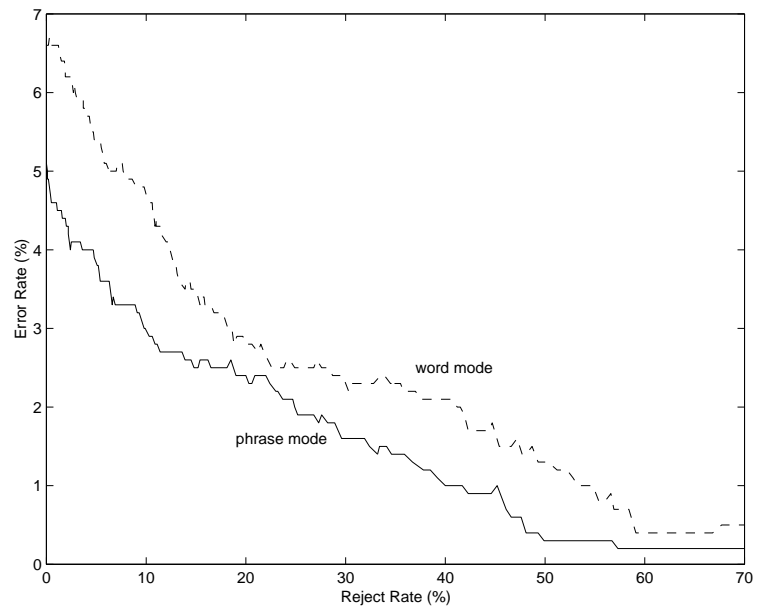


Figure 6.4: Reject to error characteristic of phrase recognition

Chapter 7

Conclusions

7.1 Summary

Following the motivation of human cognitive models, two new concepts of active and recursive recognition are introduced for handwritten text recognition. A cognitive word perception model of time course and *satisficing decision* which adopts *epistemic utility theory* is applied to the off-line handwritten word recognition problem in this dissertation.

A computational model and recursive system architecture are proposed to extend a conventional recognizer into an interactive and dynamic one. A hierarchical character recognition scenario and its application to recursive word recognition are described. Also, using a stroke period obtained from the proposed prime stroke analysis, a computationally efficient word recognition algorithm and a dynamic phrase recognition method are designed.

A hierarchical classification method has dynamic usage of hierarchical feature

space that preserves the benefits of the multi-resolution model. Recognition action begins with globally extracted features and focuses on smaller sub-images of the character successively on each pass. The sub-images are identified by a active searching algorithm called *gaze planning* which guides the successive passes through a series of informative sub-images. This hierarchical classification algorithm allows a high recognition rate with a small dimensional set of features. Finer resolution of features is simulated by sampling smaller sub-images within the same feature dimension in each pass.

A lexicon driven word recognizer is designed by applying the proposed hierarchical character recognizer which is modified to be controllable in recognition resolution externally. The word recognition engine operates actively and has different degrees of classification ability which adapts to input image quality. Both accuracy support and liability exposure of a recognition result are formulated using the character recognition confidence and matching transform distance. A lexicon complexity derived from the matching transform distance is utilized as a decision metric, which measures the difficulty of a given lexicon set projected to the classification ability of a character recognizer. Recognition, decision making and recursive updating from a closed loop architecture is designed for a recursive recognition scenario. This proposed model recursively enhances the degree of classification accuracy until a decision condition is satisfied.

A prime stroke analysis which detects a writer's spacing style is formulated to extend word recognition to a phrase recognition system. An analytic handwriting spacing model is generated in terms of strokes which are located in the primary writing band excluding flourishes and ligatures. The analysis captures writing styles of spacing between characters and words using statistical formulations. In a given image, the stroke period which is extracted from this stroke analysis is used

as a metric to detect word break points and to limit combinations of segmented strokes for word recognition.

An efficient phrase recognition system which finds the best matching path between a phrase image and a phrase lexicon which allows partial matching is presented. A method for estimating the bound of the number of characters without actual recognition and a computational model for word gap classification are also presented. The estimated number of characters and word gap confidence are used as basic parameters to generate a filtered subset of lexicon entries before recognition.

The proposed hierarchical classification model achieves 98% and 96% of top choice correct rates in NIST digit and character sets. Also it proves to be 300 times faster than a traditional k -nearest neighbor classification method and 10 times faster than a neural network method. The proposed word recognition model achieve 96% of top choice correct rate in a word recognition application tested on the city name word images collected from USPS mail stream respectively. Also 10% improvement in recognition performance has been obtained for phrase recognition application of street name lines of USPS mail pieces comparing to the performance of a word mode recognition.

7.2 Contribution

It is believed that this dissertation has made the following contributions in the field of handwriting recognition:

- A new character classification approach which operates actively in hierarchical feature space and which preserves multi-resolution benefits.

- A new word recognition architecture which combines an active recognition engine and a cognitive decision maker recursively.
- A prime stroke analysis which detects writer's spacing style.
- A dynamic phrase recognition paradigm which overcomes the problem of incomplete lexicon entries and existence of unwanted objects.

Bibliography

- [1] E. Cohen, J. J. Hull and S. N. Srihari. Control Structure for Interpreting Handwritten Address. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(10):1049–1055, Oct 1994.
- [2] E. Cohen, J. J. Hull and S. N. Srihari. Understanding Handwritten Text in a Structured Environment: Determining ZIP codes from Address. *International Journal of Pattern Recognition and Artificial Intelligence*, 5:221–264, 1991.
- [3] J. C. Simon O. Barat and N. Gorski. A System for the Recognition of Handwritten Literal Amounts of Checks. *Int. Conf. on Document Analysis System, Kaiserslautern, Germany*, pages 135–155, 1994.
- [4] G. Kim and V. Govindaraju. Bankcheck Recognition Using Cross Validation between Legal and Courtesy Amount. *Automatic Bankcheck Processing, World Scientific*, pages 195–212, 1997.
- [5] S. N. Srihari, Y.-C. Shin, V. Ramanaprasad and D. S. Lee. Name and Address Block Reader. *Proceedings of IEEE*, 84(7):1038–1049, July, 1996.
- [6] S. N. Srihari. From Pixels to Paragraph: the Use of Contextual Models in Text Recognition. In *International Conference on Document Analysis and Recognition, Japan*, pages 416–423, 1993.

- [7] V. Govindaraju, A. Shekhawat and S. N. Srihari. Interpretation of Handwritten Addresses in US Mail Stream. In *The Third International Workshop on Frontiers in Handwriting Recognition, Buffalo, New York*, pages 197–206, 1994.
- [8] Y. Belaid, A. Belaid and E. Turolla. Item Searching in Forms: Application to French Tax Form. In *International Conference on Document Analysis and Recognition, Montreal, Canada*, pages 744–747, 1995.
- [9] R. O. Duda and P. E. Hart. *Pattern Classification and Scene Analysis*. John Wiley & Sons, Inc, New York, 1979.
- [10] R. Schalkoff. *Pattern Recognition: Statistical, Structural and Neural Approaches*. John Wiley & Sons, Inc, New York, 1992.
- [11] S. Mori, C. Y. Suen and K. Yamamoto. Historical Review of OCR Research and Development. *Proceedings of the IEEE*, 80:1029–1058, 1992.
- [12] G. Nagy. At the Frontiers of OCR. *Proceedings of the IEEE*, 80:1093–1100, 1992.
- [13] C. Y. Suen. Computer Recognition of Unconstrained Handwritten Numerals. *Proceedings of the IEEE*, 80:1162–1180, 1992.
- [14] V. K. Govindan and A. P. Shivaprasad. Character Recognition - A Review. *Pattern Recognition*, 23(7):671–683, 1990.
- [15] J. Mantas. An Overview of Character Recognition Methodologies. *Pattern Recognition*, 19(6):425–430, 1986.
- [16] C. Y. Suen, M. Berthod and S. Mori. Automatic Recognition of Handprinted Characters - The state of the Art. *Proceedings of the IEEE*, 68:469–487, 1980.

- [17] O. D. Trier, A. K. Jain and R. Taxt. Feature Extraction Methods for Character Recognition - A survey. *Pattern Recognition*, 29(4):641–662, 1996.
- [18] G. Srikantan S. W. Lam and S. N. Srihari. Gradient-Based Contour Encoding for Character Recognition. *Pattern Recognition*, 29(7):1147–1160, 1996.
- [19] M. H. Glauberman. Character Recognition for Business Machines. *Electronics*, pages 132–136, February 1956.
- [20] C.-S. Lin and C.-L. Hwang. New Forms of Shape Invariants from Elliptic Fourier Descriptors. *Pattern Recognition*, 20(5):535–545, 1987.
- [21] F. P. Kuhl and C. R. Giardina. Elliptic Fourier Features of a Closed Contour. *Computer Graphics and Image Processing*, 18:236–258, 1982.
- [22] C. T. Zahn and R. C. Roskies. Fourier Descriptors for plane closed curves. *IEEE Transactions on Computers*, 21(3):269–281, 1972.
- [23] G. H. Granlund. Fourier Preprocessing for Hand Printed Character Recognition. *IEEE Transactions on Computers*, 21(2):195–201, 1992.
- [24] R. R. Bailey and M. Srinath. Orthogonal Moment Features for Use With Parametric and Non-Parametric Classifiers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(4):369–398, 1996.
- [25] C.-H. Teh and R. T. Chin. On Image Analysis by the Methods of Moments. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 10(4):496–513, 1988.

- [26] Y. S. Abu-Mostafa and D. Psaltis. Recognition Aspects of Moment Invariants. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 6(6):698–706, 1984.
- [27] S. Maitra. Moment Invariants. *Proceedings of the IEEE*, 67(4):697–699, 1979.
- [28] M.-K. Hu. Visual Pattern Recognition by Moment Invariants. *IRE Transactions on Information Theory*, IT-8:179–187, February 1962.
- [29] Y. Mallet, D. Coomans, J. Kautsky and O. DeVel. Classification Using Adaptive Wavelets for Feature Extraction. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(10):1058–1066, 1997.
- [30] P. Wunsch and A. F. Laine. Wavelet Descriptors for Multiresolution Recognition of Handprinted Characters. *Pattern Recognition*, 28(8):1237 – 1249, 1995.
- [31] A. Jain and D. Zongker. Feature Selection: Evaluation, Application and Small Sample Performance. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(2):153–197, 1997.
- [32] P. D. Gader and M. A. Khabou. Automatic Feature Generation for Handwritten Digit Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(12):1256–1261, 1996.
- [33] K. Kira and L. A. Rendell. Feature Selection Problem: Traditional Methods and a New Algorithm. In *Proceedings of Tenth National conference on Artificial Intelligence, San Jose, California*, pages 129–134, 1992.
- [34] J. Park, V. Govindaraju and S. N. Srihari. OCR in A Hierarchical Feature Space. In *IEEE International Conference on Systems, man, and Cybernetics, San Diego, CA*, pages 4324–4329, Oct. 1998.

- [35] M. Revow, C. K. Williams and G. E. Hinton. Using Generative Models for Handwritten Digit Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(6):592–606, 1996.
- [36] J. T. Favata and G. Srikantan. A Multiple Feature/Resolution Approach to Handprinted Digit and Character Recognition. *International Journal of Imaging Systems and Technology*, 7:304–311, 1996.
- [37] J. Rocha and T. Pavlidis. Character Recognition Without Segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(9):903–909, 1995.
- [38] J. Rocha and T. Pavlidis. A Shpae Analysis Model with Applications to a Character Recognition System. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(4):393–404, 1994.
- [39] J. T. Favata, G. Srikantan and S. N. Srihari. Handprinted Character/Digit Recognition using a Multiple Feature/Resolution Philosophy. In *Proceedings of the Fourth International Workshop on Frontiers in Handwriting Recognition, Taiwan*, pages 57–66, 1994.
- [40] S. J. Smith, M. O. Bourgoïn, K. Sims and H. L. Voorhees. A Shpae Analysis Model with Applications to a Character Recognition System. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(4):393–404, 1994.
- [41] J.-M. Bertille. An Elastic Matching Approach Applied to Digit Recognition. In *Proceedings of International Conference on Document Analysis and Recognition*, pages 82–85, 1993.
- [42] F. Kimura and M. Shridhar. Handwritten Numerical Recognition based on Multiple Algorithm. *Pattern Recognition*, 24(10):969–983, 1991.

- [43] L. Lam and C. Y. Suen. Structural Classification and Relaxation Matching of Totally Unconstrained Handwritten ZIP-Code Numbers. *Pattern Recognition*, 21(1):19–31, 1988.
- [44] T. M. Ha and H. Bunke. Off-Line, Handwritten Numeral Recognition by Perturbation Method. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(5):535–539, 1997.
- [45] Y. S. Huang and C. Y. Suen. A Method of Combining Multiple Experts for the Recognition of Unconstrained Handwritten Numerals. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(1):90–94, 1995.
- [46] D.-S. Lee and S. N. Srihari. A Theory of Classifier Combination: The Neural Network Approach. In *Proceedings of International Conference on Document Analysis and Recognition, Montreal, Canada*, pages 42–45, 1995.
- [47] T. K. Ho J. J. Hull and S. N. Srihari. Decision Combination in Multiple Classifier System. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(1):66–75, 1994.
- [48] Y. Lu and F. Yamaoka. Integration of Handwritten Digit Recognition Results Using Evidential Reasoning. In *Proceedings of the Fourth International Workshop on Frontiers in Handwriting Recognition, Taiwan*, pages 456–463, 1994.
- [49] S. C. Huang, K. Liu and C. Y. Suen. A Neural Network Approach for Multiclassifier Recognition Systems. In *Proceedings of the Fourth International Workshop on Frontiers in Handwriting Recognition, Taiwan*, pages 235–244, 1994.
- [50] F. Yamaoka, Y. Lu, A. Shaout and M. Shridhar. Fuzzy Integration of Classification Results in a Handwritten Digit Recognition System. In *Proceedings*

of the Fourth International Workshop on Frontiers in Handwriting Recognition, Taiwan, pages 255–264, 1994.

- [51] L. Xu, A. Crzyzak and C. Y. Suen. Methods of Combining Multiple Classifiers and Their Applications to Handwriting Recognition. *IEEE Transactions on System, Man and Cybernetics*, 22(3):418–435, 1992.
- [52] C. Y. Suen, C. Nadar, T. A. Mai, R. Legault and L. Lam. Recognition of Totally Unconstrained Handwritten Numerals based on the concept of multiple experts. In *Proceedings of the First International Workshop on Frontiers in Handwriting Recognition*, pages 131–143, 1990.
- [53] P. Whaite and F. P. Ferrie. Autonomous Exploration: Driven by Uncertainty. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(3):193–205, 1997.
- [54] P. Whaite and F. P. Ferrie. From Uncertainty to Visual Exploration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13(10):1038–1049, 1991.
- [55] D. J. C. Mackay. Information-Based Objective Functions for Active Data Selection. *Neural Computation*, 4:590–604, 1992.
- [56] R. Bajscy and M. Campos. Active and Exploratory Perception. *Image Understanding*, 56(1):31–40, 1992.
- [57] R. Bajscy. Active Perception. *Proceedings of the IEEE*, 76(8):996–1005, 1988.
- [58] A. W. Senior and A. J. Robinson. An Off-Line Cursive Handwriting Recognition System. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(3):309–321, 1998.

- [59] G. Kim and V. Govindaraju. A Lexicon Driven Approach to Handwritten Word Recognition for Real-Time Applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(4):366–379, 1997.
- [60] P. D. Garder, M. Mohamed and J.-H. Chiang. Handwritten Word Recognition with Character and Inter-Character Neural Networks. *IEEE Transactions on Systems, Man and Cybernetics - Part B*, 27(1):158–164, 1997.
- [61] M. Mohamed and P. Gader . Handwritten Word Recognition Using Segmentation-Free Hidden Markov Modeling and Segmentation-Based Dynamic Programming Techniques. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(5):548–554, 1997.
- [62] J. T. Favata. Character Model Word Recognition. In *International Workshop on Frontiers in Handwriting Recognition 96, U.K.*, pages 437–440, 1996.
- [63] R. M. Bozinovic and S. N. Srihari. Off-Line Cursive Script Word Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(1):68–83, 1989.
- [64] J. C. Simon. Off-Line Cursive Word Recognition. *Proceeding of the IEEE*, 80(7):1150–1161, 1992.
- [65] M. Y. Chen, A. Kundu and S. N. Srihari. Variable Duration Hidden Markov Model and Morphological Segmentation for Handwritten Word Recognition. *IEEE Transactions on Image Processing*, 4(12):1675–1688, 1995.
- [66] M. Y. Chen, A. Kundu and J. Zhou. Off-Line Handwritten Word Recognition Using a Hidden Markov Model Type Stochastic Network. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(5):481–496, 1994.

- [67] A. Gilles. Cursive Word Recognition Using Hidden Markov Models. In *Proceedings of United States Postal Service Advanced Technology Conference, Washington D.C.*, pages 557–563, 1992.
- [68] S. Madvanath and V. Govindaraju. Holistic Lexicon Reduction. In *Proc. Third Int. Workshop Frontiers in Handwriting Recognition, Buffalo, NY*, pages 132–141, 1993.
- [69] S. Madvanath. Using Holistic Features in Handwritten Word Recognition. In *Proc. U.S. Postal Service Advanced Technology Conference, Washington D.C.*, pages 183–199, 1992.
- [70] L. D. Harmon. Handwriting Reader Recognizers Whole Words. *Electronics*, 35:29–31, 1962.
- [71] R. Farag. Word Level Recognition of Cursive Script. *IEEE Transactions on Computers*, c-28(2):172–175, 1979.
- [72] S. Madvanath E. Kleinberg and V. Govindaraju. Empirical design of a multiple classifier thresholding/control strategy for recognition of handwritten street names. *International Journal of Pattern Recognition and Artificial Intelligence*, 1996.
- [73] S. Madvanath and V. Govindaraju. Serial Classifier Combination for Handwritten Word Recognition. In *Proceedings of International Conference on Document Analysis and Recognition, Montreal, Canada*, pages 911–914, 1995.
- [74] B. Plessis, A. Sicsu, L. Heutte, E. Menu, E. Lecolinet O. Debon, and J.-V. Moreau. A Multi-Classifer Combination Strategy for the Recognition of

- Handwritten Cursive Words. In *Proceedings of International Conference on Document Analysis and Recognition*, pages 221–226, 1993.
- [75] U. Mahadevan and R. C. Nagabushnam. Gap Metrics for Word Separation in Handwritten Lines. *International Conference on Document Analysis and Recognition, Montreal, Canada*, pages 124–127, 1995.
- [76] G. Seni and E. Cohen. External Word Segmentation of Off-Line Handwritten Lines. *Pattern Recognition*, 27(1):44–52, 1994.
- [77] U. Mahadevan and S. N. Srihari. Hypotheses Generation for Word Separation in Handwritten Lines. *Fifth International Workshop on Frontiers in Handwriting Recognition, U.K.*, pages 453–456, 1996.
- [78] G. Kim and V. Govindaraju. Handwritten Phrase Recognition as Applied to Street Name Images. *Pattern Recognition*, 31(1):41–51, 1998.
- [79] U. Mahadevan and S. N. Srihari. Robust Vision through Imperfect Eyes: A System for Processing City, State, and Zip codes in Handwritten Address. In *Proceedings of Computer Vision and Pattern Recognition*, 1998.
- [80] G. Kim, V. Govindaraju and S. N. Srihari. Architecture for Handwritten Text Recognition Systems. In *Proceedings of Sixth International Workshop on Frontiers in Handwriting Recognition (IWFHR-VI), Taejon, Korea*, pages 113–122, 1998.
- [81] G. Kim, V. Govindaraju and S. N. Srihari. Extension of Handwritten Word Recognition Method to Street Name Images. In *Proceedings of Sixth International Workshop on Frontiers in Handwriting Recognition*, pages 221–226, 1996.

- [82] G. Kim and V. Govindaraju. Recognition of Handwritten Phrases as Applied to Street Name Images. In *Proceedings of Computer Vision and Pattern Recognition*, pages 459–464, 1996.
- [83] G. Kim, V. Govindaraju and S. N. Srihari. A Segmentation and Recognition Strategy for Handwritten Phrases. In *Proceedings of International Conference of Pattern Recognition*, pages 510–514, 1996.
- [84] S. N. Srihari, V. Govindaraju and J. T. Favata. Unconstrained Handwritten Text Recognition. In *Symposium on Document Image Understanding Technology*, pages 226–235, Oct, 1995.
- [85] J. L. McClelland and D. E. Rumelhart. An Interactive Activation Model of Context Effects in Letter Perception: Part 1. An Account of Basic Findings. *Psychological Review*, 88:375–407, 1981.
- [86] M. S. Seidenberg. The Time Course of Information Activation and Utilization in Visual Word Recognition. *Reading Research*, 5:199–252, 1985.
- [87] D. Marr. *Vision*. Freeman, San Francisco, 1981.
- [88] I. Levi. *Decisions and Revisions*. Cambridge University Press, London, U.K., 1984.
- [89] I. Levi. *The Enterprise of Knowledge*. MIT Press, Cambridge, MA, 1980.
- [90] M. A. Goodrich, W. C. Stirling and R. L. Frost. A Theory of Satisficing Decisions and Control. *IEEE Transactions on System, Man and Cybernetics - Part A*, 28(6):763–779, 1998.

- [91] W. C. Stirling. Coordinated Intelligent Control Via Epistemic Utility Theory. *IEEE Control Systems Magazine*, 13:21–29, Oct. 1993.
- [92] W. C. Stirling and D. R. Morrell. Convex Bayes Decision Theory. *IEEE Transactions on System, Man and Cybernetics*, 21(1):173–183, 1991.
- [93] S. Takatsu. Decomposition of Satisficing Decision Problems. *Information Science*, 22:139–148, 1980.
- [94] T. Matsuda and S. Takatsu. Characterization of Satisficing Decision Criterion. *Information Science*, 17:131–151, 1979.
- [95] T. Matsuda and S. Takatsu. On a Qualitative Theory of Satisfactory Control. *Information Science*, 4:291–313, 1972.
- [96] A. Marzal and E. Vidal. Computation of Normalized Edit Distance and Applications. *IEEE Transactions on. Pattern Analysis and Machine Intelligence*, 15(9):926–932, 1993.
- [97] R. A. Wagner and M. J. Fischer. The String-to-String Correction Problem. *J. Assoc. Comput. Mach.*, 21(1):168–173, 1974.
- [98] G. V. Trunk. A Problem of Dimensionality: A Simple Example. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1(3):306–307, 1979.
- [99] H. Freeman. Computer Processing of Line-Drawing Images. *Computing Surveys*, 6(1):57–97, 1974.
- [100] C.-C. Lu and J. G. Dunham. Highly Efficient Coding Schemes for Contour Lines Based on Chain Code Representations. *IEEE Transactions on Communications*, 39(10):1511–1514, 1991.

- [101] J. J. Hull. A Database for Handwritten Text Recognition Research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, 1994.
- [102] M. D. Garriss and R. A. Wilkinson. *NIST Special Database 3; Binary Images for Handwritten Segmented Characters (HWSC)*. John Wiley & Sons, Inc, New York, 1979.
- [103] R. A. Wilkinson, J. Geist, S. Janet, P. J. Grother, C. J. C. Burges, R. Creecy, B. Hammond, J. J. Hull, N. W. Larsen, T. P. Vogl and C. L. Wilson. The U.S. Bureau of Census and The National Institute of Standards and Technology, Technical Report #NISTIR 4912. In *The First Census Optical Character Recognition Systems Conference*, 1992.
- [104] E. Remias, G. Sheikholeslami and A. Zhang. Block-Oriented Image Decomposition and Retrieval in Image Database System. In *Proceedings of International Workshop on Multi-media Database Management Systems*, New York, 1996.