

A STOCHASTIC FRAMEWORK FOR FONT-INDEPENDENT DEVANAGARI OCR

by

Suryaprakash Kompalli

January 24th 2007

A dissertation

submitted to the Faculty of the Graduate School

of State University of New York at Buffalo

in partial fulfillment of the requirements

for the degree of

Doctor of Philosophy

Department of Computer Science and Engineering

©Suryaprakash Kompalli, 2007

All Rights Reserved

To My Family

Acknowledgments

I would like to express my deep appreciation to Dr. Venu Govindaraju, my advisor and chair of the dissertation committee. His expertise, understanding, and consistent support have been central to this research. He led me to the frontiers of OCR and motivated me to tackle the challenging problems in this area.

I would also like to thank Dr. Sargur Srihari, member of my dissertation committee. His lectures on Pattern Recognition helped me immensely in developing research ideas crucial to this dissertation.

I am also grateful to Dr. Peter Scott, member of my dissertation committee. His excellent courses on Machine Vision, and Animate Vision Principles formed a solid background to my research work and introduced me to novel applications in my chosen field.

Special thanks to Srirangaraj Setlur whose guidance was extremely helpful at every stage of this research. Thanks also to Dr. Kalika Bali of Microsoft Research, India and Dr. Sriganesh Madhavath of HP Labs, Bangalore whose early discussions on language models provided an important dimension to this dissertation. Thanks also to Dr. Kumar Chellapilla of Microsoft Research, Seattle and Dr. George Thoma of NIH, for facilitating the enriching experience of working on industry-standard projects.

Thanks to Dr. Zhixin Shi, Ramanaprasad Vemulapati, Philip Kilinskas, and Edward Sobczak who provided immense assistance and technical guidance to create a real-time OCR system. Thanks also to Mary Jane Gallo, and Eugenia Smith, for their kind administrative support.

My deepest thanks to fellow students Anurag Bhardwaj, Faisal Farooq, Ifeoma Nwogu, Karthik Prasanna, Dr. Sergey Tulyakov, Sankalp Nayak, Sumit Manocha, Swapnil Khedekar, Praveer Mansukhani, and Vartika Singh for the many discussions on Pattern Recognition, Image Processing, and other areas of computer science that provided an ideal backdrop for this research. This work was funded by NSF grants IIS-0112059, IIS-0535038, and by the Mark Diamond Research Fund of University at Buffalo.

Abstract

Font-independent OCR solutions for Latin and Oriental scripts are commercially available and widely used in Digital library applications. However, accurate OCRs are still not available for Devanagari, a script used by over 400 million people in more than forty languages including Hindi and Sanskrit. Challenges in Devanagari OCR include: (i) Large number of classes resulting from a relatively large character set, (ii) Character shapes made of complex primitives that cannot be easily segmented using conventional character segmentation approaches, (iii) Variable representations of the same character in different fonts, and (iv) Preponderance of documents that are of poor print quality or poor quality paper that causes unpredictable character distortions even in single font documents. Prior techniques in Devanagari OCR mostly use font-dependent structural features like average height, width, and stroke thickness to segment characters and employ template matching for classification. A drawback with these methods is that segmentation errors cannot be easily corrected by classifiers. An alternative approach would be to classify characters without performing segmentation, where character shapes are modeled using shape descriptors learnt from multiple character samples, or by using font-specific templates. Given the large number of characters used in Devanagari script, a significant number of samples for each character are usually not available for training. Therefore, the models cannot represent all the shape variations seen in printed documents, leading to poor accuracy.

We address this challenge by segmenting the characters into components which are horizontally or vertically juxtaposed, and connected along non-linear boundaries. Most techniques in the literature have approximated the segmentation process by using sliding windows or projection profiles in a single direction. We adopt a Block Adjacency Graph (BAG) representation, where each

node of the BAG represents a part of the character image and the edges represent their interconnections. Characters are segmented by selecting subgraphs while also accommodating the natural breaks and joints in characters and the various ways in which alphabets can join.

Instead of the common approach of using font-dependent rules to guide the segmentation and classification process, we have developed a recognition driven segmentation method that generates multiple segmentation results for each character. Word hypotheses are generated by integrating the image recognition results with a language model that encodes frequencies of alphabets and syllabic characters. This is in contrast with the previous use of language models primarily as a post-processing technique. We take advantage of the syllabic-alphabetic nature of the Devanagari script by designing a stochastic framework where the primitives are syllabic characters made of one or more alphabets. We use dictionary lookup to enhance the word hypotheses. On a publicly available, multi-font test set of 10,606 words, we have achieved top choice word accuracy of 75%, and top-5 choice word accuracy of 85%. This is a significant improvement over the performance of previous techniques on the same test set.

Contents

1	Introduction	1
1.1	Background	1
1.2	Problem Definition	2
1.3	Devanagari Script and Terminology	3
1.4	Related Work	6
1.4.1	Prior Work in Classification and Segmentation	6
1.4.2	Related Post-processing Approaches	10
1.5	Motivation	11
1.6	Challenges	13
1.7	Proposed Method	14
1.7.1	Data Collection and Baseline System	14
1.7.2	Study of Data Structures	15
1.7.3	Recognition Driven Segmentation	15
1.7.4	Language Model	16
1.8	Dissertation Outline	16
2	Prior Work in OCR and Language Modeling	17
2.1	Prior Devanagari OCR Methods	17
2.1.1	Character Based Devanagari OCR	18

2.1.2	Component Based Devanagari OCR	19
2.2	Prior Techniques in Latin and CJK Recognition	22
2.2.1	Utility of the Recognition-driven Segmentation Model	23
2.2.2	Recognition-driven Segmentation Techniques in Latin OCR	24
2.2.3	Techniques in CJK OCR	27
2.3	Prior Techniques in Language Modeling	28
2.3.1	Dictionary based models	28
2.3.2	Non-dictionary Models	29
2.3.3	Language Models in Devanagari OCR	30
2.4	Discussion	31
3	Data Collection and Baseline OCR	32
3.1	Document Image Analysis and Data Collection	33
3.2	Baseline OCR	36
3.2.1	Segmenting Words Into Components	36
3.2.2	Component Classifier	37
3.2.3	Composing Words from Component Results	39
3.2.4	Generating Training and Test Samples	39
3.3	Testing the Baseline OCR	40
3.3.1	Segmenting Words into Components	40
3.3.2	Core Component Classification	42
3.3.3	Word Recognition	42
3.3.4	Considerations for Proposed System	43
4	Data Structures for Word Recognition	45
4.1	Representation Techniques	45
4.1.1	Bitmap Encoding	46

4.1.2	Chain-code Representation	47
4.1.3	Lossless Compression of Images	47
4.1.4	Horizontal/Vertical Run-length Encoding	48
4.1.5	Block Adjacency Graph (BAG)	50
4.2	Image Processing Using Run-length Representation	51
4.2.1	Slant Angle Detection and Correction	52
4.2.2	Contour Smoothing Using BAG	54
4.2.3	Baseline Estimation	55
4.2.4	Character Segmentation	57
4.3	Using BAG in OCR	59
5	Recognition Driven Segmentation of Devanagari Characters	60
5.1	Segmentation by Block Adjacency Graphs	62
5.1.1	Construction of BAG	63
5.1.2	Character Segmentation Using BAG	65
5.1.3	Connecting Broken Junctions in the BAG	67
5.2	Using Component Classifiers to Guide Segmentation	68
5.2.1	Feature Extraction Techniques	72
5.2.2	Classifier Design	72
5.2.3	Obtain Threshold Confidences for Classifiers	73
5.2.4	Ranking Segmentation Results	75
5.2.5	Necessity of Additional Models to Generate Words	76
6	Stochastic Modeling to Generate Words	81
6.1	Linguistic Restrictions in Hindi	82
6.1.1	Script Writing Grammar	82
6.1.2	Restrictions on Syllable Co-occurrence	84

6.2	Algorithms to Model Linguistic Phenomena	87
6.2.1	Corpus Analysis of Hindi N-grams	88
6.2.2	Stochastic Finite State Automaton	89
6.2.3	Implementing Ngram Statistics in the SFSA	92
6.3	Generating Words from Components	95
7	Experiments	98
7.1	Character Recognition	99
7.2	Word Recognition	101
8	Conclusions	105
8.1	Summary	105
8.2	Contributions	106
8.3	Future Work	108
8.3.1	Recognition Driven Segmentation	108
8.3.2	Pruning Top-n Word Choices	108
8.3.3	Applying the OCR to Other Languages	109
A	Theoretical Relation of Scripts and Sounds	110
A.1	Articulation: The Generation of Sounds	110
A.2	Phonological Syllable	112
A.3	Orthography	113
A.3.1	The Orthographic Syllable in Devanagari	114
A.4	Phonological Studies in Hindi	115
B	Manual for “ILT Indic Digital Library Toolkit”	120
B.1	Truthing Interface	121
B.1.1	Menu Options	121

B.1.2	Truthing Interface Tool bar Options	122
B.1.3	Image Panel	122
B.1.4	Transliteration Panel	123
B.2	Devanagari Script Analysis Interface	123
B.2.1	Menu Options	123
B.2.2	Truthing Control	124
B.2.3	Devanagari Script Analysis	124
B.2.4	Query Panel	124
B.3	Multi-Lingual Translation Interface	125
B.3.1	Menu Options	125
B.3.2	Text Controls	129

Bibliography	133
---------------------	------------

List of Figures

1.1	Character zones in a Devanagari word (top) and splitting of a Devanagari word (bottom)	4
1.2	Devanagari text samples (a) Descenders, (b) Fragmented conjuncts have the same separation that is seen between characters in words, (c) Fused characters	8
1.3	Character segmentation in Latin and Chinese OCR: (a) Sliding window technique for fused character segmentation (Technique from [88]), (b) Segmenting Chinese characters into straight line primitives (Technique from [82]), (c) Devanagari characters (top) and non-linear primitives (bottom)	9
1.4	Stroop test to study the use of sound cues in reading	12
1.5	(a) Devanagari words in different fonts. (b) Devanagari characters having alphabets connected from left to right or top to bottom.	13
2.1	Segmentation involved in character and component based recognition	18
2.2	Character recognition using decision tree	20
2.3	Segmenting conjuncts and descenders to obtain components [8]	21
2.4	Errors due to font-dependent segmentation.	22
2.5	OCR paradigms based on segmentation methodology [15]	23
2.6	Examples of Latin OCR techniques which follow different OCR paradigms	25

2.7	(a): Devanagari characters which need to be segmented along non-linear windows, (b) Devanagari characters which need to be segmented from top to bottom, or from left to right	26
2.8	Example illustrating shape-based matching.	29
3.1	DIA stages used to obtain paragraphs, lines, and words from a Devanagari docu- ment [38].	34
3.2	Examples of documents in the ILT dataset	35
3.3	Segmenting words into components	38
3.4	Parsing a word on query: Extracting component truth (Left), Extracting compo- nent boundaries (right), Combining the two results to obtain component images (bottom).	40
3.5	Word accuracy on different documents using the segmentation driven recognizer. Documents of type 4 are printed in a consistent font, while the other documents have variation in font and print quality. The accuracy varies not only across differ- ent document types, but is very different even in documents coming from the same type.	43
4.1	(a) Original image, (b) HIPS representation [87]	46
4.2	(a) Devanagari alphabet, (b) Horizontal run lengths, 36 in number, (c) Vertical run length, 34 in number.	48
4.3	Obtaining block adjacency graph: (a,b,c) Using horizontal runs (Hruns) to create BAG. (d,e) Using vertical runs to create BAG. The BAG is seen in (c,e) as a set of nodes connected by black edges.	51
4.4	Slant angle detection using run-lengths [87].	53
4.5	Morphological opening	54

4.6	Image smoothing using BAG: (a,b) BAG with small blocks, (c,d) Spurious blocks are removed. Narrow joins and small blocks are not deleted, but merged with larger blocks.	55
4.7	The use of BAG in character segmentation	58
5.1	Recognition driven segmentation	62
5.2	Using BAG for recognition	64
5.3	Creating the BAG data structure for a Devanagari word	69
5.4	Results of recognition driven segmentation using Block Adjacency Graph.	71
5.5	Classifier receiver operator characteristics for a full consonant: (a) Neural network with gradient features, (b) Neural Network with GSC features, (c) K-Nearest neighbor with GSC features	74
5.6	Results from the recognition driven segmentation performed in Figure 5.4; Classifier scores are not shown here	76
5.7	Examples of results from recognition driven segmentation. Word images are on the left and component results are on the right.	78
5.8	Examples of results from recognition driven segmentation where some core components correspond to fused or conjunct characters. Word images are on the left and component results are on the right.	79
5.9	Examples of recognition driven segmentation where some core components have been over-segmented. Word images are on the left and component results are on the right.	80
6.1	component sequences (left) and possible words (right).	82
6.2	Removing component sequences by applying linguistic restrictions.	83
6.3	Finite State Automata that can accept the classifier labels in Figure 5.4 and emit the correct word पूछती	91

7.1	Overview of results from each stage of recognition	99
7.2	Classifier accuracies for 8,029 consonant/vowel images, Average accuracy: Neural network with gradient features: 94.36%, Neural network with GSC features: 95.25%, and K-nearest neighbor with GSC features: 96.97%	100
7.3	Word accuracy on different documents (a,b): Using the recognition driven segmentation approach, (c): Using the segmentation driven approach. The accuracy of segmentation driven approach varies widely across different documents, while the recognition driven segmentation approach provides more consistent accuracy .	103
A.1	Syllable structure	113
A.2	Different types of phonological scripts	114
B.1	Main window options	121
B.2	Truthing Interface Menu: (a) File Options and (b) Image Options	121
B.3	Menu Options in the Devanagari Script Analysis Interface	123
B.4	Empty truthing tool: (a) Image Panel, (b) Transliteration Panel, (c) GUI Keyboard, (d) Truth Panel, (e) Tool bar and Menu options	126
B.5	Truthing interface: Please refer to Figure B.4 for legend. Red outline is for the word being truthed.	126
B.6	Devanagari Script Analysis, Truthing Interface: (f) Thumbnail panel, (g) Tool bar, and (h) Tab pane; W: Word level truth, C: Character level truth, and Q: Query Panel. 127	
B.7	Devanagari Script Analysis, Segmentation Analysis: (a) Images of the word and its' components, (c) Word and component truth, (d) Tree layout of the truth.	127
B.8	Devanagari Script Analysis, Query Panel: The display shows a query for the component “ka” क : (a) Query transliteration input, (b) Query display, and (c) Query options	128

B.9	Error Report Creator: (a) Error Type, (b) Truth File Name, (c) Line and Word Number, (d) Error Description, and (f) Stack Trace.	131
B.10	Menu Options for the Multi-Lingual Translation Interface	131
B.11	Text Controls for the Multi-Lingual Translation Interface	131
B.12	Multi-Lingual Translation Interface: Multi-Lingual Translation Interface: (I) <i>Transliteration window</i> to enter transliterated English text for scripts, (II) <i>Annotation window</i> to view parsed text along with annotation, (III) <i>Script window</i> displays transliterated text without the associated annotations. (IV) GUI keyboard for text entry, (V) Text controls and tool bar. (a) Text to which annotation has been added, (b) Annotation for pink colored region, (c) Text without any annotation, (f,i,g,h) Corresponding parsed text.	132

List of Tables

1.1	Devanagari alphabets with their IPA (International Phonetic Association [48, 70]) pronunciation and their phonetic attributes.	5
1.2	Types of Devanagari characters and constituent components	6
3.1	Character segmentation errors	41
3.2	Core component classifier: (I): Neural network classifiers with manual pre-classification, (II): Neural network classifiers with automatic pre-classification	42
3.3	Comparison of Recognition Results	43
3.4	Types of characters and corresponding frequencies evaluated from the ILT [58] image data set, and the nearly 1000 times larger EMILLE [5] text corpus.	44
5.1	Devanagari components recognized by our system	70
6.1	Formation of alphabets from components	85
6.2	Bigram and trigram perplexity computed on the EMILLE dataset	89
6.3	Word recognition results after applying language models and dictionary lookup. Word results are sorted according to the score returned by the SFSA.	97
7.1	False reject rate (FRR) and false accept rate (FAR) obtained at the equal error rate for consonant/vowel classifiers, Testing is done on 8,029 consonants and vowels, 479 consonant-descenders and 686 conjuncts	100

7.2	Recognition results from 479 descender and 686 conjunct characters	101
7.3	Comparison of conjunct recognition using segmentation driven and BAG based techniques	102
7.4	Accuracy and average string edit distance using different n-gram models. The test set has 10,606 words	104
7.5	Comparison of Recognition Results	104
A.1	Place of articulation. *These dont exist in English	112
B.1	Query options: Examples are results of queries executed over the domain कारन कर्ता कर्म क्यारि	125

Chapter 1

Introduction

1.1 Background

Optical character recognition (OCR) is the process of automatically converting machine printed documents to editable text. OCR combines the fields of pattern recognition, computer vision, and natural language processing to “read printed text”, a task which previously could be accomplished only by manual transcription. OCRs have facilitated several applications like automated labeling of documents [9, 104], creating digital libraries from historical and contemporary literary collections [3, 29, 47], and remastering paper documents for digital distribution [13, 83].

Latin and Oriental script OCR currently report high accuracy rates on printed documents, with several easily available commercial implementations. Similar advances have not been reflected in Devanagari script recognition. While Latin OCR techniques handle multiple fonts with high accuracy and have been integrated into several applications [69], Devanagari OCRs presently focus on processing single font documents [6, 18, 64]. This thesis develops a methodology for font-independent recognition of Devanagari text.

A subarea of pattern recognition is the use of language models in enhancing recognition results. Recognition systems are often meant to be used on natural speech, handwriting, or publishing.

Linguistic patterns can be inferred from real data and used to assist the recognition process. For example if an OCR is processing documents from the medical domain, the word sequence “brain cell” can be considered more likely than “brain call”. Similarly, in English text, the alphabet “q” is always followed by the alphabet “u”. Such knowledge can help in correcting errors that would otherwise pass undetected. The detection and correction of errors can become better if the language model can more accurately represent the patterns in real data. Previous recognition systems have considered words or alphabets as linguistic units [12, 44, 51, 86], thereby modeling word or alphabet distributions. Word models are typically applied after applying models that operate on sub-words. We design a language model by using syllables as the linguistic units, and show that they can give better results than alphabet models. Further sections of the chapter are devoted to problem definition, the basics and terminology of Devanagari script, related work, our motivation, and proposed methodology.

1.2 Problem Definition

Converting the image of a machine printed document into an editable format involves several stages: (i) Noise removal that erases specs, blobs and other artifacts to obtain a “clean” document image [1, 107, 108], (ii) Binarization that converts a gray-scale image into a bitonal image [98], (iii) Correcting the document slant [17, 54], (iv) Separating text regions from lines, photographs and other non-alphabetic data [16, 19, 53] and (v) Recognizing the text. This thesis concentrates only on the problem of recognizing text from word images and defines the multi-font Devanagari OCR problem as follows:

- **Input:** Binary image of a printed Devanagari word.
- **Output:** List of word results along with their scores.

A significant point is that the problem is defined to be independent of visual statistics that can be extracted from a paragraph or line of text. Many Devanagari OCR techniques require a paragraph or line of text to be presented to the OCR. Such OCRs extract structural statistics like character height, character width, number of vertical strokes etc from character images in a paragraph or line and use the statistics while recognizing each word. Our viewpoint is that relying on information embedded in these measures is akin to the extraction of font information, and is highly error-prone. We will show that these statistics are consistent only in the single font scenario and that they are unreliable for multi-font documents. Images of paragraphs or lines can contribute to word recognition when using a word or sentence based language model. However, our language model deals with units within a word and does not cross word boundaries.

To enable the processing of document images with our OCR system, we have implemented a few previously established algorithms for noise removal, binarization, text-artifact separation, and line/word separation [55]. The paucity of publicly available Devanagari data sets and usable implementations of prior OCR techniques has led to the design of our own data sets and the use of an in-house baseline OCR for comparison purposes. The data sets, tools and OCRs developed as part of the thesis have been made available to the research community¹.

1.3 Devanagari Script and Terminology

This section outlines characteristics of Devanagari and linguistic definitions that are relevant to our OCR problem. Appendix A provides more details on the subject.

Devanagari belongs to a class of scripts which are known as alphasyllabaries, alphabetic-syllabic, or abugida [24]. According to linguistic definitions (Chapter 1, 31 of [24]), a characteristic of these scripts is that they have a basic set of symbols each of which represents a specific sound. One or more of these symbols may be combined to form a single visual unit. Further, each

¹The data and tools are available online at: <http://www.cedar.buffalo.edu/ilt>. The OCR is located at: <http://phoebe.cedar.buffalo.edu/ocr/>

visual unit corresponds to the onset and nucleus of a syllable, a well defined unit of pronunciation (Formal definition given in section A.2). For example, क and य are two symbols in Devanagari and are pronounced as *k* (*ka* without the ending *a* sound) and *ya* respectively. They can be fused together to write क्य , which is pronounced as *kya*. Let us consider a script which belongs to a different class - (Say) Latin, which belongs to the class of alphabetic scripts. When Latin is used to write English, alphabets are supposed to be printed as isolated units and are connected with each other only in cursive writing; further, the pronunciation of Latin alphabets can change with context, E.g: *o* is different in *book* and *dog*.

The research community uses a diverse terminology while referring to the linguistic units discussed above. Phonological studies on syllabic-alphabetic scripts refer to the visual unit as an *orthographic syllable* [91], while the OCR domain refers to them as *characters* [6, 36]. The basic set of symbols are referred to as *glyphs* in some phonological studies [91] and as *alphabets* both in OCR and phonological research [6, 24, 36, 100]. We use the term “character” for each visual unit and the term “alphabet” for each symbol in the basic set. In our OCR design, certain alphabets are segmented into smaller parts to assist recognition. We use the term “components”, to mean the parts into which an alphabet image can be segmented during the OCR process. Figure 1.1 shows an example of how the terminology is used to refer to different parts of a Devanagari word.

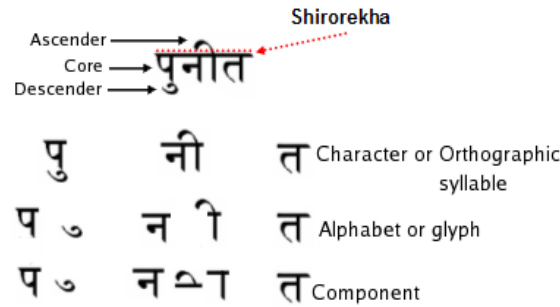


Figure 1.1: Character zones in a Devanagari word (top) and splitting of a Devanagari word (bottom)

The Devanagari alphabet set contains consonants, half-consonants, vowels, vowel modifiers and digits (Table 1.1). Devanagari digits are similar in usage to their counterparts in Latin script.

	Voiceless Plosives/Stops				Voiced Plosives/Stops				Nasals		
	Un-Aspirated		Aspirated		Un-Aspirated		Aspirated				
Velar	क	क	क ^h	ख	ग	ग	ग ^h	घ	ङ	-	ङ
Palatal	च	च	च ^h	छ	ज	ज	ज ^h	झ	ञ	ञ	ञ
Retroflex	ट	-	ट ^h	ठ	ड	-	ड ^h	ढ	ण	ण	ण
Dental	त	त	त ^h	थ	द	-	द ^h	ध	न	न	न
Labial	प	प	प ^h	फ	ब	ब	ब ^h	भ	म	म	म

	Palatal		Retroflex		Dental		Labial		Glottal	
Sonorants	य	य	र	-	ल	ल	व	व	ह	-
Sibilants	श	श	ष	ष	स	स				

(a) Consonants, consonant half-forms and their pronunciation

	Front		Central		Back	
High	इ	ई	उ	ऊ	ऋ	ॠ
Mid	ए	ऐ	ओ	औ	ऌ	ॡ
Low	अ	आ	इ	ऋ	उ	ऊ

(b) Vowels, vowel modifiers and their pronunciation

० १ २ ३ ४ ५ ६ ७ ८ ९
। अँ : ँ ं ः

(c) Digits and special symbols

Table 1.1: Devanagari alphabets with their IPA (International Phonetic Association [48, 70]) pronunciation and their phonetic attributes.

Devanagari characters, however may contain one or more alphabets whereas Latin characters are synonymous with Latin alphabets. Effectively, a large number of Devanagari characters may be formed by combining different alphabets: A vowel or consonant may be considered as a character (Eg: क अ औ), conjunct or compound characters are produced by combining half-consonants with consonants (Eg: क + य → क्य), and consonant-modifier and conjunct-modifier characters are produced by combining consonants and conjuncts with vowel modifiers (Eg: क + ँ → का, क्य + ँ → क्यं).

Devanagari is used to write several languages like Hindi, Marathi, Sanskrit etc. We will be studying the use of Devanagari in writing Hindi. Most consonant and vowel forms are written from left to right and joined together into words by a horizontal line, called the shirorekha (Fig-

ure 1.1). We refer to these as core components in this thesis. Some modifiers may be written above the shirorekha and we refer to them as ascender components or ascenders. Modifiers written below the consonants are called descender components or descenders. In our work, we refer to different characters by using a nomenclature that indicates the location of constituent components (Table 1.2).

Character type	Modifier type	Sample Character	Alphabet(s)	Component(s)
Consonant/vowel	-	य	य	य
Consonant-descender	Descender	यु	य ॡ	य ॡ
Consonant-ascender	Ascender	ये	य ॠ	य ॠ
Consonant-core	Core	या	य ा	य ा
Consonant-ascender	Ascender/core	यी	य ी	य ा ॠ
Conjunct	-	क्य	क य	क य
Conjunct-descender	Descender	क्यु	क य ॡ	क य ॡ
Conjunct-ascender	Ascender	स्थे	र थ ॠ	र थ ॠ
Conjunct-core	Core	क्या	क य ा	क य ा
Conjunct-ascender	Ascender/core	क्यों	क य ी	क य ा ॠ

Table 1.2: Types of Devanagari characters and constituent components

1.4 Related Work

1.4.1 Prior Work in Classification and Segmentation

Most Devanagari OCR techniques split words into smaller units, though implementations vary on the level of segmentation. While some approaches use characters as classification units [18, 71], others segment a character into components before classification [6, 64]. When characters are used as classification units, horizontal and vertical profile of each word is examined to remove the shirorekha and separate the characters. In most Devanagari words, the shirorekha is very well-defined and projection-profile techniques are fairly accurate for this stage. Each character image is classified using decision trees, with the class space being pruned at each node [18]. These systems need sufficient training samples for each character class, and do not include infrequent conjunct and

consonant-modifier characters in their class space. Images obtained after removing the shirorekha correspond to Devanagari characters. The images are input to a classifier that returns a character class label. The word result is obtained by concatenating character class labels.

Analysis of the EMILLE corpora [5] shows that Devanagari text has a sparse character distribution, with a large number of classes having few samples. Character classifiers often omit the rare characters from their class space, or use the few available samples as training data. In either cases, this leads to an adverse reduction of overall accuracy. Consider for example the ILT image dataset which has 32,413 characters [58]: 537 characters occur less than 7 times each, while 426 classes have just over 60 samples each. The former set of 537 characters will easily classify as rare characters and would be left out of the class space. Due to the sheer number of classes, they constitute 12% of the total text, and ignoring these characters or training on the very few samples would degrade OCR results. Sufficient samples of rare characters cannot be collected by simply increasing the amount of data collection. Because multiple characters can be formed from Devanagari alphabets, the number of alphabet combinations increases as we observe more data. Chapter 3 describes the aspect of data collection and character frequencies in Devanagari.

Another challenge for character-based classifiers is that font changes in Devanagari are more significant at the character level than at the component level. The shape of a component does exhibit some variation across different fonts, but the difference in fonts is most apparent when components are combined to form characters. Descender components are often added at different locations based on the font - For example Figure 1.2(a) shows samples where the modifier “ृ” is attached differently on core components क and ह. Print variations also alter the positions of components, leading to characters that may be disjoint in some words and connected in others (Figure 1.2(b)). Collecting representative samples of all these problems and modeling all such variations is not easy.

Some recognizers segment a character before classification [6, 64] to take advantage of the fact that character classes are large in number, but are composed of a smaller set of components.

The number of characters can range in the hundreds (973-5537 depending on the language/data set, Table 3.4), while the number of components is much less (we use 52 components, Table 5.1). Attributes like location and number of vertical strokes, average width, and average height are used to distinguish conjuncts and consonant-descenders from consonants/vowels. These features are used in a rule-based technique to segment conjuncts and consonant-descenders. This approach is not robust to font and print variation and incorrect segmentation may occur in many cases. For example, we examine descender segmentation [8], which is done by using average character height to predict a region of interest (ROI). The ROI is examined using run-length analysis to separate descenders. However, variation in height often causes over-segmentation. Figure 1.2(a) shows how the descender “ु” can be correctly separated from “क” and “ह”, but at the same time the character “इ” is over-segmented. Rule-based segmentation schemes also fail to account for fused characters where two or more alphabets are illegally touching each other (Figure 1.2(c)). The accuracy of these approaches is limited by the accuracy of segmentation. Additionally, sequential application of segmentation followed by component recognition does not lend itself to post-processing.

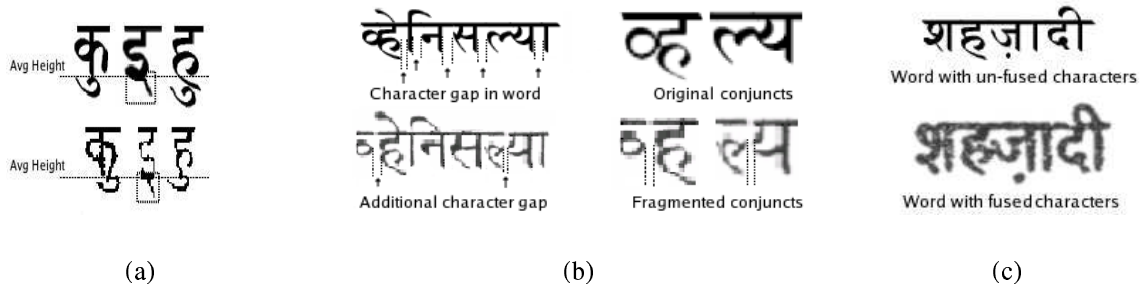


Figure 1.2: Devanagari text samples (a) Descenders, (b) Fragmented conjuncts have the same separation that is seen between characters in words, (c) Fused characters

Researchers in Devanagari OCR have so far dealt with single font systems, but the approaches cannot scale to font-independent OCR. In multi-font Devanagari OCR, problems are mainly caused by the large number of character classes and the unpredictable character shapes. Latin and CJK recognition have been able to handle similar problems, and report high accuracies on multi-font

documents. We study some of the techniques in Latin and CJK recognition and examine if they can be applied to Devanagari.

1) Segmentation Techniques for Latin Scripts: Sliding window with selective segmentation has been used in Latin OCR to separate fused characters [15, 88] (Figure 1.3(a)). A rectangular window is used to horizontally divide fused characters into two parts. Classifiers are used to reject or accept segmentation results. In Devanagari, conjunct and descender characters need to be segmented at non-linear boundaries in two directions: horizontally or vertically. On the other hand, reported schemes either use a linear segmentation window or segment characters along a single direction.

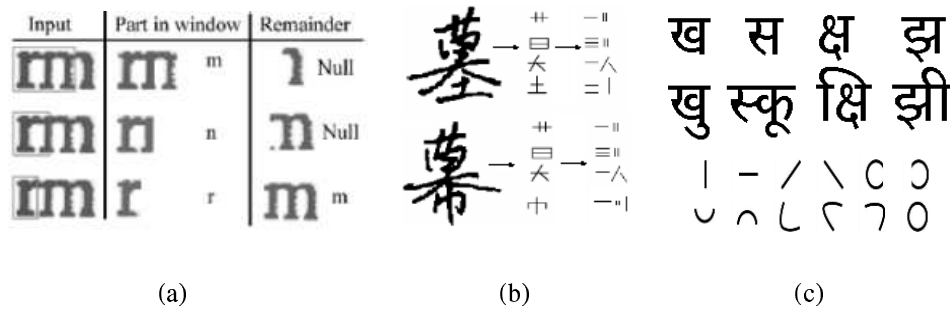


Figure 1.3: Character segmentation in Latin and Chinese OCR: (a) Sliding window technique for fused character segmentation (Technique from [88]), (b) Segmenting Chinese characters into straight line primitives (Technique from [82]), (c) Devanagari characters (top) and non-linear primitives (bottom)

2) Classifiers for Large Chinese Character Set: One approach seen in Chinese OCR is to divide the large character set (3,000-7,000 characters) into groups [82, 96]. Each group is associated with the presence or absence of a characteristic line (linear primitive). At run-time, linear primitives extracted from a test character are used to limit the recognition hypothesis to one of these groups, and detail features are used to classify within the group. Another technique to recognize the large character set is to model a character as combinations of linear primitives [94, 105]. Results from classifying such linear primitives are used to derive character hypotheses (Figure 1.3(b)).

While these systems rely on identifying linear primitives from a character, Devanagari characters are composed of non-linear primitives (Figure 1.3(c)). No technique exists to obtain non-linear primitives from Devanagari characters.

1.4.2 Related Post-processing Approaches

In Devanagari, character formation from components is governed by a set of rules called script writing grammar or script composition syntax. Sinha demonstrated the utility of these rules in correcting simulated OCR errors [84]. Examples of script writing rules are: (i) Vowel modifiers should be associated with a consonant, and (ii) only one vowel modifier must be present in a character. Dictionary lookup is used by Bansal et al to prune word recognition results [7].

It is common knowledge that probabilistic models can enhance recognition results before using a dictionary. Examples include alphabet n-grams which have been accommodated in Latin OCR [12, 44], word n-grams and trigger pairs which have been analyzed for speech recognition systems [80]. Kukich provides an excellent survey of some of these techniques [59]. At the sub-word level, attempts have been made to design models based on morphemes, or suffix-prefix relations [42, 86].

Probabilistic language models are intended to be algorithmic representations of natural language phenomena. For instance, word n-gram models are based on the natural phenomena of related words appearing together in a paragraph of text. Most of the probabilistic language models used in OCR have concentrated on studying alphabets, or words as units of the language model. Devanagari OCR techniques have been limited to the rule-based implementations of script writing grammar or dictionary based comparison. In this thesis, we explore the use of syllables and sound cues in processing language.

1.5 Motivation

Psychological studies show that humans use auditory perception as a support system for reading [41, 97]. To describe this phenomena, we cite an experiment in which users are presented with words printed in different colors (Figure 1.4(a), 1.4(b)). The user is required to press keys corresponding to the color in which each word is written (Figure 1.4(d)). The time taken to press the sequence of keys is recorded for each line. The observation in this study is that users take less time to press the keys when the underlying sound is similar to the color of the word, but take more time when the sounds and colors are dissimilar. The original test and psychological theory is detailed in [10, 93] and arrives at the conclusion that humans automatically use phonological or sound models to read text. Independent studies have confirmed that this conclusion holds true in several languages like English, French, Hebrew etc. where specific sounds are associated with alphabets [14, 76]. Human studies outlined in [27, 90] show that auditory perception is used even to read pictographic scripts like Japanese and Chinese which have no concept of an alphabet, or a sound being directly associated with an alphabet. These studies further suggest that the linguistic units used by humans correspond more closely to syllables than to sub-syllable units like alphabets or extra-syllable units like prefixes. In [70] it is claimed that a set of rules defined using syllables can represent the phonological structure of Hindi (Details in section A.4). The study includes experiments which show that native readers of Hindi very quickly reject non-words that violate these rules, but have difficulty in rejecting non-words which conform to the rules.

While the use of sound models in human reading is well established, there are some difficulties in incorporating sound information in a computational model in OCR systems. In many scripts, the mapping from written symbol to the spoken sound is not direct. In alphabetic scripts like Latin, the pronunciation of the alphabet changes with context (E.g. water and owl) and additional computational models are required to map written symbols to sound [91]. These variations in pronunciation occur very rarely in syllabic-alphabetic scripts like Devanagari. In Devanagari, each

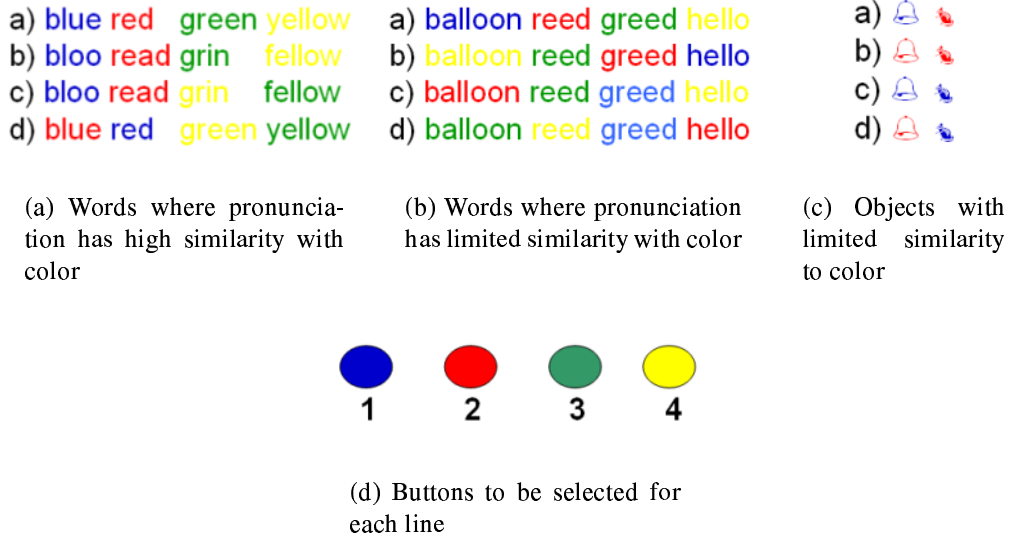


Figure 1.4: Stroop test to study the use of sound cues in reading

alphabet is assigned a fixed pronunciation (Table 1.1). For example: ण is said to be a *Labial*, i.e. it is generated from the lips; ण is a *dental*, generated from the teeth and so on. Since this pronunciation does not change with context, the structure and statistics of Devanagari phonetic units like syllables and phonemes can be extended to characters.

The recognition-driven design for our OCR is inspired by the analysis of different character segmentation paradigms, originally proposed by Casey et. al: “dissection techniques for segmentation, recognition driven segmentation and holistic recognition” [15]. It is well known that dissection techniques have the intractable drawback of propagating segmentation errors. Holistic recognition has been popular with lexicon-driven HWR (handwriting recognition) [57, 112], but has not been widely used in OCR. Recognition driven segmentation has long been the focus of both Latin and Chinese OCR, with high accuracy rates being reported in both approaches [62, 88, 102]. Most of the attempted techniques in Devanagari OCR have followed dissection based approaches, with limited success. We attempt to design a Devanagari OCR system that follows the recognition driven segmentation paradigm, combining all stages of recognition into a single integrated model.

1.6 Challenges

Some of the research problems in Devanagari OCR are:

(a) A large number of in-frequent character classes are seen in DL collections (973 character classes in ILT image dataset [58], 5,537 in EMILLE text corpus [5]). In current OCR approaches, many characters are misclassified or not included in the classifier training set (Section 1.4.1).

(b) Font variations change the relative positions of combining alphabets, significantly altering the shapes of consonant-descender and conjunct characters (Figure 1.5(a)). Current techniques for descender and conjunct character segmentation do not address these variations (Section 1.4.1). Sequential application of segmentation and component recognition does not lend itself to post-processing.

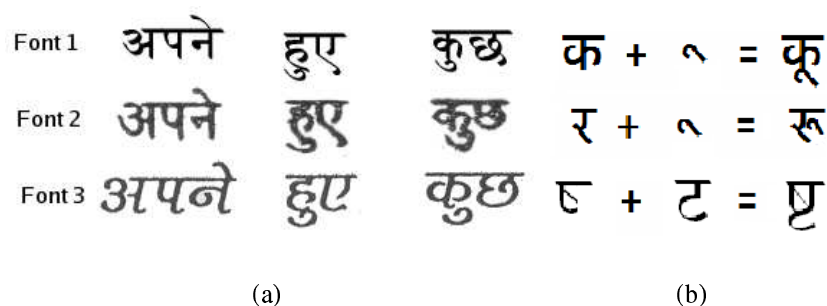


Figure 1.5: (a) Devanagari words in different fonts. (b) Devanagari characters having alphabets connected from left to right or top to bottom.

(c) Poor printing methods may lead to fragmented conjuncts or fused characters (Figure 1.2) which are not handled by current systems.

(d) State of the art techniques for Latin and Oriental script OCR focus on identifying linear primitives from a character, or segmenting characters along a single direction. Devanagari characters are composed of non-linear primitives that may be fused vertically or horizontally (Figure 1.3). Devanagari OCR requires a scheme that allows segmentation of characters along non-linear boundaries in a vertical or horizontal direction. Figure 1.5(b) shows samples of Devanagari characters

which cannot be segmented using linear boundaries or by scanning characters in a single direction.

(e) Post-processing stages using syllable or phoneme restrictions can be adapted into Devanagari OCR. A related issue is that these restrictions fit into a probabilistic framework and others like script writing grammar need rule-based models, and have not been integrated into a single system.

1.7 Proposed Method

In order to address these challenges, we propose a new methodology for Devanagari OCR: (i) Use non-linear boundaries to segment components that may be present at the center, top and bottom of character region, (ii) Design component classifiers which can tolerate segmentation errors, and (iii) Generate word results by combining the rule-based script writing grammar and phonetic language models into a single stochastic framework.

1.7.1 Data Collection and Baseline System

The availability of standard, public-domain data sets is a must to train recognition systems and compare results with the community. Whereas Latin and CJK recognition technologies have benefited immensely from the ready availability of data resources [39, 43, 56], Devanagari OCR is faced with lack of public domain training and test data sets. Our first step has been to collect data sets for Devanagari and provide them for public use.

Standard techniques are available for certain stages in Devanagari OCR like shirorekha removal and are known to be highly accurate. We find it prudent to examine the strengths of contemporary techniques before designing a new system. A segmentation driven recognizer has been designed to provide a baseline system for comparison purposes.

1.7.2 Study of Data Structures

Efficiency of the data structures used for image representation is central to recognition tasks. Data structures impact the amount of memory needed to store an image as well as on how suitable it is for image recognition. One may tend to confuse image representation with feature extraction [99], but these are very separate problems. The latter deals with representing image characteristics like skeletons, structural descriptors or moments that can abstract out pixel level information. The former refers to a data structure that can recreate the original image in a lossless manner.

A previous study compares bitmap, chain-code, fax-compressed, and run-length representation representations on the basis of slant correction and contour smoothing [87]. We extend that work by considering the Block Adjacency Graph representation and evaluate the important aspect of character segmentation. Devanagari characters are segmented by taking advantage of the Block Adjacency Graph (BAG) representation, so that segmentation boundaries can accommodate the natural breaks and joints in a character and the various ways in which alphabets can join even in non-standard fonts.

1.7.3 Recognition Driven Segmentation

The BAG based character segmentation provides multiple segmentation hypotheses. Classifiers can be used to rank these character images and identify ones that closely represent Devanagari alphabets and others that do not. In case of Devanagari, characters need to conform to a set of linguistic conditions which are defined as script composition syntax or script writing grammar. The common approach in Devanagari recognition has been to use hand-crafted, sequential rules for segmentation, classification, and matching with script writing grammar. We design a recognition driven segmentation method using the hypothesis and test paradigm.

1.7.4 Language Model

Motivated by the syllabic-alphabetic nature of Devanagari script, we implement a language model that encodes n-gram frequencies of syllables or phonemes. Previous work in linguistic post-processing techniques have considered alphabets or words as primitives. Our thesis explores a stochastic framework where the primitives are syllables which can include one or more alphabets. Devanagari can be used to write several languages like Sanskrit, Hindi, Marathi. Even though the script remains the same, there are subtle variations in the phonetics across the different languages. This thesis studies the phonetic aspects of Hindi and designs a language model that can enhance Devanagari OCR results.

1.8 Dissertation Outline

Analysis of related prior work is presented in chapter 2. Chapter 3 describes the collection of ground-truthed data and the implementation of a baseline OCR system. Discussion of the Block Adjacency Graph representation and other data structures is provided in chapter 4. Chapter 5 outlines the design of a recognition driven segmentation technique for Devanagari characters. The stochastic language model to obtain word recognition results is outlined in chapter 6. Experiments on our publicly available dataset and conclusions from our work are presented in chapters 7 and 8 respectively.

Chapter 2

Prior Work in OCR and Language Modeling

The problem of OCR has been central to the domain of pattern recognition, with several researchers providing varying solutions to decipher printed text. While section 1.4 provides an introduction to some of the prior art in this field, a more detailed view is presented here.

2.1 Prior Devanagari OCR Methods

Most researchers in Devanagari OCR split words into smaller parts, but the level of segmentation differs across implementations - some techniques separate all the constituent components from a character and classify each component [6, 64, 38] while others treat character shapes as distinct classes [18, 85] (Figure 2.1). A Devanagari word can be segmented into characters by removing the shirorekha, which is often very well defined. On the other hand, segmenting Devanagari characters into components is more challenging, primarily due to the following factors: (i) Characters can have more than one component, (ii) Components can be fused horizontally or vertically, (iii) Components are connected along non-linear boundaries, and (iv) Changes in font type and print

quality vary the relative locations of components. On account of these factors, Devanagari characters have a more complex shape than their constituent components. The complexity of character shapes implies that classifying characters is not easy in comparison with classifying components. In the remaining parts of this section, we will discuss prior techniques in Devanagari OCRs by dividing them into those which use character classifiers and those which use component classifiers.

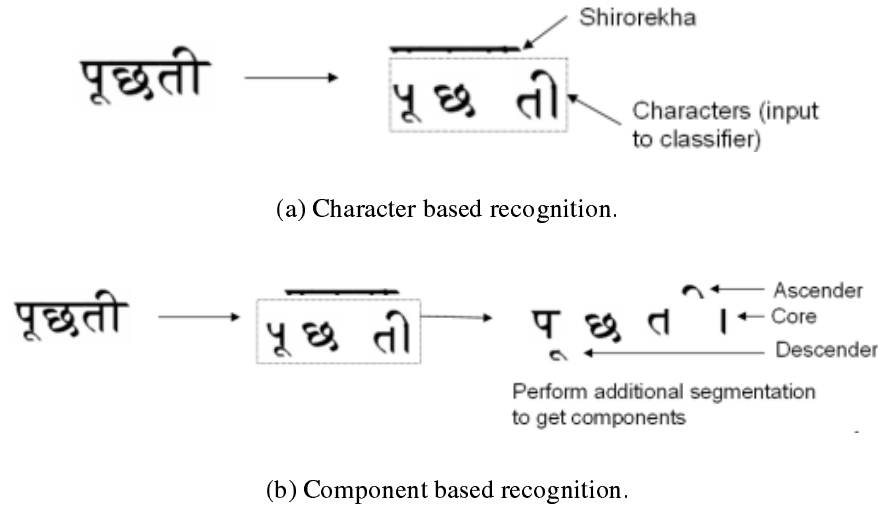


Figure 2.1: Segmentation involved in character and component based recognition

2.1.1 Character Based Devanagari OCR

Character based recognition approaches usually involve only a single stage of segmentation which is the removal of the shirorekha (Figure 2.1(a)). Chaudhary [18] and Khedekar [38] use projection profiles of words to detect and remove the shirorekha. Horizontal projection profiles sum the intensity values along each row of the image. Since the shirorekha extends across the width of the word, a very high value of the horizontal projection profile will be obtained at rows corresponding to the shirorekha. In [17], Chaudhuri outlines an approach where shirorekha detection and skew correction is performed by analyzing the upper envelope of words.

Once the shirorekha is removed, classifiers are used to match the resulting images with the closest character class. Sinha and Mahabala [85] outline a syntactic recognition system for characters. The system generates a picture language for Devanagari characters by extracting structural features and their relative positions. During recognition, structural features of the input character are compared with the stored picture language description, and the closest comparison is selected as the classifier result. Pal and Chaudhuri [18] employ a decision tree structure to classify character images. The character space is pruned at each node of the decision tree by checking for the presence or absence of a specific structural feature. An example is shown in figure 2.2. Pal and Chaudhuri observe that structural features are not sufficient to cover the entire range of Devanagari character classes, and some of the terminal nodes in their decision tree structure have multiple character classes. They use template matching of binarized character images to disambiguate among the multiple character classes.

These character-based systems need sufficient training samples for each character class, and choose to not include infrequent conjunct and consonant-modifier characters in their class space. Devanagari text has a sparse character distribution, with a large number of classes having few samples - consider for example the ILT image dataset which has 32,413 characters [58]: 537 characters occur less than 7 times each, while 426 classes have just over 60 samples each. The former set of 537 characters are usually considered as rare characters and left out of the class space. Thus, the character-based system is incapable of modeling the large number of shapes and leads to a reduction in overall accuracy.

2.1.2 Component Based Devanagari OCR

Component based character recognition has three salient segmentation steps: Removing the shirorekha and ascenders, segmenting conjunct characters, and segmenting descender characters. The images obtained after completing the segmentation steps are matched with component classes by using classifiers. Words are generated by recombining the component results into characters. The

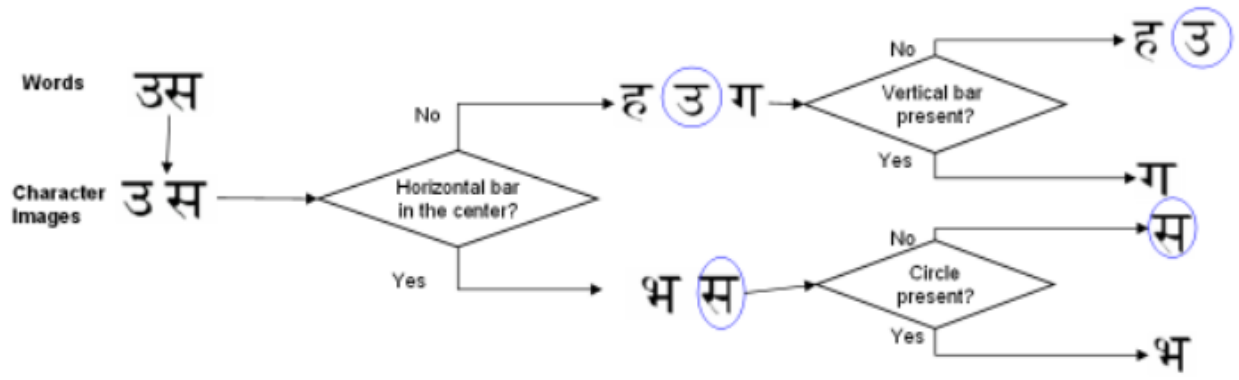


Figure 2.2: Character recognition using decision tree

horizontal projection profile discussed in the previous section is found to be stable, and is used in most Devanagari OCR implementations. After shirorekha removal, the ascenders are easily identified as image segments which lie above the shirorekha.

Bansal and Sinha perform conjunct and descender character segmentation by examining the structural features of character images [8]. They first perform shirorekha and ascender removal using horizontal projection profiles. They examine different structural features to segment conjunct and descender shapes. If the height of a character is greater than a predefined threshold, the character is assumed to contain a descender. Horizontal profile of each row is computed from such characters, and the image is segmented at a row having the lowest value of horizontal profile. Using classifiers, the lower part of the image is matched with descender classes and the upper part is matched with consonants. Images whose aspect ratio and width is greater than the corresponding thresholds are examined for conjunct character segmentation. The authors scan each such image from right to left, searching for a thin junction point. A thin junction point is defined as an image segment having height less than a predefined threshold. In case a junction is obtained, the resulting image is segmented at the column corresponding to the junction. Figure 2.3 shows examples of conjunct and descender segmentation using these methods. Ma and Doermann [64] employ the same methodology for conjunct and descender segmentation.

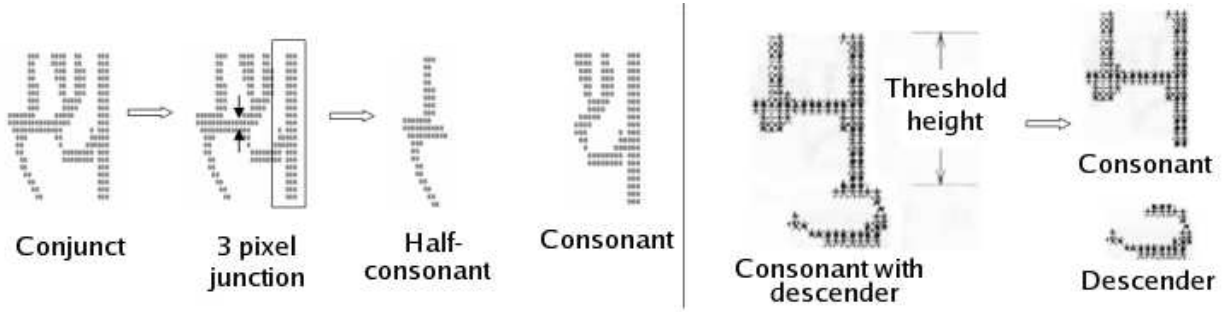
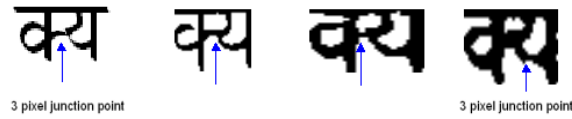


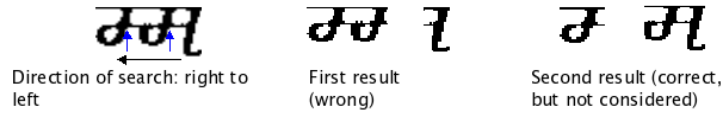
Figure 2.3: Segmenting conjuncts and descenders to obtain components [8]

Once character images are segmented, classifiers can be used to map the image segments to component classes. Bansal classifies images by using structural features in a blackboard framework [6]. In this approach, various structural features like number and position of vertical strokes, horizontal strokes, loops, or curves are used to guide the classification. Ma and Doermann [64] classify component images by template matching using generalized Hausdorff image matching. The result of component recognition is in the form of a sequence of components. Since a character is composed of many components, several characters may be formed by different groupings of components. However, all component sequences do not form valid characters and the formation of characters from components is guided by a set of rules called script composition grammar [84] (Section 2.3.3). Both Bansal and Ma use such rules in processing components and generating the final word hypothesis.

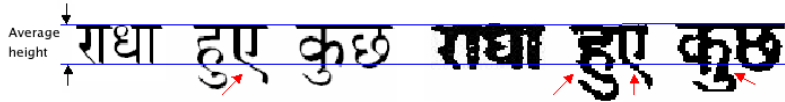
The approach used in designing a component-based Devanagari OCR has so far followed sequential application of segmentation, classification, and post-processing. Due to sequential processing, errors in each stage are propagated and the accuracy of this methodology is limited by the accuracy of segmentation. Figure 2.4(b) shows a sample where the segmentation criteria may be satisfied at two segmentation hypotheses. However, current systems are designed to consider only one segmentation result and hence make errors in such cases. Segmentation methods have so far depended on structural features which are sensitive to font and print variations. Character shapes vary across fonts, and the structural features of a specific font may not necessarily hold true



(a) Segmenting different samples of क्य which is written by combining क with य



(b) Error in segmenting म्म which is written by combining म् with म्. Blue arrows indicate possible segmentation points.



(c) Use of average height as threshold to segment characters. Red arrows indicate erroneous segmentation.

Figure 2.4: Errors due to font-dependent segmentation.

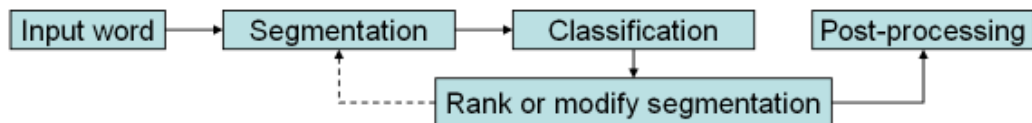
for characters of a different font. Poor print quality also adds to the unpredictability of character shapes. Figure 2.4(a) shows an example where a thin junction is found on certain samples of the conjunct क्य, but is not found on other samples. Figure 2.4(c) shows that the use of average height as a criteria for character segmentation can cause incorrect results.

2.2 Prior Techniques in Latin and CJK Recognition

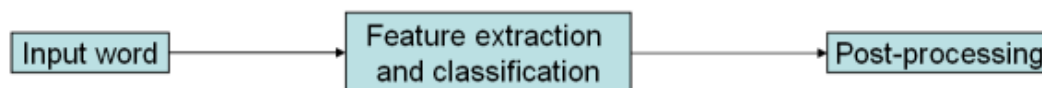
Most of the challenges in Devanagari OCR are due to the large number of character classes that have multiple components connected to each other. Latin and CJK (Chinese, Japanese, Korean) recognition systems have previously dealt with similar problems. In this section we look at both machine print and handwriting recognition approaches in Latin and CJK scripts.



(a) Dissection technique or segmentation driven technique



(b) Recognition driven technique



(c) Holistic or hybrid technique

Figure 2.5: OCR paradigms based on segmentation methodology [15]

2.2.1 Utility of the Recognition-driven Segmentation Model

Casey groups OCRs into three categories based on the segmentation methodology used [15]: “dissection techniques for segmentation, recognition driven segmentation and holistic recognition”. Figure 2.5 illustrates the major steps in these three types of techniques, and Figure 2.7 shows examples implemented along the different paradigms. Most of the Devanagari OCRs follow the first paradigm by dividing a word into smaller units and then classifying the segmented images. This sequential application of segmentation and classification is the main disadvantage of prior Devanagari OCRs. In the dissection framework, errors that occur while segmenting characters into components are carried forward and cannot be corrected by downstream processing.

Holistic techniques do not segment a word image and perform recognition using higher level features like loops, cusps, word envelops etc extracted from the entire word image [61, 66]. Words in the printed domain show more regularity in their holistic features. While these features can contribute to recognition at specific stages, additional discriminative information is essential to

recognize printed text.

Recognition driven segmentation has long been the focus of both Latin and Chinese OCR, with high accuracy rates being reported in several implementations [62, 88, 102]. In recognition driven segmentation, there is the concept of lexicon driven and lexicon-free recognition. The former design has been more popular in handwriting recognition [2, 31, 57]. In lexicon-driven systems, it is assumed that the recognizer results will be limited to a known set of words. To explain the lexicon-driven approach, we briefly describe the recognition system which has been originally reported in [31]. The recognizer first uses several training samples of each word and learns a model for each entry in the lexicon. At run time, the test word is matched with each lexicon model using feature extraction and classification methodology that has been applied during training. The result is taken as the word whose model provides the best match with the test sample. As the number of lexicon entries increase, the discriminative power required in each word model also increases. Accuracy is known to drop dramatically with increase in the lexicon size, and gets to impractically low values with lexicons having thousands of entries [40, 101].

Lexicon-driven recognizers often reduce the number of word hypotheses by using additional models or contextual information during runtime. For example, bank check recognizers limit their recognition to numerical spellings, postal address recognizers limit their recognition to city names etc. This thesis deals with Devanagari OCR where the input is only a single word image where we do not have any models to reduce the lexicon. Our effective lexicon is of the order of a few thousand words. Therefore, the best solution for our problem is a lexicon-free, recognition driven segmentation approach. In the next section, we will look at previous approaches in recognition-driven segmentation for Latin/CJK OCR and examine if they can be adapted to Devanagari OCR.

2.2.2 Recognition-driven Segmentation Techniques in Latin OCR

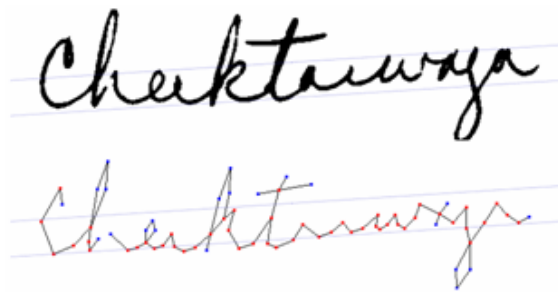
In Latin OCR, recognition-driven segmentation techniques have been used to address the problem of illegally fused alphabets [62, 88]. Song recognizes fused Latin alphabets by using a rectangular



(a) Segmentation driven technique [15]



(b) Recognition driven technique [15]



(c) Holistic or hybrid technique [106]

Figure 2.6: Examples of Latin OCR techniques which follow different OCR paradigms

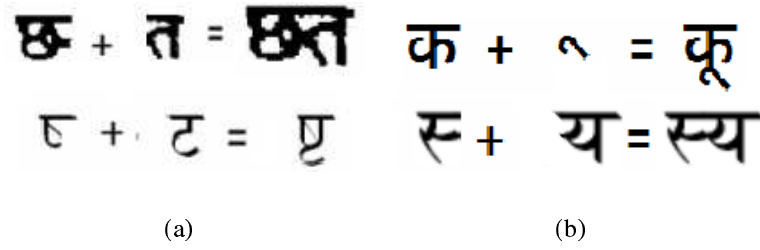


Figure 2.7: (a): Devanagari characters which need to be segmented along non-linear windows, (b) Devanagari characters which need to be segmented from top to bottom, or from left to right

window to divide the image into two parts - one that is inside the window and the remaining that is outside the window [88]. Each part is classified by using an alphabet classifier, and the classification result is rejected or accepted using a pre-defined threshold on the classifier score. When the classifier result is rejected, the right boundary of the window is increased by a small amount and a new segmentation hypothesis is considered. This method cannot segment Devanagari characters that are joined along non-linear boundaries (Figure 2.7(a)). The same authors introduce character adaptive masking to obtain a non-linear boundary. In this approach, a character hypothesis is obtained for each rectangular window. A prototype image of the specific character class is placed in the window and the pixels falling under the prototype image are set to a background color. The drawback of this method is that the selected prototypes have fixed fonts, and inconsistent boundaries are usually generated while recognizing characters with a different font.

One method to obtain non-linear segmentation boundary has been reported in Lee et. al. [62], where each pixel of an image is considered to be the node of a graph. Successive segmentation hypotheses are created by starting from the leftmost pixel in the image and selecting one node (pixel) at a time. This method is capable of segmenting characters along a non-linear boundary. While non-linear boundaries are possible in the approach developed by Lee, the non-linear boundary cannot be more than a few pixels in width. These methods are designed to handle segmentation in only one direction. However, Devanagari characters have to be segmented in at least two direc-

tions - from top to bottom and left to right (Figure ??). Similar constraints are observed in other segmentation techniques, many of which are described in the survey by Casey [15].

In recognition-driven segmentation approaches, multiple hypotheses are generated during the segmentation process. In the techniques discussed so far, consecutive hypotheses are obtained by pixel level operations on the segmentation boundaries: typically, the boundary is shifted by one or two pixels to the right. When characters are fused, the points of fusion form natural breaks and joins. Generating segmentation hypotheses at such breaks and joins is sufficient for recognition purposes. Segmentation techniques that use pixel level operations cannot accommodate such knowledge, and generate more segmentation hypotheses than required. In subsequent chapters, we design a segmentation approach that uses the Block Adjacency Graph representation to generate fewer hypotheses.

2.2.3 Techniques in CJK OCR

CJK scripts have large number of characters (3,000-7,000 characters [56, 96]), and therefore are similar to the large number of characters seen in Devanagari. Most of the CJK characters are composed of straight line primitives and recognition techniques often divide characters into pre-defined groups based on the presence or absence of a characteristic straight line primitive [82, 96]. At run-time, linear primitives extracted from a test character are used to limit the recognition hypothesis to one of these groups, and detailed features are used to classify within the group. Another technique to recognize the large character set is to model a character as a combination of linear primitives [94, 105]. During run-time, the candidate image can be classified by identifying the radicals present in it. These systems rely on identifying linear primitives from a character, whereas Devanagari characters are composed of non-linear primitives. Moreover, no technique exists to obtain non-linear primitives from Devanagari characters.

2.3 Prior Techniques in Language Modeling

In order to assist recognition, natural language phenomena is often emulated using computational models. E.g: Alphabet co-occurrence statistics have been accommodated in OCR applications by using the n-gram model [12, 44, 72]. The utility of language models need not be limited to recognition systems. An example is the stochastic finite state automata used by Forcada to design a predictive text entry scheme for reduced keyboard devices [32]. In devices like cell phones, the number of keys is less than the number of alphabets required for the language. The number of keystrokes required to enter a word using a reduced keyboard is far greater than the keystrokes needed on a normal keyboard. Forcada designed a SFSA that modeled suffix/prefix relations in the Catalan language and predicted users intended text after observing the already typed text. This reduced the keystrokes needed for text entry and created a faster interface.

Kukich provides an excellent survey of techniques that can correct word errors in text [59]. We discuss some of the currently prevalent techniques in language modeling, and outline a novel method based on modeling sound cues to prune recognition results.

2.3.1 Dictionary based models

Dictionary lookup is a language model where one creates a list of known words and matches the recognition results to a word from the list. Dictionary lookup is implicit in lexicon-driven recognition systems (Section 2.2). In case of lexicon-independent recognizers, dictionary lookup is implemented as a post-processing step [7, 81, 86, 103]. The common approach is to compare recognition results with each entry in the dictionary using a distance metric such as the Levenshtein edit distance [63]. The dictionary entry with the least distance is selected as the correct result.

It is likely that during dictionary lookup, multiple words in the dictionary are found to have the same edit distance. Constraints are imposed on the dictionary lookup process [81, 86] to disambiguate among multiple matches. For example in [81], instead of searching for a dictionary entry

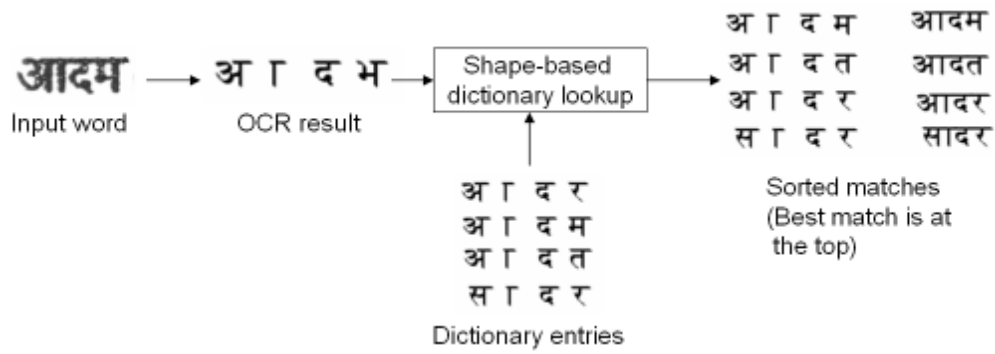


Figure 2.8: Example illustrating shape-based matching.

having the minimum Levenshtein edit distance, a shape based edit distance is used to perform the dictionary lookup. The edit distance is designed such that the matched entry has a shape that is similar to the alphabet classes provided by image based recognition. An example is illustrated in Figure 2.8. Dictionaries can often be very large in size, and the lookup process can be computationally expensive. There is considerable background work on fast dictionary lookup to reduce computational complexity [7, 103].

2.3.2 Non-dictionary Models

One view of language models is a system that can consider a sequence of alphabets as history and the provide conditional probabilities of the following sequence of alphabets. It is widely believed that these non-dictionary models can be used to enhance recognition results before applying dictionary comparison [42, 44, 59, 86]. In addition, non-dictionary models can be closely integrated with the underlying image-based recognition. Recognizers which use dictionary-based models have to construct the entire word hypothesis before taking advantage of the language model. Lexicon-driven recognizers are an exception, but are not suitable for Devanagari OCR (Section 2.2).

A non-dictionary language model is the n-gram model which can be defined on different linguistic units: alphabets [12, 44], sub-words like suffixes or prefixes [72], and words [80]. Another

model uses the Finite State Automata in accepting results as belonging to a language or rejecting them for further processing [50, 32]. Most of these non-dictionary models have considered alphabets, or words as units of the language model.

2.3.3 Language Models in Devanagari OCR

Dictionary lookup for Devanagari OCR has been implemented by Bansal in [7]. The dictionary is partitioned into groups of words, with each group identified by two alphabets, called “tags”. During recognition of a word image, the image classification part outputs a word hypothesis. Groups of two alphabets are formed from this word result. Words having the same tag as these groups are selected as the possible dictionary entries to be matched against the word result given by image classification. Information from structural features is used while matching dictionary entries with recognition results. The idea is to inhibit dictionary entries which require a substitution between characters having dissimilar shapes and structural features.

Sinha outlined a series of rules called script writing grammar or script composition rules, which specify how components combine to form alphabets and how alphabets combine to form characters [84]. All component based Devanagari OCRs [6, 38, 64] implement some form of the script writing grammar. The rules used to combine components into alphabets can vary across implementations. This is usually because specific implementations vary in the method by which they segment characters into components, or vary in the number and type of components supported. For instance if the component ˘ is not recognized in an OCR, the alphabets ओ ऐ े etc cannot be recognized. In another case, an OCR implementation might segment the vowel modifier ँ into components ˘ and ँ while another method might not perform such segmentation. In the former approach, a rule would be used to combine components ˘ and ँ into the vowel modifier ँ, whereas in the latter approach such a rule would not be needed.

2.4 Discussion

Devanagari OCR systems using character level classifiers cannot handle all the character combinations possible in Devanagari [18, 85]. Due to the paucity of character samples, such systems need to train classifiers using the limited number of samples or drop large number of characters from the class space. Both options lead to a significant decrease in accuracy. In case of component-based approaches [6, 64, 38], segmentation errors are propagated and there is no means to correct the errors once segmentation has been performed. Prior Devanagari OCR approaches have been designed for processing single-font documents and use structural features for either segmentation or classification. The accuracy of these systems is therefore limited by the accuracy of segmentation or the ability to extract structural features.

Methods used to handle the large number of CJK classes are applicable on characters having linear primitives [82, 94, 96, 105], while Devanagari is composed of several non-linear primitives. Latin recognition techniques have used approaches that can accommodate classifier scores into the segmentation process [88] and some are even capable of segmenting at non-linear boundaries [62]. However, these approaches are limited to segmentation in only one direction, and cannot scale to Devanagari characters which contain components that need to be segmented along two directions. We outline a BAG-based methodology that can segment Devanagari characters vertically or horizontally along non-linear boundaries and also show that a lower number of segmentation hypotheses can be generated by using intelligent data structures (Chapter 4, 5).

In the area of language models, Devanagari OCR approaches have so far used only script writing grammar and dictionary lookup. No effort has been made to apply non-dictionary models like n-grams which can enhance recognition results before using dictionary matching. In Chapter 6, we show that script writing grammar can be integrated with n-gram models in a stochastic framework. We also show that using syllables as units in the n-gram model can yield us better recognition results than previously used models which rely on alphabet units.

Chapter 3

Data Collection and Baseline OCR

OCR development requires large amount of training data with ground-truth to achieve acceptable levels of performance. Data should also be standardized and publicly available to compare OCR performance. The non-availability of such data sets has been a drawback for Devanagari OCR. In some scripts, simulated data has been proposed as an alternative to using real data for OCR evaluation. However, even the development of such simulated models requires established data sets and baseline OCR systems [4, 52]. Given the absence of such resources in Devanagari, we have developed tools to collect a representative data set of real Devanagari documents. Closed data sets have been used to evaluate previous OCR designs, and the implementations also are not in the public domain (Section 2.3). To make a fair comparison, we implement some of the previous approaches and create a baseline OCR system. This chapter outlines our approach in collecting public-domain Devanagari data sets, implementing a baseline OCR, and testing of the baseline OCR.

3.1 Document Image Analysis and Data Collection

Document Image Analysis (DIA) includes binarization, skew correction, text block extraction, line and word separation. This section briefly discusses DIA stages that we have implemented to semi-automatically collect a testing and training data set. The first step in our data collection is to scan Devanagari documents from contemporary books, magazines, and newspapers at 300dpi resolution with 256 grayscale levels. The documents are binarized using a neighborhood adaptive thresholding technique [49]. Skewed documents are identified and corrected by a human operator. Document smearing and horizontal projection profiles are used to separate text and graphics. Devanagari script has a distinctive header line (Shirorekha) which gives regularly spaced sharp peaks in the projection profiles. Regions of the image which generate such distinctive projection profiles are identified as text blocks. Each peak of the projection profile corresponds to the approximate location of a text line. The approximate locations are refined using a region growing approach (Figure 3.1).

The steps mentioned in this section have been combined into an interface that allows users to: (i) Open a scanned document, (ii) Correct the boundaries obtained from line and word separation, (iii) Type the truth of each word, and (iv) Save the word locations, truth and other meta-data into an XML file for future use. Our data collection tools, document images, and associated truth have been published as the ILT (Indian Language Technologies) data set [38, 58]. The documents have been manually grouped into four types. The first three types of documents contain noise and font variation (Figure 3.2(a)). Documents of type four are printed in a consistent font, and the text is relatively free of noise (Figure 3.2(b)). The complete set of collected document images, and associated truth has been made available on the internet¹.

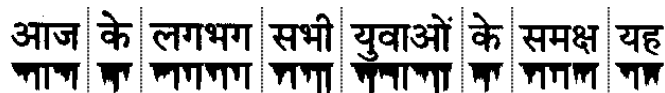
¹URL: <http://www.cedar.buffalo.edu/ilt>



(a) Text-image separation

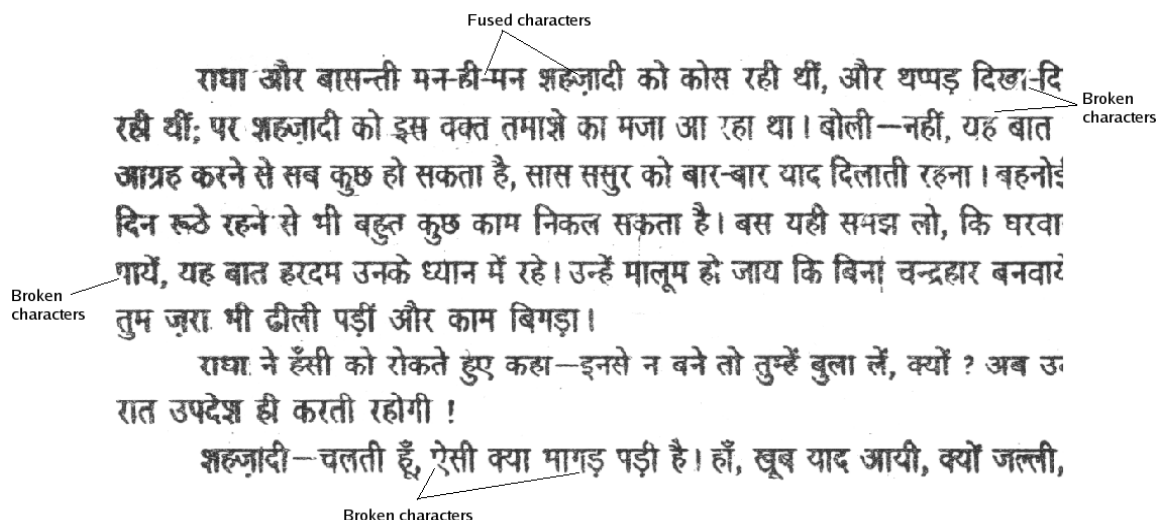


(b) Region growing technique to refine lines



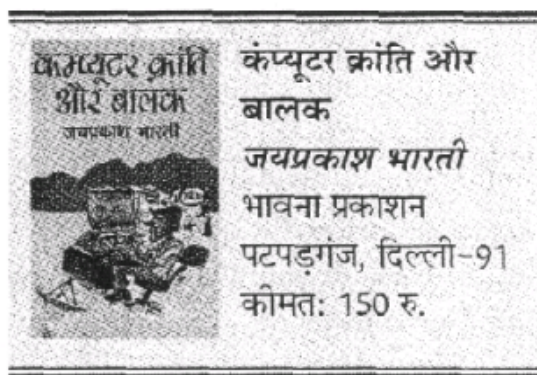
(c) Word segmentation using vertical profile

Figure 3.1: DIA stages used to obtain paragraphs, lines, and words from a Devanagari document [38].



(a) Document with several fused characters and broken letters. These can occur due to printing errors or aging in library storage. Page taken from [75]

लेखक ने उत्तर भारत के पांच राज्यों का भ्रमण करके अपने अनुभवों को सूत्रबद्ध किया है। जाहिर है, गांव-समाज जैसे विषय पर शोध पुस्तकालयों में बैठ कर या तथाकथित विशेषज्ञों से विचार मांग कर नहीं किए जा सकते। 'प्राण' ने जिन पांच राज्यों को चुना वे हैं—उत्तर प्रदेश, मध्य प्रदेश, राजस्थान, हरियाणा और पंजाब। ये शायद 'पंच' शब्द की अवधारणा के मुताबिक चुने गए 'पांच राज्य' थे; जो कि नाकाफी थे क्योंकि उत्तर भारत के हिंदीभाषी राज्यों में बिहार जैसे उर्वर प्रवेश को न शामिल करना संतुलन को बनाए नहीं रखता। राजनैतिक प्रपंच और सामाजिक वैषम्य एक अलग तथ्य है। परंतु सभी राज्यों का ग्रामीण ढांचा कमोबेश एक समान ही



■ कौशलेंद्र प्रपन्न

अपनी बालकथाओं व कविताओं के लिए मशहूर जयप्रकाश भारती ने एक अलग ही विषय पर महत्वपूर्ण

(b) Document with different fonts. Page taken from [74]

Figure 3.2: Examples of documents in the ILT dataset

3.2 Baseline OCR

Devanagari OCR techniques have been using classifiers at the character level [18] and at the component level [6, 64]. While both designs need a representative set of training samples, the first one requires training samples for character classes. Frequency analysis of the ILT image dataset shows that 15% of the characters are from classes which occur less than 25 times each. The number of character classes increases as one collects more data (Table 3.4) and sufficient training samples are not easily available for each character class. Therefore, classifiers often choose to not include infrequent characters in their class space, leading to a reduction in overall accuracy.

OCR techniques which use component classifiers follow the rationale that Devanagari characters are large in number, but these characters are composed of a relatively smaller set of components (Table 1.2). These techniques generally include the following steps: (i) Segmenting words into component images, (ii) Mapping segmented images into component classes, and (iii) Generating word results by combining the component hypotheses. Prior Devanagari OCRs which use component classifiers apply these stages sequentially on each word image (Section 2.1.2). In remaining parts of this section, we describe our design of an OCR along a similar sequential methodology. We use this implementation as a baseline system for the recognition driven design that is presented in later chapters.

3.2.1 Segmenting Words Into Components

Our first step in segmentation is to identify the shirorekha, which is a prominent horizontal line connecting the characters in a word. Approximate location of the shirorekha is identified using peaks in the horizontal projection profile of a word image (Figure 3.3(a)). Each row corresponding to the longest horizontal profile is examined for black-white transitions. A row is elected to belong to the shirorekha if the number of black-white transitions is less than a pre-defined threshold. After locating these rows, we identify and remove ascenders, which are the modifiers written above a

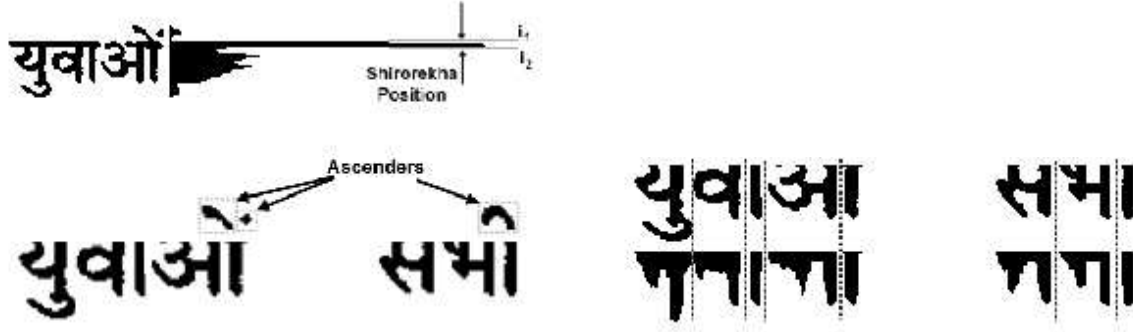
shirorekha.

Once the shirorekha and ascenders are removed, we obtain images corresponding to core or core-descender characters (Figure 3.3(b)). Characters containing descenders are identified using average character height as a threshold. Descenders are modifiers written in the lower part of a word. Therefore, if the height of a character is greater than the average height, it is selected as a potential candidate for descender segmentation. Run-length analysis is performed to identify the candidate characters which have a thin junction point (Figure 3.3(c)). If a thin junction point is not found, it is assumed that the potential candidate does not have a descender. If a junction point is obtained, the character is segmented at the junction point. The upper and lower parts are respectively selected as a core and descender hypothesis.

We have found that descenders are segmented improperly when average height is estimated using characters from a single word. We obtain some improvement in segmentation by using the average character height computed from a line of text. The use of structural features computed from a paragraph or line of text is a common practice among prior Devanagari OCR techniques [6, 64]. In the next section we discuss the feature extraction and classification techniques used to map the segmented images to corresponding classes.

3.2.2 Component Classifier

Images obtained from our segmentation stage correspond to: ascenders, descenders, core components, and conjunct characters. We have selected a class space having seven ascenders and descenders, 38 consonants, and 83 conjuncts. The classifier uses gradient features extracted from the segmented images. All images are normalized to 32×32 and the gradient magnitude and direction is computed using a Sobel operator [92]. A gradient map is computed by applying an empirical threshold on the gradient magnitudes. The gradient map is placed in a 3×3 grid and the gradient directions are quantized into 8 bins. Concatenation of the gradient vectors for the 9 cells forms a 72-element feature vector.



(a) Shirorekha detection and ascender separation

(b) Vertical profile to separate core components



(c) Top: Average height identifies 3 possible descenders, Bottom: Run length accepts the first two, and rejects the last segmentation hypothesis.

Figure 3.3: Segmenting words into components

A nearest neighbor classifier is used for the seven ascender and descender classes. Core components are pre-classified into four groups based on the presence of a vertical bar: no vertical bar (Eg: छ ट ह), vertical bar at the center (Eg: र फ क), right (Eg: व त म) or at multiple locations (Eg: क्य स स्त). Similar structural features have been previously used to perform Devanagari component classification [6, 64]. In our case, structural features are used to select a component group and four neural networks are used for classification within the group. Each network is fully connected, uses the logistic activation function, and is trained using standard back-propagation. The design of our training and test set is described in section 3.2.4. Once classification is performed, we have a sequence of class labels for each component in a word image. In the next section, we discuss the

procedure which composes these class labels into word hypotheses.

3.2.3 Composing Words from Component Results

The segmentation steps divide each character into multiple components; E.g.: के \rightarrow क + ऎ and को \rightarrow क + ऎ + र . The effect of this division is that after using the component classifier, the results for each character appear as a sequence of component class labels. Sinha defines a set of script composition rules to combine such component classes into characters [84]. These rules use positional information and the class labels of segmented images to form the character. We encode these rules in a decision tree to form characters from component results. Once characters are obtained, words are formed by placing each character one after the other from left to right. Some of the errors in segmentation or classification will lead to sequences of component class labels that violate the script composition rules. Such illegal component class labels can be rejected by our system. Dictionary lookup is performed using a lexicon of 4,291 words. The entry which gives minimum Levenshtien edit distance with classifier output is taken as the corrected result.

3.2.4 Generating Training and Test Samples

The DIA tools in section 3.1 provide word images and associated truth as part of the ILT dataset. While word images are sufficient to test the OCR system, component images are needed to train and test the classifiers. We have implemented a technique to automatically extract component images.

Using script composition rules, the word truth of each image is parsed to obtain class labels of the constituent ascender, descender and core components. The segmentation stages outlined in section 3.2.1 is used to obtain the bounding box of each component in the word image. In the absence of errors in the image segmentation stage, bounding boxes map directly to component class labels (Figure 3.4). All component images are extracted and assigned appropriate class labels.

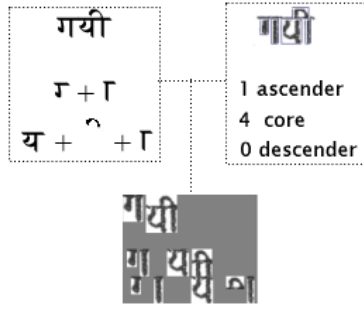


Figure 3.4: Parsing a word on query: Extracting component truth (Left), Extracting component boundaries (right), Combining the two results to obtain component images (bottom).

The images extracted for each class are manually examined to remove samples that have excessive noise and samples generated by erroneous mapping. The cleaned component truth contains on an average 400 samples of each component and is divided into a 60% test and 40% training set. The number of component samples are sufficient for training and testing purposes.

3.3 Testing the Baseline OCR

In this section, we discuss experiments on the baseline system. The stages of character segmentation, core component classification, and word recognition are analyzed using the ILT dataset.

3.3.1 Segmenting Words into Components

Our OCR outputs the number and type of components present in a word image. A difference between these numbers and the number of core, ascender and descender components in the word truth indicates a segmentation error. The difference is manually verified to distinguish between segmentation errors and classification errors. The interface used for manual verification is described in section B.2.3 (Figure B.7). Errors in any one stage are not added as errors in subsequent stages; E.g.: when a descender is not segmented, the error is not included in core component segmentation, even though core component recognition can be affected by such a segmentation error. Table 3.1

Component segmentation stage	Characters that have respective components	Number of errors	Error rate
Core	31,252	1,720	5.50%
Ascender	15,269	2,249	14.73%
Descender	1,844	1,077	58.39%

Table 3.1: Character segmentation errors

shows the percentage of segmentation errors in each stage. In what follows, we summarize the problems with the different segmentation stages:

(a) Ascender segmentation: The shirorekha has proved to be a prominent feature of each word, and is correctly identified in all words. The system is able to accurately segment the components which occur above the shirorekha, and labels such components as ascenders. The 14.73% error rate in segmentation is due to ascenders which are fused with neighboring ascenders, Eg: The character श्री looks like श्री when ascenders get fused with each other. Some ascenders like the *nukta* (Dot on a character: ण्) are small, and difficult to distinguish from noise.

(b) Core component separation: Conjoint characters are sometimes printed with a large distance between half and full consonants. This gap is frequently confused with character separation in words (Figure 1.2(b)). Poor print quality in Devanagari also creates fused characters, which are mis-recognized as core components (Figure 1.2(c)). This contributes to an error rate of 5.5% in core component segmentation.

(c) Descender segmentation: The position of descender changes with alphabet height and font. Eg: Descender ्र connects at different positions on characters क ह, to form: कु हु (Figure 1.2(a)). Due to these variations, average character height in a line often fails to identify characters having descenders. Most descender segmentation errors are due to under-segmentation, where a descender is not detected in a character. We have observed that altering the thresholds to correct under-segmentation in certain fonts leads to over-segmentation in other fonts.

3.3.2 Core Component Classification

Neural network classifiers corresponding to core components have an average accuracy of 84.77% (Table 3.2). The pre-classification stage which selects one of these neural networks has an accuracy of 86.88%. The pre-classifier relies on structural features and is not robust to font variations in the test set. Eg: ह and र correspond to the “No bar” group, but their font variations ह and र are wrongly pre-classified and sent to the “Multiple Bar” and “Middle Bar” neural networks respectively. To account for common confusions, some core components are placed in two groups, Eg: र may be recognized by both “Middle bar” and “No bar” neural networks. Despite this, accuracy reduces to 74.53% due to pre-classification error (77% decrease in case of “Multiple bar”).

Core component type (Position of vertical bar)	I (%)	II (%)
No bar	80.04	73.15
Middle	90.83	82.82
Right	75.00	72.16
Multiple	93.20	69.99
Average accuracy:	84.77	74.53
Average error rate:	15.23	25.47

Table 3.2: Core component classifier: (I): Neural network classifiers with manual pre-classification, (II): Neural network classifiers with automatic pre-classification

3.3.3 Word Recognition

The OCR recognizes 53.85% of all words using image-based approaches. Dictionary matching improves these results by 14.76%, bringing the overall recognition rate to 61.8%. For documents which contain little variation in font and print quality, the word recognition rate is as high as 74%. This compares favorably with previously reported results on single font data sets (Table 3.3). The recognition results vary across different document types, and on certain documents which have excessive font variation, the accuracy is as low as 42% (Figure 3.5). The variation of accuracy across document types reduces the overall accuracy to 54%.

OCR Type	Font restriction, Classifier	Accuracy	
		Character	Word
Segmentation driven OCR	Single font, component classifier [6]	93%	67%
	Two font, component classifier [64]	80%	87%
	Two font, character classifier [18]	90%	91%
	Multi-font, component classifier (This chapter)	84%	53%

Table 3.3: Comparison of Recognition Results

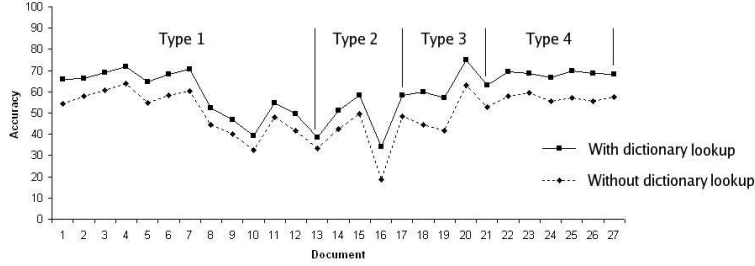


Figure 3.5: Word accuracy on different documents using the segmentation driven recognizer. Documents of type 4 are printed in a consistent font, while the other documents have variation in font and print quality. The accuracy varies not only across different document types, but is very different even in documents coming from the same type.

3.3.4 Considerations for Proposed System

Our experiments show that the shirorekha is a very prominent feature of Devanagari words, and can be correctly identified by projection profiles. Once the shirorekha is detected, ascenders can be located correctly but it is difficult to segment fused ascenders. The image information in fused ascenders is not sufficient to assist segmentation. These cases are often corrected by dictionary lookup.

For some fonts, the word accuracy from our method is 74%, while in other cases the accuracy is as low as 42%. The fall in accuracy is mainly because thresholds or heuristics are used in pre-classification and character to component segmentation. Since these structural features are hand-crafted, they work for specific fonts, but are unreliable for multi-font documents. While the accuracy of our neural network classifier is 84.77%, using structural features for pre-classification reduces accuracy to 74.53% (Table 3.2). It is found that statistics based on structural features like

average height or run length cannot be relied upon to segment all components.

An attempt has been made to design our baseline OCR using the same segmentation and component classifier techniques reported in previous designs [8, 64]. Closed data sets have been used in evaluating most of the former systems, making the comparison difficult². However, if reported word accuracies were to be compared, in some fonts we get word accuracy of 74%, far exceeding the 67% accuracy reported in [64].

Analysis of text corpora shows that 88% of all characters contain consonants/vowels or consonant-ascenders, and upto 12% contain consonant-descenders and conjuncts (Table 3.4). Structural features used to guide the segmentation criteria are not sufficient to prevent errors in 88% of the text. Sequential use of segmentation, classification, and word creation does not allow one stage to correct errors in previous stages. In subsequent chapters, we will describe techniques that can address these shortcomings.

Character Type	Example	ILT			EMILLE		
		Classes	Number of samples and their percentage		Classes	Number of samples and their percentage	
Consonant/vowel without ascender or descender	क स अ	48	14,451	44.59%	50	10,871,001	41.40%
Consonant/vowel with ascender	कि की हैं	378	14,117	43.55%	994	12,097,099	46.07%
Consonant with descender	कु सू ढ़	107	1,806	5.57%	399	12,27,331	4.67%
Conjunct without ascender or descender	क्य स्त	117	688	2.12%	793	761,112	2.90%
Conjunct with ascender	क्यों स्ति	288	1,152	3.55%	2,770	1,113,610	4.24%
Conjunct with descender	क्यू स्तु	18	38	0.12%	455	62,040	0.24%
Special symbols	ऽ : क्र	6	100	0.31%	10	65,376	0.25%
Digits	२ ३	10	25	0.08%	10	46,228	0.17%
Punctuation/Latin symbols	- : ,	11	36	0.11%	56	15,518	0.06%
Grand Total:		973	32,413	100.00%	5,537	26,259,315	100.00%

Table 3.4: Types of characters and corresponding frequencies evaluated from the ILT [58] image data set, and the nearly 1000 times larger EMILLE [5] text corpus.

²The technique reported in [64] uses scanned pages from the Oxford Hindi dictionary as test set while others use scanned copies of computer printouts [6], and photocopied documents [18].

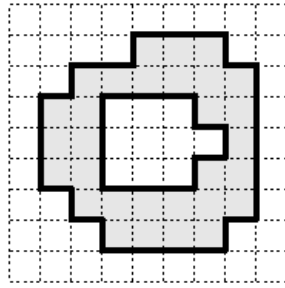
Chapter 4

Data Structures for Word Recognition

In recognition tasks, data structures can provide higher-level information regarding the image. This chapter discusses the use of higher-level image data to assist image recognition. We compare run-length representation of images with other widely used encoding schemes like bitmap, chain-code, and fax-compressed format. These representations have been previously compared on the basis of slant correction and contour smoothing algorithms [87]. We extend that work by presenting the Block Adjacency Graph (BAG) data structure and designing a BAG-based algorithm for character segmentation. We also discuss algorithms that can perform slant correction, contour smoothing, and character segmentation; all of which are important stages in word recognition.

4.1 Representation Techniques

By image representation, we refer to a data structure that can recreate the original image in a lossless manner. It should however faithfully reflect deliberate changes that may be imposed due to image transformation stages like noise removal, slant correction, smoothing etc. Hence, this study is not about representations of image characteristics or parts of images which fall under the domain of feature extraction (For example skeletons, gradient descriptors or moments) [99]. This



(a)

255	255	255	255	255	255	255	255	255	255
255	255	255	255	0	0	0	255	255	255
255	255	0	0	0	0	0	0	255	255
255	0	0	255	255	255	0	0	255	255
255	0	0	255	255	255	255	0	255	255
255	0	0	255	255	255	0	0	255	255
255	255	0	0	0	0	0	0	255	255
255	255	255	0	0	0	0	255	255	255
255	255	255	255	255	255	255	255	255	255

(b)

Figure 4.1: (a) Original image, (b) HIPS representation [87]

study is also not about lossy compression techniques that offer no benefit to word recognition. This discussion assumes that the word image in question has been binarized, and is concerned with different ways to represent a black-and-white image. In this section, we will mention some popular encoding schemes and briefly discuss their use in word recognition.

4.1.1 Bitmap Encoding

One way of representing a word image is to use an array where each entry corresponds to one pixel of the image. The value of the entry corresponds to the color, or grayscale value of the pixel. Figure 4.1 shows a sample bitmap representation where black colored pixels are represented by the value 0 and white by the value 255. Non-binarized gray scale images can be represented by having numbers between 1 and 254. This is a straightforward, though memory-intensive record of the image. The information at each pixel is imperative in performing certain operations like cross-correlation or gradient computation. Bitmaps provide a direct method to access such information. We will shortly discuss certain higher-level representations which require special procedures to provide pixel level data.

4.1.2 Chain-code Representation

Chain-codes encode the boundary information of an image as a list of exterior and interior contours. The contour can be an array of consecutive pixels, where each pixel is represented by its x and y coordinates. Another representation of the contour is a variant of the Freeman chain-code encoding [34]. Information such as the length and type of the contour, size of the bounding box may be derived from the boundary information or stored as part of the data structure.

Chain-code representation is well-suited for slant-angle detection, and has been adapted into holistic word recognition systems that do not need segmentation [65]. However, changing even a single point on the contour can cause contour breaks. Special care has to be taken while performing image manipulation like slant correction or smoothing. Segmenting an image is even more cumbersome since contours have to be split into several pieces, interior contours can become exterior contours, and at each segmentation point extra contour elements have to be added. Due to these issues, image operations which use contours require multiple copies of contour sections and have a high memory overhead.

Operations that require bitmap information cannot be easily performed using the contour. For example, the task of determining whether a pixel is inside or outside a contour is not a straightforward operation. As outlined in [73, 77], several specialized routines are required to process contours and obtain pixel level data.

4.1.3 Lossless Compression of Images

Lossless compression of images is commonly used as a representation scheme that can be transmitted over low-bandwidth connections. While these encoding schemes optimize bandwidth usage, the data structure is not designed for image processing. Consider the basic form of fax-compressed representation [46], where each row of the image is scanned and the lengths of white pixel runs and the lengths of black pixel runs are encoded using Huffman code [22]. The beginning of a new

row is represented by a special new-row symbol. This results in a very compact representation of the entire image.

These compression-oriented representations are not convenient for recognition because special functions are needed to manipulate the underlying image. For example, when run lengths are encoded using Huffman coding, it is rather complicated to change the length of a run - changing the length of a single run immediately affects the position of all the runs to the right. To obtain pixel level data, one has to reverse the compression operation, and often generate a bitmap representation. It is time-consuming to reverse the compression format for every image operation.

4.1.4 Horizontal/Vertical Run-length Encoding

Run length encoding [79] keeps track of horizontal runs (Hruns) or vertical runs (Vruns) of black pixels, storing the coordinates of the first pixel and the run length. Figure 4.2 shows the vertical and horizontal run-lengths of a Devanagari alphabet. There are 36 horizontal runs and 34 vertical runs in this image. The minimum information needed to represent images using runs is significantly less than that needed for bitmap representation. For example in case of Hruns, we only need 36 tuples of the form $\langle row\ number, column\ number, run\ length \rangle$

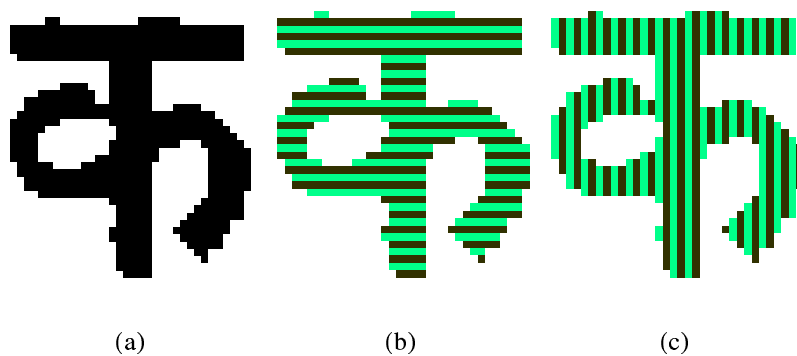


Figure 4.2: (a) Devanagari alphabet, (b) Horizontal run lengths, 36 in number, (c) Vertical run length, 34 in number.

Run-length representations allow us to derive additional image attributes that can assist image recognition. For example, one can associate neighborhood information with each run by storing the number of runs directly above or below (in case of Hruns), or the number of runs to the left and right (in case of Vruns). A concept of the *type* of each run has also been outlined in [87]. The type is an integer value which depends on the number of neighbors: for example, an isolated Hrun (a Hrun with no neighbors) has *type* 0, a Hrun with one top neighbor and no bottom neighbors has *type* 3 etc. In general, $type = 3 \times \text{number-of-top-neighbors} + \text{number-of-bottom-neighbors}$, where more than two neighbors still count as two.

Data Structure 1 Implementation of HRUN

```
struct Hrun {
    int r;                // Row of the horizontal run
    int c1, c2;           // First and last columns of
                        // the run
    int len;              // Length of the run
    struct Hrun *prev, *next; // Previous and next runs on
                        // the same row
    struct Hrun *first_abv, *last_abv; // Links to touching runs above
    struct Hrun *first_bel, *last_bel; // Links to touching runs below
};
```

To implement the run-length representation, we use Hrun and Vrun data structures similar to the ones presented in [87] (Data structures 1, 2). A linked list or array is used to store all the Hruns or Vruns of an image. These run-length representations can be constructed by only one scan of the image bitmap. The time complexity is $O(n)$, where n is the number of pixels in the image. Information about a specific pixel can be obtained by scanning through each run, and identifying the start and end pixels of the run. This is a linear time operation and can be implemented in a straightforward fashion. Adding or deleting pixels from the representation is more involved since it can create new runs or merge neighboring runs. However, changes will be required only in the runs belonging to a single column or row.

Data Structure 2 Implementation of VRUN

```
struct Vrun {  
    int col;                                // The column of the  
                                           // vertical run  
    int r1, r2;                            // First and last rows of  
                                           // the run  
    int len;                               // Length (height) of  
                                           // the run  
    struct Vrun *prev, *next;              // Previous and next runs  
                                           // on the same column  
    struct Vrun *first_left, *last_left;   // Links to touching runs  
                                           // on left  
    struct Vrun *first_right, *last_right; // Links to touching runs  
                                           // on right  
};
```

4.1.5 Block Adjacency Graph (BAG)

BAGs can be created by classifying runs as merging, splitting, or continuing runs. Adjacent runs which are classified as “continuing” are merged into a block. A graph can be constructed by connecting the centroids of neighboring blocks. Information pertaining to image structure can be inferred from edges in the graph. BAGs can be constructed using both Hruns and Vruns. For simplicity, we first discuss Hruns and then demonstrate extensions for Vrun-based BAG. The first use of BAG has been in the domain of forms processing [110], where the merging criteria creates blocks by combining adjacent horizontal runs having the same starting and ending columns. The resultant BAG is not an accurate representation of the original document, but offers significant space reduction. The underlying graph representation proves useful in identifying the semantics of different regions in a printed form.

We generate the BAG data structure in a method originally reported by Xue in [106]. Hruns are classified by the number of Hruns above (*abv*) and number of Hruns below (*bel*): in splitting Hruns $abv \leq 1$ and $bel > 1$, and in merging Hruns $abv > 1$ and $bel \leq 1$. The remaining Hruns are considered “continuing”. Blocks are constructed by combining the continuing Hruns

into a single structure. A splitting Hrun is added to the block that occurs above it, and a merging Hrun is added to the block that occurs below it. The motivation is that this merging criteria follows natural breaks and joins in the image; hence it can segment a written symbol into its constituent strokes. This intuition often proves to be correct as the resulting graph follows character shape, simulating the way a human might trace the constituent strokes of an alphabet/word. For example, the character क contains a loop on the left side and a half circle on the right. These features are also seen in the BAG in figure 4.3(c).

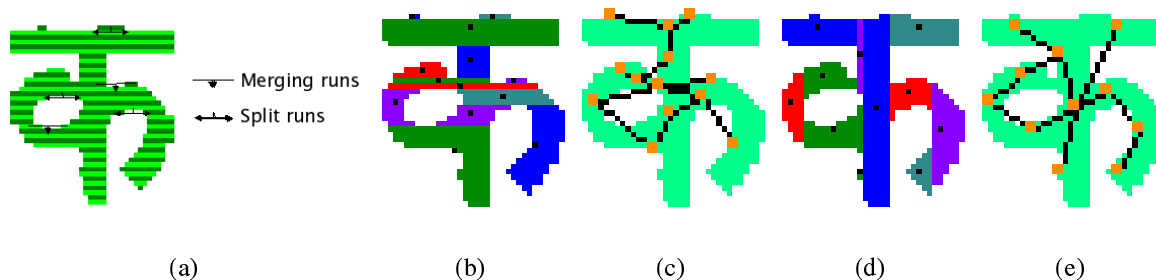


Figure 4.3: Obtaining block adjacency graph: (a,b,c) Using horizontal runs (Hruns) to create BAG. (d,e) Using vertical runs to create BAG. The BAG is seen in (c,e) as a set of nodes connected by black edges.

The process of computing the BAG requires only one visit to each Hrun. The computation has time complexity of $O(n)$, where n is the number of Hruns. The original Hrun data structure has been extended to create a representation for each block in the BAG (Data structure 3). The BAG is represented as an array or a linked list of blocks.

4.2 Image Processing Using Run-length Representation

To convert bi-tonal images into BAG encoding using Hruns, each row of the image is scanned from left to right. Sequences of black pixels separated by white pixels are identified as the Hruns of that row. BAG encoding proceeds from the Hruns in the method outlined in section 4.1.5. Extracting

Data Structure 3 Implementation of each block in the BAG

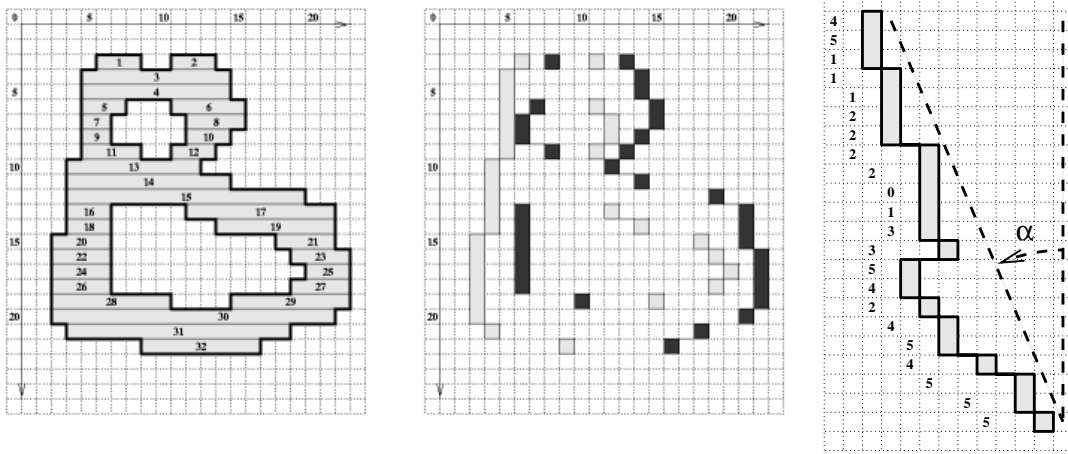
```
struct Block {
    struct Hrun *first_run;           // First Hrun in the block
    struct Hrun *last_run;           // Last Hrun in the block
    int height;                       // Height of the block
    int area;                         // Area of the block
    int x, y;                         // Centroid of the block
    int in, out;                      // In degrees and out degrees
    struct Block *above;              // Linked list or array of
                                    // blocks above
    struct Block *below;              // Linked list or array of
                                    // blocks below
}
```

Vruns and the corresponding BAG is similar, except that the image is examined one column at a time and the scan direction is top to bottom. We now demonstrate various image processing stages that can be implemented by taking advantage of the Hrun and BAG representations. Characters are used as examples for some of the stages discussed below. These stages can be automatically extended to word images without any change in implementation.

4.2.1 Slant Angle Detection and Correction

In this section, we describe the slant angle correction using Hruns that was originally reported in [87]. Firstly, *left contour* of the image is defined as a collection of left endpoints of Hruns. A *right contour* is similarly defined using the right endpoints of Hruns. Figure 4.4(a) shows an image of the letter “B” represented by horizontal runs and its 7 left and 7 right contours. For each contour point, its’ two upper and two lower contour points are selected as context. The context now contains five contour points. Four lines are formed by considering pairwise combinations of the consecutive points. The slope of each line is computed, the 4 slopes are summed and assigned as slope s of the point in consideration. A contour point is said to have near-vertical slope if the $s \leq thresh$, where *thresh* is a pre-computed threshold. When a left or a right contour has at

least two consecutive points with near-vertical slopes, the contour is selected as a *near-vertical line* (Figure 4.4(b)).



(a) Horizontal run-length representation and left (in gray) and right (in black)

(b) Example of a near-vertical line

Figure 4.4: Slant angle detection using run-lengths [87].

The clockwise angle made by a line connecting the end points of a near-vertical line with the vertical axis is assigned as the *slant angle* α of the near-vertical line. The slant angle of each near-vertical lines is multiplied with its height to obtain a weighted slant angle. Average of the weighted slant angles is considered as the average slant angle of the image. Image slant is now corrected by taking the slope of near-vertical lines closer to 0. This is achieved by shifting all runs in the horizontal direction by a number of pixels depending on the average slant angle and the row of the run. During slant correction, the offset of all runs of a particular row is the same, and the offset needs to be calculated only once for all Hruns of a particular row. Additionally, slant correction can be done by only updating the value of $c1$ and $c2$ in the data structure of each Hrun.

4.2.2 Contour Smoothing Using BAG

Contour smoothing is a process by which noisy pixels on the boundary are removed. This process can also be extended to interior contours, by filling up small holes and smoothing irregularities on convex surfaces. Previous data structures support morphological operations like averaging, masking, opening etc [57, 87]. Morphological operations are applied uniformly over the entire image and can cause undesirable changes in character structure. These operations typically use a structuring element to model the noise which needs to be removed from a character. However, the likely structure of a character is not modeled in these operations. The methods cannot easily distinguish between undesirable noise and the thin junctions and holes which are a natural part of character structure (Figure 4.5).



Figure 4.5: Morphological opening

We design an alternate approach where a pixel is modified not only by using structuring elements, but by also considering its' connectivity to other parts of the character. As described in section 4.1.5, each pixel of an image can be mapped to a block in the BAG structure. The connectivity of this pixel is represented using the edges of the BAG structure. If the corresponding block has multiple edges, it indicates that the block is connected to many other blocks of the graph. This usually means that the pixel is part of a joint that connects two or more parts of the image (Figure 4.6(b)).

In our implementation, we compute a average stroke width from the Hruns: $sw = (\text{sum of lengths of each run})/(\text{number of runs})$. The average size of blocks in BAG representation is calculated as: $a = (\text{sum of areas of each block})/(\text{number of blocks})$. Blocks which have an area smaller than a and are connected by only one edge are flagged as noise and removed. This removal can

be accomplished by deleting the block from the data structure, a constant time operation. Figure 4.6 shows how noisy blocks can be removed. Since BAG construction represents merging and splitting runs, thin sections at the trailing ends of a character are grouped into a single large block (Figure 4.6(b.4)). Morphological operations cannot naturally come up with this association and ad-hoc restrictions have to be programmed to prevent such regions from being eroded.

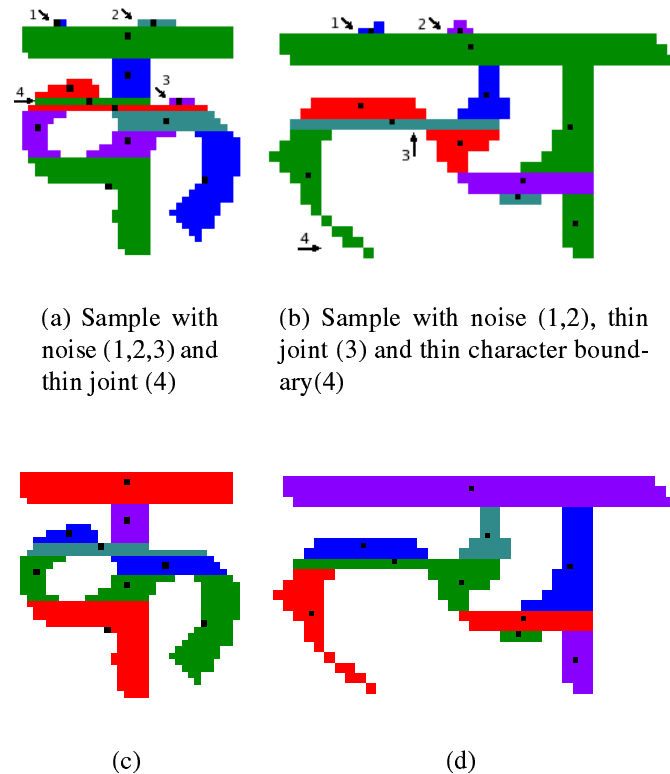


Figure 4.6: Image smoothing using BAG: (a,b) BAG with small blocks, (c,d) Spurious blocks are removed. Narrow joins and small blocks are not deleted, but merged with larger blocks.

4.2.3 Baseline Estimation

Baseline is a line that provides an approximate alignment of characters as they are written in a word. Sometimes, more than one baselines are derived from a word: Upper baselines, middle baselines and lower baseline. A typical approach of baseline detection is to consider all the lower

extremities in a word and use linear regression to construct a straight line that best fits these points. Lower extremities of a word can be obtained from Hrun representation by collecting all runs that have no neighbors below. In contrast, identifying lower extremities using a bitmap or contour representation is computationally more intensive. However, the baselines obtained using the lower extremities of a word can be thrown off by artifacts. Generated baselines tend to be lower than expected in several words and the lowering of baselines is more in words with multiple descenders.

A more natural baseline can be created by incorporating the concept of the “mass” of each character. Such a baseline would be closer to the center of mass of each character in a word and at the same time bend marginally towards descenders. It has been shown in [87] that such baselines can be designed by using Hrun representation. In the approach, an initial centerline is obtained by linear regression through all pixels in the word. Hruns that lie below the initial centerline are then selected. Linear regression is performed on the selected Hruns to obtain a final baseline estimate.

The Hrun based approach fails to obtain correct baselines in the case of short words with multiple descenders. In such words, the centerline computed as the first step tends to move towards the descenders. This results in several Hruns not being selected while computing the baseline. To counter the subsequent inaccuracy in baseline estimation, the authors in [87] propose different methods for short and long words. We use the BAG representation to eliminate the need of a separate procedures for words of different length. We first compute the BAG of a word and estimate the center of mass of each block in the BAG. Linear regression over the centers of mass is performed to get a center line. Blocks which fall below the centerline are selected. Linear regression is performed on the centroid of these blocks, and the result is taken as final estimate of the baseline. The approach in [87] performed linear regression on Hruns, while we perform regression on the centroids of blocks. The advantage of our approach is that the centroids of blocks are similar to the centroids of strokes in a character. Effectively, our technique performs regression on the centroids of those strokes which lie at the bottom of a word, or which extend from the bottom of a word

to the center. A reliable baseline is found irrespective of whether the word is short or long. Our approach also uses a lesser number of points for linear regression.

4.2.4 Character Segmentation

Word recognition systems often use segmenting modules that divide words into smaller parts. One approach in these segmentation algorithms is that the segments generated from a word should closely represent legal characters [57, 62]. In some other approaches, the results of segmentation correspond to pieces of characters [31]. The underlying motivation is that modeling the shapes of characters and parts of characters is a simpler task in comparison with modeling entire word shapes. Recognition can be performed by first classifying the characters or parts of characters and then constructing the final word from classifier results.

Previously, OCR researchers have performed segmentation using sliding window techniques. As discussed in section 2.2, window based segmentation techniques cannot account for non-linear character boundaries and are able to segment only in one direction. We design a BAG-based segmentation technique that can incorporate non-linear character boundaries and perform segmentation in at least two directions - top to bottom and left to right. We also show that the BAG based technique gives a lesser number of segmentation hypotheses in comparison with previous methods.

The intended results of segmentation are images that correspond to characters or pieces of characters. We believe that the following conditions can restrict segmentation hypotheses to images that represent characters: (i) Boundaries of the segmented images should occur at natural breaks and joins, and should not follow pre-defined cuts, (ii) Each image result from segmentation should have only one connected component. We outline a BAG based approach that can incorporate these restrictions into character segmentation.

The first restriction is met by building a BAG using the Hrun merging condition presented in section 4.1.5. The merging criteria ensures that each block in the BAG follows natural breaks and joins in a character. At least one subset of blocks can be selected in a way that the image formed

by combining the selected blocks closely represents one of the constituent characters. Figure 4.7 shows an example of how the BAG of a fused character can be broken into sub-graphs that are very similar to BAGs of the individual characters.

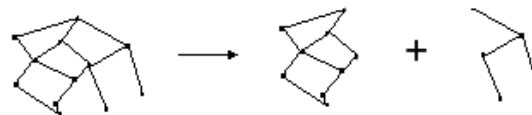
We impose the second condition as a restriction on selecting nodes which form the segmented images: “The segmentation hypotheses has to be a subgraph selected from the BAG such that none of the selected nodes are dis-connected”. In the absence of this condition, multiple character segmentation hypotheses can be obtained by selecting arbitrary combinations of nodes from the BAG. Arbitrary selection of subgraphs from any BAG can generate a prohibitively large number of hypotheses. This condition ensures that the number of segmentation hypotheses is never more than the number of blocks in the BAG.



(a) Individual characters and their BAG



(b) BAG of fused character



(c) Graph of the fused character contains graphs that are similar to the graphs of individual characters.

Figure 4.7: The use of BAG in character segmentation

Representations like bitmaps and chain codes can support only window based segmentation. Other techniques use bitmap representations to segment characters by applying a neural network for selective segmentation [35, 109]. Outputs of the neural network determine regions of the image which correspond to characters and those that do not. This method requires training of a neural network or some other classifier to obtain the segmentation hypothesis. The BAG technique can perform character segmentation without the use of such classifiers.

4.3 Using BAG in OCR

Our comparison of the different representation techniques shows that BAG is capable of assisting in the following stages of image recognition: slant correction, contour smoothing, baseline estimation, and character segmentation. The BAG structure can be built in linear time and most image operations on the BAG have linear time complexity. These block units reduce the run-time needed for operations like baseline estimation that would otherwise operate on each pixel. BAG allows us to use various merging criteria to form such blocks, enabling different applications like form dropout [110], feature extraction and character modeling [106], and character segmentation (This thesis). In the subsequent chapter, we discuss how BAG representation can be used to design a Devanagari OCR based on the recognition driven segmentation paradigm.

Chapter 5

Recognition Driven Segmentation of Devanagari Characters

Each character in a Devanagari word can contain different numbers and types of components. A component is called a core, ascender, or descender depending on its position within a word (Figure 1.1). For example, in the word पूछती; प, छ, त, and र are core components. ू is a descender and ी is an ascender. Devanagari characters can be of different types depending on the components present in it. Conjunct characters contain core components fused from left to right, E.g.: क्य is made of क and य. Consonant-descenders contain a core and a descender fused from top to bottom, E.g.: कु is made of क and ू. Some characters contain a core and an ascender, E.g.: के is made of क and ी. In a word, characters are placed from left to right and connected on the top by a horizontal line called the shirorekha.

As discussed in chapter 2, Devanagari OCRs may be designed by segmenting a word into characters, or components. Segmenting a word into characters is easier than segmenting into components. Character classes are large in number and components are relatively fewer. Therefore, the task of classifying all possible character shapes is relatively difficult when compared to the classification of components (Section 2.1.1). In this chapter, we will use component classifiers and

develop a new recognition driven methodology to segment Devanagari characters.

Prior Devanagari OCRs segment characters into components by using font-dependent heuristics, and apply classifiers on the segmented images. We find that while heuristics can accurately segment the shirorekha and identify the ascenders, errors are caused while segmenting conjuncts and descenders. The stages of segmentation and classification are applied sequentially, and there is no feedback process to correct segmentation using classifier results or a language model. Segmentation errors are carried forward and cannot be corrected by downstream processing (Drawbacks outlined in section 3.2). This cascading of errors has been addressed by Latin and CJK OCRs designed along the recognition-driven segmentation paradigm [62, 88, 102]. However, in section 2.2, we have seen that Latin/CJK OCR approaches are designed to segment characters using linear windows, or perform segmentation in only one direction (left to right). Devanagari characters have to be segmented from left to right and top to bottom, and contain characters which have to be segmented along non-linear boundaries. We will show that BAG representation can be used to segment Devanagari characters both horizontally or vertically and accommodate non-linear segmentation boundaries. Our model can also combine segmentation hypotheses with classifier scores and n-gram statistics to generate word results.

In our methodology, recognition of Devanagari words is divided into three tasks: (1) Segment characters along non-linear boundaries from top to bottom or left to right, (2) Use classifier scores to rank the segmentation hypotheses, (3) Combine the different classifier results into a word. An overview of our methodology is shown in figure 5.1. Stages (a,b,i,h) of the methodology contribute to task (1) and are outlined in section 5.1. Stages (c,d,e,f,g,j) cover task (2) and are discussed in section 5.2, and stage (k) which covers task (3) is outlined in chapter 6.

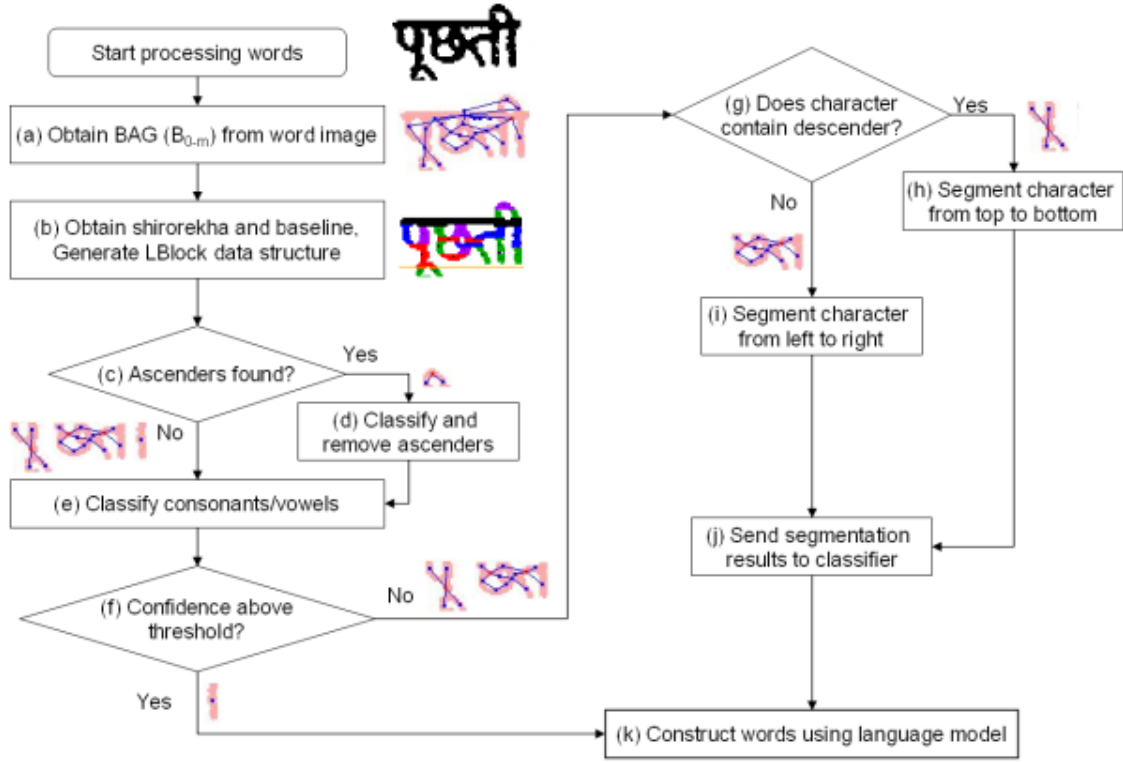


Figure 5.1: Recognition driven segmentation

5.1 Segmentation by Block Adjacency Graphs

One approach to segment Devanagari characters is to divide them along natural breaks and select adjacent, non-linear image segments. Block Adjacency Graph (BAG) is a common technique used to represent adjacency information needed for this process [89]. Previous implementations which use graphs include Chinese and Latin OCRs which extract features like loops, cusps, joints from the BAG and use the features to classify the graph [78, 106, 110, 112]. Variations of graphs have been used for image manipulation tasks like thinning [67, 95]. While these approaches indicate the utility of BAGs, they do not employ BAGs for segmentation. In our approach, we use graph representation to segment characters into images of components. Classification is done using features extracted from these images.

5.1.1 Construction of BAG

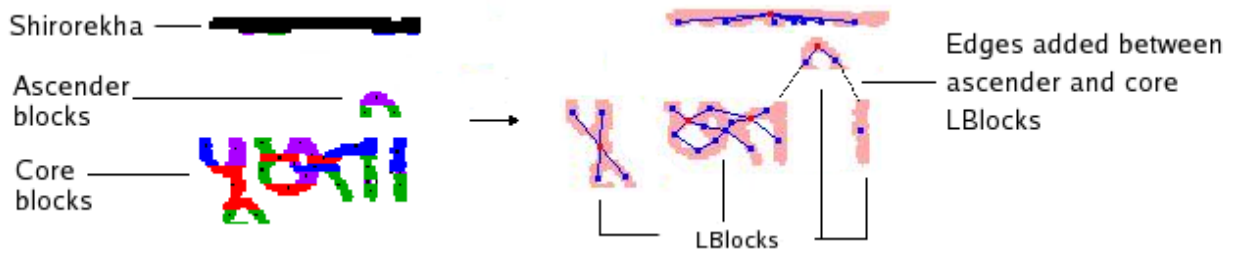
We examine word images along horizontal runs (Hruns) to build a Block Adjacency Graph (BAG) (Section 4.1.5). The process first creates a Hrun representation from the word image. Hruns having more than one adjacent run at the top or bottom are marked as merging or splitting Hruns. Blocks are constructed by grouping Hruns that are separated by such merging or splitting Hruns. A graph can be constructed by the following two steps: (i) The centroid of each block is stored as a node, (ii) Edges are added between centroids of neighboring blocks. The BAG thus represents the word image as a connected set of discrete nodes which follow the natural character/word contour (Figure 5.2(a)). The blocks are stored in an array *block_array*. A few analytical operations are now performed on the BAG to obtain positional information: average stroke width of the word *sw* is computed using the Hruns (Section 4.2.2), and baselines are computed by linear regression over the centroids of each block in the BAG (Section 4.2.3). The longest horizontal block in the BAG is identified as shirorekha.

LBlock Data Structure:

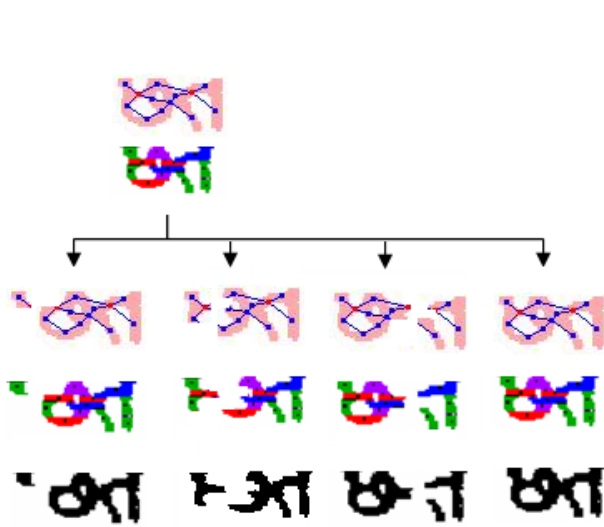
At this stage, the shirorekha has been separated from the word. We compare the position of each of the remaining blocks with the shirorekha. Blocks that are above shirorekha are flagged as belonging to ascenders and those that are below are flagged as belonging to core components. A component can contain multiple blocks, E.g.: The ascender in figure 5.2(b) appears as three blocks connected by edges. We create a data structure called LBlock (Linked Block, Data structure 4) that groups together such connected blocks as a single component. All blocks of the BAG are stored in an array and edges on each block are traversed using breadth first search (BFS [22]) to construct Lblocks (Procedures 5, and 6). Edges are added between LBlocks to indicate connections between a core component and an ascender component. Figure 5.2(b) shows the construction of four LBlocks from the BAG of a word.



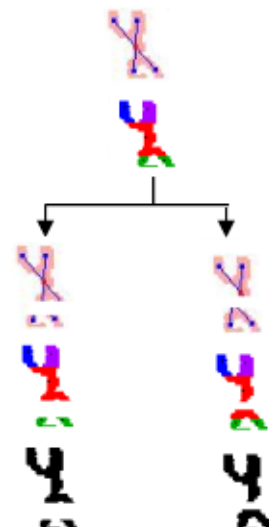
(a) Obtaining BAG from a word. Left: Original word, Center: Blocks in the BAG, baseline and shirorekha detection, Left: Graph representation of the word



(b) Generating LBlock data structure



(c) Segmenting a character horizontally by selecting nodes from left to right



(d) Segmenting a character vertically by selecting nodes from bottom to top

Figure 5.2: Using BAG for recognition

Data Structure 4 Implementation of Linked Blocks in the BAG

```
struct LBlock {
    Block* first_block, last_block;    // Pointers to the first
                                        // and last block
    BlockList* blocks;                // Link list of blocks
    int left, top, right, bottom;      // Bounding box
    int area;                          // Area in pixels
    int pos;                           // Set to 0, if it is
                                        // above shirorekha,
                                        // set to 1 if below.
    int tag, visited;                 // Data to represent
                                        // processing information
};
struct BlockList{
    Block* current, next;              // Pointers to manage a link list
}
```

5.1.2 Character Segmentation Using BAG

Each LBlock is a BAG and we generate multiple segmentation hypotheses by selecting different combinations of blocks from the BAG. Ascender LBlocks are not segmented. A naive approach is to exhaustively select subgraphs from the BAG. The number of hypotheses generated by exhaustive selection is exponential in terms of the number of nodes, and will be prohibitively expensive. We prevent an exhaustive enumeration by taking advantage of the fact that segmentation hypotheses correspond to components, and that component shapes are restricted by two conditions: (i) Components are mostly joined from left to right or top to bottom, (ii) Devanagari components cannot have disconnected images. The equivalent conditions for subgraph selection are: (i) A node cannot be selected if it is to the right (bottom for descender characters) of a node that has not been selected, (ii) Subgraphs represent components and cannot have dis-connected nodes.

The conditions on subgraph selection ensure that the number of segmentation hypotheses is never more than the number of blocks in the BAG. The hypotheses are also limited to subgraphs which structurally represent Devanagari components. Figure 5.2(c) shows an example where

Procedure 5 Generate linked blocks: create_LBlocks(struct Block* block_array[], struct LBlock* lblock_array[])

```
1: for Each block  $b$  in block_array do
2:   Set  $b.visited$  to false
3:   if  $b$  is on the shirorekha then
4:     Create an LBlock for the shirorekha store it in  $cb$ 
5:     Set  $b.tag = 0$ 
6:   end if
7:   if  $b$  is above the shirorekha then
8:     Set  $b.tag = 1$ 
9:   end if
10:  if  $b$  is below the shirorekha then
11:    Set  $b.tag = 2$ 
12:  end if
13: end for
14: for Each block  $b$  in block_array do
15:  if  $b.visited$  is set to true ||  $b.tag == 0$  then
16:    Set  $b$  to the next block
17:    Goto step 2
18:  end if
19:  Create a new LBlock and store it in  $cb$ 
20:  Set  $cb.first\_block$  to  $b$ 
21:  Set  $b.visited$  to true
22:  Call bfs( $cb, b$ ) // Procedure 6
23:  Add  $cb$  to lblock_array
24: end for
```

Procedure 6 Breadth first search: bfs(struct LBlock* cb , struct Block* b)

```
1: for Each edge  $e$  of block  $b$  do
2:   Set  $\hat{b}$  to block at the end of  $e$ 
3:   if  $\hat{b}.visited$  is set to true ||  $\hat{b}.tag \neq b.tag$  then
4:     Set  $e$  to the next edge
5:     if  $e$  is null then
6:       exit // All edges have been traversed
7:     end if
8:     Goto step 2
9:   end if
10:  Set  $\hat{b}.visited$  to true
11:  Add  $\hat{b}$  to the linked list of  $cb$ 
12:  Call bfs( $cb, \hat{b}$ )
13: end for
```

blocks are selected from left to right to segment a character horizontally. Figure 5.2(d) shows an example where blocks have been selected from top to bottom to segment a character vertically. Since the BAG data structure is consistent with breaks in the character image, we are able to segment Devanagari characters along non-linear boundaries. In contrast, segmentation techniques have so far used sliding windows (Section 2.2.2). Such methods can only segment an image from left to right. In comparison with our approach, window-based techniques generate a larger number of segmentation hypotheses. BAG computation and analytical operations discussed here use only the word image and do not depend on font information. As discussed in section 2.1.2, prior Devanagari OCR methods use external font information to perform segmentation.

5.1.3 Connecting Broken Junctions in the BAG

During segmentation, it is assumed that LBlocks generated in the section 5.1.1 correspond to ascenders, cores or fused components. However, noise in printed text can cause broken characters, which result in LBlocks that do not correspond to such components. Figure 5.3(b) shows that five LBlocks are obtained from a sample of the word सुख. This word has only two core characters सु and ख, and should contain only two LBlocks. In this section, we describe a method that is used to correct this aberration. This method is performed on the BAG of a word image before performing the segmentation explained in section 5.1.2.

In a correctly printed document that is free of image degradation, there would be a thin junction between the images corresponding to LBlocks 2 and 3 (Figure 5.3(b)). Similarly, images corresponding to LBlocks 3 and 4 would have two junctions. This would have given only two LBlocks for our sample word. In BAGs, missing junctions appear in the form of sharp edges that are close to each other but belonging to different LBlocks. We use chain code processing to identify and connect such nearly-contiguous sharp edges in neighboring LBlocks. We first obtain maximum-curvature points from the BAG. Maximum-curvature point is a location on the outer contour where curvature changes by more than 90 degrees. Chain codes are used to examine the outer contour of

each LBlock and determine points of maximum-curvature. Many maximum-curvature points are seen in regions of LBlocks that have sharp edges. These regions are usually at the ends of disconnected junctions. If two adjacent LBlocks contain maximum-curvature points that are closer to each other by less than twice the stroke width, we add an edge between the two LBlocks. The process is outlined in procedure 7.

Procedure 7 Connect broken LBlocks: `refine_LBlocks(struct LBlock* lblock_array[], int a)`

```

1: Sort the LBlocks in lblock_array from left to right based on their centroid.
2: for Each pair of neighboring LBlocks  $Lb_1$  and  $Lb_2$  in lblock_array do
3:   Compute the maximum-curvature points  $mc_1$  for  $Lb_1$ 
4:   Compute the maximum-curvature points  $mc_2$  for  $Lb_2$ 
5:   for Each point  $m_1$  from  $mc_1$  do
6:     for Each point  $m_2$  from  $mc_2$  do
7:       if  $dist(m_1, m_2) < a$  then
8:         Introduce an edge from a block in  $Lb_1$  to a block in  $Lb_2$  { $a$ : Average stroke width}
9:         Goto line 2
10:      end if
11:    end for
12:  end for
13: end for
14: Copy all the blocks in lblock_array to blocks_array
15: Delete all entries from lblock_array
16: Call create_lblocks(lblock_array, blocks_array) {Procedure 5}

```

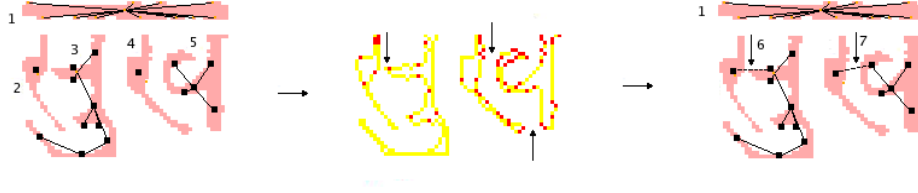
LBlocks corresponding to ascenders are omitted from the refinement process. The limitation with ascenders is that they are very small in size and there is little image information to analyze ascender shape.

5.2 Using Component Classifiers to Guide Segmentation

The steps described till now construct LBlocks from a word image and segment LBlocks to obtain ascender, descender and core components. In this section, we discuss the design of classifiers that can map the segmentation results to component class labels. We design three classifiers: (i) An ascender classifier having 7 classes, (ii) A core classifier with 40 classes, and (iii) A descender



(a) Obtaining BAG from a word: Left: Original word, Center: Blocks in the BAG, baseline and shirorekha detection, Left: Graph representation of the word



(b) Connecting broken characters: Blocks 2-3 and 4-5 are from broken characters and can be connected by using chain code (arrows) to form connected LBlocks 6 and 7

Figure 5.3: Creating the BAG data structure for a Devanagari word

classifier with 5 classes (Table 5.1). In sections 5.2.1 and 5.2.2, we discuss feature extraction and classification techniques employed in our design.

So far, we have assumed that apart from ascender LBlocks, each LBlock in a word is segmented using the BAG. Multiple segmentation hypotheses are generated from such segmentation, and each hypotheses is mapped to a component class label. In practice, segmentation is not necessary in ascender LBlocks, and in core LBlocks which contain only one component. Such LBlocks correspond to Devanagari consonants, vowels, and consonant-ascenders which form 88% of Devanagari text (Table 3.4). By avoiding the segmentation of corresponding LBlocks, we can save on execution time, and more importantly, prevent errors due to oversegmentation. Ascender LBlocks can be easily identified using the distinctive shirorekha, and can be excluded from segmentation. It is comparatively difficult to predict if a core LBlock contains a single component or multiple components. We use classifier confidence to make this prediction. In section 5.2.3, we will discuss how equal error rate is used to obtain threshold scores for the core classifier.

।	च	थ	र
अ	छ	व	ल
इ	ज	ध	ळ
उ	झ	न	व
ऊ	ट	प	श
ए	ठ	फ	ष
क	ड	ब	स
ख	ढ	भ	ह
ग	ण	म	क्ष
घ	त	य	त्र

Core components

ॠ ॡ ॢ ॣ । ॥ ० १ २ ३ ४ ५ ६ ७ ८ ९

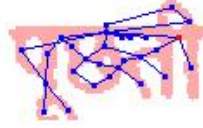
Ascenders Descenders

Table 5.1: Devanagari components recognized by our system

Typically, a low classifier score indicates that the input image is a poor match to the class space. An LBlock composed of multiple components is likely to have a shape that matches poorly with any single core component. When such an LBlock is input to our core classifier, the output is the closest matching class label with a low score. We compare the output score of each LBlock with the threshold score of the corresponding class label. If the output score is higher than the threshold, we consider that the LBlock contains a single component, and retain the top two classifier results. If the output score is lower than the threshold score, we consider that the LBlock has more than one component; the topmost classifier result is retained, and the LBlock is flagged for segmentation.

When an LBlock is selected for segmentation, structural features are examined to choose one of two segmentation hypothesis: (a) *Consonant-descenders*: If at least 10% of the blocks are below the baseline, it is considered as a consonant-descender character, otherwise it is considered to be a conjunct or fused character. Adjacent blocks are selected from top and bottom, and sent to core and descender classifiers respectively, (b) *Conjuncts/fused characters*: Adjacent blocks are selected from left and right, and are sent to core classifier.

In subsequent sections, we discuss the design of feature extraction techniques and classifiers for Devanagari components. We analyze receiver operator characteristics and obtain threshold scores for the core classifier.



(a) Word image and the BAG

		र 0.16 ~ 0.33			र 0.46 झ 0.33
		त 0.28 ~ 0.46			र 0.04 झ 0.32
		प 0.52 ~ 0.30			र 0.36 ढ 0.32
		प 0.49 ~ 0.33			छ 0.47 त 0.50
		प 0.27			छ 0.51 त 0.43
					झ 0.18

(b) Left: BAG nodes selected as segmentation hypotheses, Middle: Component images, Left: Classifier result and score

(c) Left: BAG nodes selected as segmentation hypotheses, Middle: Component images, Left: Classifier result and score

र	र झ		
प	छ ढ	~	र
त	त		
प	झ		

(d) Lattice of classifier hypotheses. Results in circles are the correct ones

Figure 5.4: Results of recognition driven segmentation using Block Adjacency Graph.

5.2.1 Feature Extraction Techniques

Feature extraction is the method of processing images to obtain data that can capture the general shape and structure of printed symbols, and at the same time ignore subtleties like font variations or noise. An introduction to different feature extraction schemes can be found in the survey papers by Connel [21] and Jain [99]. We experimented with two feature extraction techniques for Devanagari OCR:

Gradient Features: Images are normalized and a Sobel operator is used to measure the magnitude and direction of intensity changes in a small neighborhood of each pixel. A gradient map is generated by retaining gradient magnitudes which are above a threshold. The gradient map is put into a 3 X 3 grid and gradient directions are quantized into 8 bins: horizontal, vertical, and the two diagonals in each direction. Each cell contributes 8 real-valued features to the feature vector. Concatenation of the gradient vectors from 9 cells forms a 72 element real-valued feature vector that represents gradient variation in the image.

Gradient, Structural and Concavity Features: The GSC feature set, originally reported in [30] encompasses three general attributes that are hierarchical in scale and orthogonal to each other. It is a 512 binary feature vector consisting of: (a) 192 gradient features, which represent curvature information of local edges in the image, (b) 192 structural features, which represent short strokes of different angles, and (c) 128 concavity features which reflect holes and concavities pointed towards left, right, up or down.

5.2.2 Classifier Design

Two neural networks and a K-nearest neighbor classifier have been studied for Devanagari. We use three-layer, fully connected neural networks with logistic activation function. The neural networks have been trained using standard back-propagation.

Neural Network Classifier using Gradient features: This classifier has 72 input nodes, 48 hidden nodes and 40 output nodes.

Neural Network Classifier using GSC features: This classifier has 512 input nodes, 128 hidden nodes and 40 output nodes.

K-Nearest Neighbor using GSC features: A K-nearest neighbor classifier has been designed using the GSC feature set and the Kulzinsky metric. This has been used earlier for Latin handwriting recognition [111]. Given a test image, the distances of its “K” nearest neighbors are obtained. Different values of K were examined before selecting K=3, which gave a top choice accuracy of 94.27% (Experiments outlined in chapter 7).

The scores in both the neural network and K-Nearest neighbor range from 0.0 to 1.0. However, the scores do not sum to one, in that the scores are not probabilities. While it is possible to design better classifiers, the study merits a separate discussion and the interested reader is pointed to other resources [11, 26, 68].

5.2.3 Obtain Threshold Confidences for Classifiers

We use threshold classifier score to determine whether a core LBlock corresponds to a single component or to multiple components. This section describes the process used to obtain the threshold scores for each core class. A low score on the top choice result is accompanied by lower scores for the subsequent choices. Hence only the top-choice score is used for comparison purposes. For a given LBlock, if the classifier score for the top-choice class is less than the threshold score of that class, the LBlock is segmented further using the process given in section 5.1.2. Characters containing conjuncts and descenders need to be segmented, while consonants/vowels contain only one core component and should not be segmented.

In our problem, a false positive occurs when a character having conjuncts/descender is not selected for segmentation, while a true positive occurs when a consonants/vowels is not selected for segmentation. A low threshold would increase the number of false positives and a high thresh-

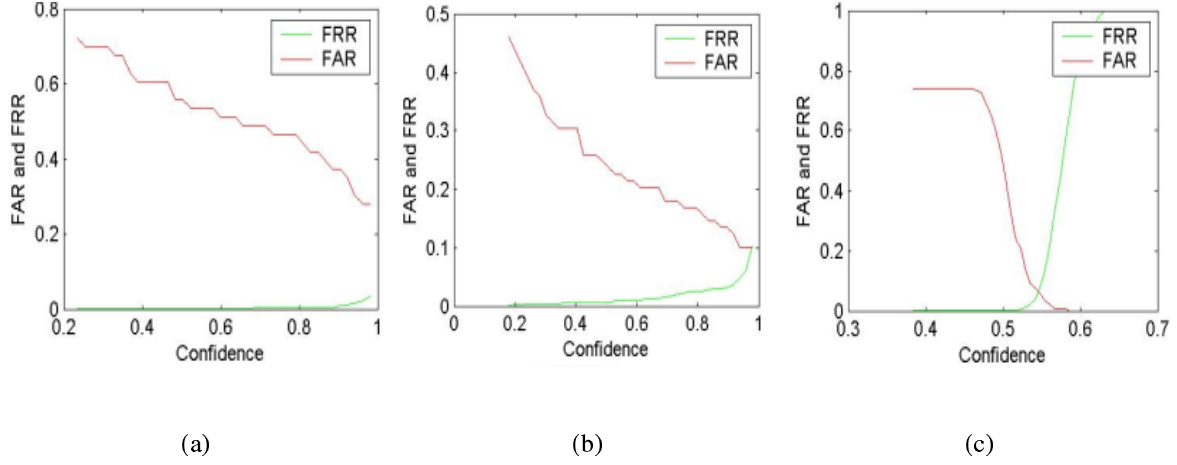


Figure 5.5: Classifier receiver operator characteristics for a full consonant: (a) Neural network with gradient features, (b) Neural Network with GSC features, (c) K-Nearest neighbor with GSC features

old would decrease true positives. Analysis of receiver operator characteristics (ROC) is an accepted technique of arriving at thresholds that optimize the true positive and false positive rate of a system [23, 28, 60]. We obtain a separate threshold confidence for each class by analyzing the respective receiver operator characteristics (Figure 5.5).

To analyze ROC curves, we construct an image test set containing consonants, vowels, conjuncts and consonant-descenders. The classifier result for each sample is represented as: $\langle I, truth, r, c \rangle$, where I : input image, $truth$: true class label, r : top choice class label ($r \in \{\text{consonants/vowels}\}$), and c : top choice score. The result is a true positive if $truth \in \{\text{consonants/vowels}\}$ and $c \geq \hat{c}_r$, where \hat{c}_r is the threshold confidence for class r . The result is a false positive if $truth \notin \{\text{consonants/vowels}\}$ and $c \geq \hat{c}_r$. For class r , $True_r$ and $False_r$ are the number of true positives and false positives respectively. False Accept Rate (FAR) and False Reject Rate (FRR) is computed for each class by taking different values of \hat{c}_r (Equation 5.1). The confidence corresponding to equal error rate is selected as threshold confidence of the class. True positive and false positive rates for these thresholds show that the GSC based neural network and KNN classifiers give the best performance (Table 7.1). In a second test, only consonant/vowel images have been used

as input to the classifiers and all three classifiers are found to give high accuracies (Figure 7.2). However, GSC based classifiers outperform the gradient based classifier in the important stage of separating conjuncts and descenders from vowels.

$$\begin{aligned}
 FRR &= 1 - \frac{\sum_{r=1}^n True_r}{\text{Number of consonant/vowel samples}} \\
 FAR &= \frac{\sum_{r=1}^n False_r}{\text{Number of conjunct and consonant-descender samples}} \\
 n &= \{\text{Number of consonant/vowel classes}\}
 \end{aligned} \tag{5.1}$$

5.2.4 Ranking Segmentation Results

Our recognition driven segmentation method returns one result for each ascender LBlock in the word and multiple results for the remaining LBlocks. This section discusses two issues which arise because of these multiple results, and explains how component results are reordered before applying a language model. The multiple results lead to two issues: (i) In characters having two components, two class labels are generated at each segmentation result. It often happens that the correct class for one of the component is obtained from a particular segmentation hypothesis, and the correct class for the second component is obtained at a different segmentation hypothesis. (ii) The same class label may be output by the classifier for multiple segmentation hypotheses.

We will use an example to illustrate the issues mentioned above. In a sample of the word पूछती four segmentation hypotheses are obtained for the first LBlock (Figure 5.4(b)). The correct result for this LBlock is ष and ू. However, ू is obtained as the class label for the first two segmentation hypotheses and ष is obtained in the last two segmentations. The class labels ष and ू are not obtained at a single segmentation. We create additional results by exhaustively enumerating all possible pairwise combinations of the class labels and scores. In case a class label occurs at multiple segmentation hypotheses, the higher score is considered to enumerate the pairs. In our example, six pairs are generated, bringing the total results for the first LBlock to seven (O_1 of Figure 5.6). Thirteen segmentation hypotheses are generated for the second LBlock, of which six

$$\begin{aligned}
O_1 &= \{\text{प}, \langle \text{र} \text{ } \rangle, \langle \text{प} \text{ } \rangle, \langle \text{त} \text{ } \rangle, \langle \text{र} \text{ } \rangle, \langle \text{प} \text{ } \rangle, \langle \text{त} \text{ } \rangle\} \\
O_2 &= \{\text{झ}, \langle \text{झ} \text{ } \rangle, \langle \text{झ} \text{ त} \rangle, \langle \text{झ} \text{ ढ} \rangle, \langle \text{छ} \text{ } \rangle, \langle \text{छ} \text{ त} \rangle, \langle \text{छ} \text{ ढ} \rangle\} \\
O_3 &= \{\text{ } \} \\
O_4 &= \{\text{ } \}
\end{aligned}$$

Figure 5.6: Results from the recognition driven segmentation performed in Figure 5.4; Classifier scores are not shown here

hypotheses are shown in Figure 5.4(c). By exhaustively enumerating the classifier results from the thirteen segmentation hypotheses, we get seven results for the second LBlock (O_2 of Figure 5.6).

The results from recognition driven segmentation are organized in the form: $O_i = \{o_{i0}s_{i0}, [\langle \check{o}_{i1}\check{s}_{i1}, \hat{o}_{i1}\hat{s}_{i1} \rangle, \dots, \langle \check{o}_{ik}\check{s}_{ik}, \hat{o}_{ik}\hat{s}_{ik} \rangle]\}$, where o and s represent the class label and score, and i represents the LBlock number. If the classifier returns multiple consonant/vowel choices for a particular LBlock i , or if multiple segmentation hypotheses are obtained for LBlock i , then $|O_i| > 1$. In the case of ascenders, only the top choice classifier result is retained, making $|O_i| = 1$.

5.2.5 Necessity of Additional Models to Generate Words

The stages discussed till now provide a list of components for an input word. In this section, we enumerate the reasons why additional processing is required to obtain word results from components. A language model which implements the additional processing steps is discussed in the next chapter.

- The segmentation stage divides some modifiers into two parts. Such division is seen in modifiers which contain both a core as well as an ascender component E.g: $\text{फ } \text{र } \text{ } \text{र}$. Figure 5.7(b) shows the example of modifier फ which is recognized as the sequence of class labels र , and $\text{ } \text{र}$. We use a grammar to identify such class label sequences and output the original un-segmented modifier.
- The recognition driven approach cannot distinguish between a half-consonant and a conso-

nant while segmenting characters horizontally. Conjunct characters contain a half-consonant fused horizontally with a consonant. A fused character contains a consonant fused horizontally with another consonant due to printing error or noise. The recognition driven segmentation stage treats both the conjunct and the fused character in a similar fashion. Half-consonants are not included in the class-space of the recognizer. It is assumed that when the image of a half-consonant form is sent to the classifier, the top choice result will be the corresponding consonant. Figures 5.8(a), 5.8(b), and 5.8(c) illustrate the example of conjunct/fused characters. While processing the results from a conjunct/fused character, a language model is used to select between the consonant and half-consonant.

- Multiple component results are often obtained for core components. The recognition driven segmentation uses threshold confidences to determine whether a core component has to be segmented or not. Noisy core components are often recognized with a low confidence, and segmented to obtain multiple hypotheses. Figure 5.9(a) shows a noisy word where three out of five core components have been wrongly over-segmented. Over-segmentation is also seen in words which have a font that is markedly different from the fonts used in the training set (Figure 5.9(c)). The recognizer retains top two results from the component classifier, examples of which are seen in figure 5.8 (The second choice is shown with the “Alt” label). Retaining multiple results from the component classifier is useful since the likelihood of getting the correct class among the top-n choices of a classifier is higher than getting the correct class in the top choice. Multiple component class labels can lead to large number of word results (We shall discuss this shortly). A language model is used to rank the different word results.

In the next chapter, we will design a probabilistic finite state automata to process the core, descender and ascender component hypotheses and generate word results.



Number of Components: Ascender: 1, Core: 2, Conjunct/Fused: 0, Descenders: 0

च Core: (च, 0.539062)
 ॐ Ascenders: (ॐ, 0.523438)
 न Core: (न, 0.541504)

(a) Word with no descender, conjunct/fused character



Number of Components: Ascender: 1, Core: 3, Conjunct/Fused: 1, Descenders: 0

। Core: (।, 0.510010)
 ँ Ascenders: (ॐ, 0.472656)
 ज Core: (ज, 0.524414)
 न Core: (न, 0.541748)

(b) Word with no descender, conjunct/fused character



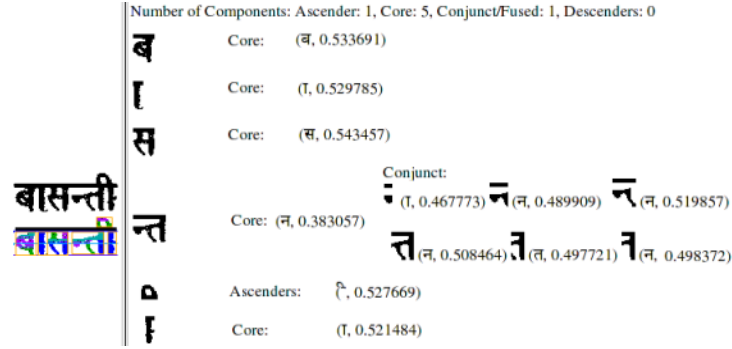
Number of Components: Ascender: 0, Core: 2, Conjunct/Fused: 0, Descenders: 1

Descender:

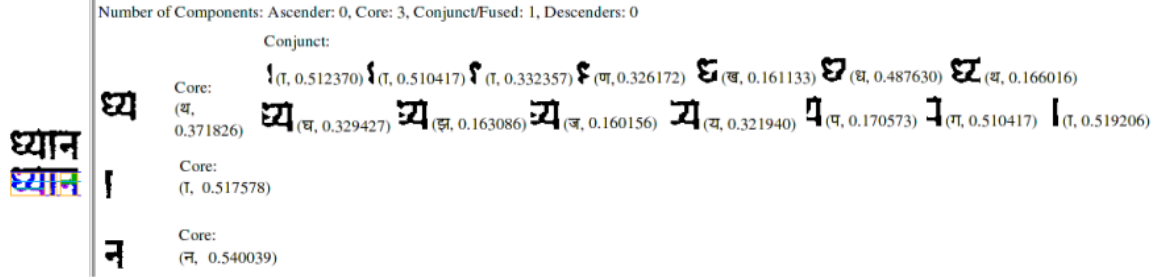
कु Core: (क, 0.502930) ँ (ग, 0.485677) ँ (न, 0.315430) ँ (क, 0.483073) क (क, 0.501628) क (क, 0.541992) क (क, 0.523763)
 ँ (ॐ, 0.425781) ँ (ॐ, 0.287760) ँ (ॐ, 0.291667) ँ (ॐ, 0.454427) ँ (ॐ, 0.314779) ँ (ॐ, 0.157878)
 ञ Core: (ञ, 0.530029)

(c) Word with a descender but no conjunct/fused character

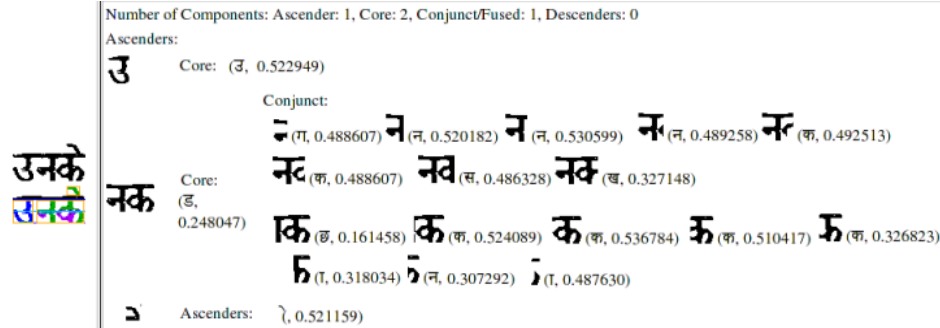
Figure 5.7: Examples of results from recognition driven segmentation. Word images are on the left and component results are on the right.



(a) A word where the half-consonant न of the consonant न is used in a conjunct



(b) A word where the half-consonant ध is connected to य



(c) A word where the consonant न is connected to consonant क due to print error

Figure 5.8: Examples of results from recognition driven segmentation where some core components correspond to fused or conjunct characters. Word images are on the left and component results are on the right.

Chapter 6

Stochastic Modeling to Generate Words

The recognition driven segmentation processes word BAGs and obtains component results for each character. Some of these components directly map to characters, while others have to be combined with their neighboring components to form a character. Once characters are formed from the components, they only have to be placed in sequence to form a Devanagari word. The process outlined in chapter 5 is capable of recognizing the 52 most common Devanagari components. If Devanagari characters could be formed by permutation alone, the number of characters would be $\sum_n P_k$, with $n = 52$ and k indicating the number of components in a character. Characters with upto 3 components occur most frequently in Devanagari, but 4 or 5 components are also seen. By setting $k = 1, 2, 3$, the possible character permutations is 1,35,252. The number of observed Devanagari characters is lesser: Only 5,537 characters are present in a Hindi text corpus having 11 million words. This is because the underlying language has restrictions about which component sequences can and cannot form a legal character.

The formation of Devanagari characters and words from components is not an unrestricted permutation, and characteristics of the underlying language can be incorporated into a model to enhance recognition results. Depending on the language being processed, different linguistic rules govern the formation of characters from components. Devanagari can be used to write several

Sequence 1	प	झ	~		पझी
Sequence 2	प	<झ >	~		Cannot form a legal word
				
Sequence 41	<प _० >	<छ त>	~		पूछती
				
Sequence 48	<त _० >	<छ त>	~		तूछती
Sequence 49	<त _० >	<छ ढ>	~		तूछढी

Figure 6.1: component sequences (left) and possible words (right).

languages - Hindi, Sanskrit, Marathi etc, and we consider the Hindi language for our study. In section 6.1, we discuss linguistic aspects which restrict the formation of characters from components. In section 6.2, we outline computational techniques to combine the component classifier results and the linguistic models into a single stochastic framework.

6.1 Linguistic Restrictions in Hindi

It is known that humans are familiar with a language or script use non-image knowledge while reading text. For example, readers familiar with Hindi can say that the word generated by sequence 48 of figure 6.1 is more probable than any of the other words. In remaining parts of this section, we describe linguistic restrictions that may be used by humans to make comparisons between words.

6.1.1 Script Writing Grammar

Script writing grammar is a set of rules that determine how Devanagari components should be combined to form characters. Sequences of components which violate these rules cannot form valid characters. For example, by applying the rule “Half-consonants cannot have vowel modifiers”, one can determine that the component sequence त_० is illegal. Similarly, using the rule “Only one vowel modifier is allowed on a consonant”, one can eliminate a different sequence of components. The utility of script writing grammar can be illustrated by considering the classifier results of our

$$O_1 = \{\text{प}, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle\}$$

$$O_2 = \{\text{झ}, \langle \text{झ} \rangle, \langle \text{झ त} \rangle, \langle \text{झ ढ} \rangle, \langle \text{छ} \rangle, \langle \text{छ त} \rangle, \langle \text{छ ढ} \rangle\}$$

$$O_3 = \{\sim\}$$

$$O_4 = \{|\}$$

(a) Original components generated from recognition driven segmentation of figure 5.4

$$O_1 = \{\text{प}, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle\}$$

$$O_2 = \{\text{झ}, \langle \text{झ} \rangle, \langle \text{झ त} \rangle, \langle \text{झ ढ} \rangle, \langle \text{छ} \rangle, \langle \text{छ त} \rangle, \langle \text{छ ढ} \rangle\}$$

$$O_3 = \{\sim\}$$

$$O_4 = \{|\}$$

(b) Applying script writing grammar

$$O_1 = \{\text{प}, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle, \langle \text{र} \rangle, \langle \text{प} \rangle, \langle \text{त} \rangle\}$$

$$O_2 = \{\text{झ}, \langle \text{झ} \rangle, \langle \text{झ त} \rangle, \langle \text{झ ढ} \rangle, \langle \text{छ} \rangle, \langle \text{छ त} \rangle, \langle \text{छ ढ} \rangle\}$$

$$O_3 = \{\sim\}$$

$$O_4 = \{|\}$$

(c) Applying syllable sequences

Figure 6.2: Removing component sequences by applying linguistic restrictions.

sample word पूछती (Figure 5.4). Recognition driven segmentation gives a lattice of component results from the word image. At least three of these results can be pruned by using the two script writing grammar rules mentioned here (Figure 6.2). The original result for our sample word had seven results each for the first and second character. By permutation, 49 component sequences can be formed and these sequences can generate 49 words. By imposing the two script writing rules, we have reduced the component hypotheses of the first and second characters to 6 and 5 respectively, reducing the number of word results to 30.

Certain rules determine how components combine to form alphabets and characters. Consider

for example the rule: “When a consonant is followed by the sequence $\text{ँ} \text{ॱ}$, the sequence will change to the consonant followed by vowel modifier ँ ”. One can use this rule to change the sequence $\text{तँ} \text{ॱ}$ into the alphabets तँ and ँ . These two alphabets in turn represent the character ती .

Prior researchers have used a decision-tree framework to implement these rules [6, 85]. We encode these restrictions in a context sensitive grammar, and use the grammar to design a computational model. Equation 6.1 encodes the combination of alphabets into characters. The grammar which encodes the combination of components into alphabets is given in table 6.1. Section 6.2.2 describes a computational model that can incorporate the grammar to prune illegal component results and to combine components into alphabets and characters.

$$\begin{aligned} \text{Devanagari character} &= V|(c^*C[v][s]) \\ \text{where } V &\text{ - Vowel, } C \text{ - Consonant, } c \text{ - Half-consonant,} \\ v &\text{ - Vowel modifier,} \\ s &\text{ - Special accent (nukta, ')} \\ \text{Devanagari word} &= (\text{Devanagari character})^+ \end{aligned} \tag{6.1}$$

6.1.2 Restrictions on Syllable Co-occurrence

As discussed in section 1.5, psychological and linguistic studies show that humans rely on their knowledge of sounds to read text [41, 97]. The use of sound models is found to hold true across multiple languages: English [10, 93], Hebrew [14, 76], Japanese [27], Chinese [90, 45], and Hindi [70]. These studies further suggest that the linguistic units used by humans correspond more closely to syllables than to sub-syllable units like alphabets or extra-syllable units like prefixes. A study by Taylor reviews research on different linguistic units that are used by humans (Chapter 8, [97]), and says the following: “In word recognition, splitting a word into syllable units permits subjects to overcome the effect of time separation between presentation of parts of the word, whereas splitting into similar sized non-syllable clusters does not.”. In the case of Hindi, it

Vowels having only one component:

अ इ ए

Consonants having only one component:

क ख घ च छ ज झ ञ ट ठ ड ढ त थ द ध न
प फ ब भ म य र ल व ष स ह

Vowels having multiple components:

आ = अ ा

ई = इ ऀ

ऐ = ए ै

ओ = आ ो

औ = ओ ौ | आ ै

Consonants having multiple components:

ग = ग ा

ण = ण ा

श = श ा

$V \in \{\text{अ इ ए आ ई ऐ ओ औ}\}$

$C \in \{\text{क ख ग घ च छ ज झ ञ
ट ठ ड ढ ण त थ द ध न
प फ ब भ म य र ल व श ष स ह}\}$

(a) Vowel (V)

(b) Consonant (C)

half-consonant label, $l \in \{\text{क ख ग घ च छ ज झ ञ ट ठ ड ढ ण त थ द ध न प फ ब भ म य र ल व श ष स ह}\}$
 $c = l | C$

(c) Half-consonant (c)

Modifiers having only one component:

ँ ॡ ॢ ॣ । ॥ ० १ २ ३ ४ ५ ६ ७ ८ ९

Modifiers having two components:

$(c^*C) \text{ ० } = \text{ ० } (c^*C)$

$(c^*C) \text{ १ } = (c^*C) \text{ १ }$

$\text{ २ } = \text{ २ } | \text{ ३ } \text{ ४ }$

$\text{ ५ } = \text{ ५ }$

$\text{ ६ } = \text{ ६ }$

$v \in \{\text{ँ ॡ ॢ ॣ । ॥ ० १ २ ३ ४ ५ ६ ७ ८ ९ फि फे फी फे}\}$

(d) Vowel modifier (v)

Table 6.1: Formation of alphabets from components

is suggested that a set of phonological rules closely represent the model used by native users to decipher Hindi text [70]. A list of the phonological rules is provided in appendix A.4, and a few of the rules are discussed in this section.

To explain the relevance of the Hindi phonetic rules in our OCR problem, we will see how the rules can be manually applied to prune component results obtained from the sample word पूछती (Figure 6.2(a)). One of the restrictions is: “Two breathy voiced consonants (BVC, ख घ छ झ ठ ढ थ ध फ भ) co-occur very rarely in a word.”. By this rationale, we can eliminate two component sequences from our result: झ ढ and छ ढ (Figure 6.2(c)). Another condition is: “Breathy voiced consonants rarely appear as conjuncts with other consonants”. This inhibits occurrence of the following components: ङ् झ , ण् झ , त् झ , ङ् छ , ण् छ , त् छ . By imposing the previous two restrictions as rules, we can eliminate five component sequences from our result (Figure 6.2(c)).

The script writing grammar in the previous section reduced three of the component hypotheses. When the script writing grammar and syllable restrictions are applied together, the number of results for the first and second character is reduced from seven to three each, effectively reducing our hypotheses from forty-nine words to only nine words. Significant gains can be made by pruning recognition results using the script writing grammar rules and phonological restrictions.

Adapting Syllable Restrictions to OCR

While phonological restrictions can be useful in pruning OCR choices, we do not use them in the form presented by linguistics. Notice the conditional term used in the following restriction: “Two BVC co-occur **very rarely** in a word” or “Breathy voiced consonants **rarely** appear as conjuncts with other consonants”. N-grams derived from a text corpus can implicitly model these restrictions. An n-gram model specifies the conditional probability of each unit of text (Called a token) c_n with respect to its history $c_1 \dots c_{n-1}$ [20]:

$$P(c_n|c_{1..n-1}) = \frac{P(c_1...c_n)}{P(c_1...c_{n-1})} = \frac{Count(c_1...c_n)}{Count(c_1...c_{n-1})} \quad (6.2)$$

Previous techniques in language modeling have used alphabets, suffix/prefix units, or words as tokens in n-grams [42, 44, 80, 72]. Most of the linguistic rules in Hindi are defined on syllables, and it seems likely that a syllable ngram model can be useful in Devanagari OCR.

Before designing a syllable n-gram model, we need to identify a technique that can provide syllables from a word. Sproat discusses computational models to identify syllables in different languages, including English, Hindi, Japanese, etc [91]. The advantage with Hindi is that it is a syllabic language, and Devanagari is a syllabic alphabetic script. The orthographic syllable, or character of Devanagari closely represents a Hindi syllable (Orthographic syllable explained in section A.3.1). A Devanagari word can be split into characters or orthographic syllables by using the grammar defined in equation 6.1. Thus, we are able to study the use of both alphabets and syllables as tokens in the n-gram modeling of Hindi.

6.2 Algorithms to Model Linguistic Phenomena

The script writing grammar discussed in section 6.1.1 has been applied in a rule-based framework for OCR post-processing [84]. Dictionary lookup has been tried for Hindi [7]. The additional linguistic restrictions described in section 6.1.2 have not been adapted in OCR designs. In this section, we will discuss corpus analysis of Hindi using N-gram model. A Stochastic Finite State Automaton (SFSA) is designed to combine components into characters and words. Lastly, we discuss how N-gram frequencies and classifier scores can be integrated into the SFSA to generate a ranked list of words.

6.2.1 Corpus Analysis of Hindi N-grams

To establish the use of different alphabet and syllable units in modeling Devanagari, we studied n-gram frequencies of the EMILLE Hindi corpus [5]. After removing punctuation, digits, and non-Devanagari alphabets the corpus has 11,381,720 words and 59,524,848 alphabets.

Each token c_i of an n-gram model can be an alphabet and will give us alphabet n-grams, or it can be syllables, in which case we get syllable n-grams. The units or tokens of the n-gram model are significant to model design and behavior. A common technique to design better models is classing in which tokens of a specific type are included in a common class, and instances of the token c are considered as instances of the class \hat{c} , $c \in \hat{c}$. In this case, the equation changes to:

$$P(\hat{c}_n | \hat{c}_{1..n-1}) = \frac{P(\hat{c}_1 \dots \hat{c}_n)}{P(\hat{c}_1 \dots \hat{c}_{n-1})} = \frac{\text{Count}(\hat{c}_1 \dots \hat{c}_n)}{\text{Count}(\hat{c}_1 \dots \hat{c}_{n-1})} \quad (6.3)$$

We use different classing techniques based on the phonetic properties of Devanagari alphabets. Trigram and bigram frequencies are obtained by applying eight tokenizing techniques:

1. **Alphabets with no classing:** No classing is performed, i.e. Each alphabet of Devanagari is considered as a separate class.
2. **Phonetic grouping of vowels and consonants:** Tokens are alphabets shown in table 1.1 with each row of consonants and vowels being grouped into one class.
3. **Breathy voiced consonants in separate classes:** Similar to case 2, except that each breathy voiced consonant (BVC, ख घ छ झ ठ ढ थ ध फ भ) is placed in a group by itself.
4. **BVCs in one class:** Similar to case 2, except that all BVCs are placed in one class.
5. **Nasals in one class and BVCs in separate classes:** Similar to case 2, except all nasal consonants (NCs, ङ ञ ण न म) are grouped into one class, and BVCs in a different class.

Tokens	Number of classes	Perplexity (Pp)	
		Bigram	Trigram
Alphabet	126	13.64	6.23
Alphabets grouped by phonetic properties	26	8.28	6.06
	36	9.00	6.10
	27	8.63	6.15
	28	9.00	6.23
	32	9.38	6.19
	41	9.71	6.15
Character	5538	18.51	3.53

Table 6.2: Bigram and trigram perplexity computed on the EMILLE dataset

6. **NCs in separate classes and BVCs in one class:** Similar to case 2, except that each NC is assigned a separate class, and all BVCs are placed in one class.
7. **NCs and BVCs in separate classes:** Similar to case 2, except that each NC and BVC is assigned a separate class.
8. **Characters with no classing:** Tokens correspond to characters in the text, and each character is a separate class.

$$Pp(D, q) = 2^{-\frac{1}{N} \sum_x \log q(x)}, q \rightarrow \text{n-gram probability} \quad (6.4)$$

The perplexity (Pp) of each language model is computed using equation 6.4, and is shown in table 6.2. A perplexity of k indicates that while predicting each token, the language model is as surprised as a uniform, independent selection from k tokens [20]. In section 6.2.2, we incorporate n-gram statistics into an SFSA to generate Devanagari words.

6.2.2 Stochastic Finite State Automaton

The general methodology of designing an SFSA is to create a finite-state machine (FSM) that accepts strings from a particular language and assign data-driven probabilities to transitions be-

tween states of the FSM. If a candidate word is accepted, a score or probability is obtained and if rejected, corrections can be made using the minimum set of transitions that could not be traversed [50]. SFSA has been used to generate predictive text for hand-held devices [32], and to model error correcting parsers for OCR [51]. In both applications, the SFSA has been built using n-grams models of English suffixes/prefixes. SFSA has also been used in handwriting recognition to model the relation between image features like loops, cusps, or lines and the underlying alphabets and words [106]. In our methodology, the structure of SFSA is modeled using Devanagari script writing grammar (Equation 6.1), and is defined as: $\lambda = (Q, O, I, F, A, \Sigma, \delta)$, where

- Q : states
- O : observations (52 component classes, table 5.1)
- I : initial states, which indicate start of a character
- F : final states, which indicate end of a character
- $A = a_{i,j}(o)$: set of transitions where $a_{i,j}(o)$ represents the probability of going from state i to j while accepting the label $o \in O$.
- Σ : emission symbols (Alphabets from table 1.1, and the null symbol ϵ)
- $\delta : A \times \Sigma$, a mapping from transitions to the emissions.

When a sequence of observations are accepted by the FSA, emissions generated by corresponding transitions give the word result for this sequence. We clarify the design of our FSA by using an example; the following transitions can accept the classifier labels given in Figure 6.2(a) and emit the correct word पूछती:

$$S \xrightarrow{c} T_c \text{ emit } c, c \in \{ प छ त \}$$

$$T_c \xrightarrow{d} T_d \text{ emit } d, d \in \{ ् \}$$

$$T_c \xrightarrow{m} NT_m \text{ emit } \epsilon, m \in \{ \text{ascender} \}$$

$$NT_m \xrightarrow{l} T_v \text{ emit } v, l \in \{ \text{vertical bar} \}, v \in \{ \text{modifier} \}$$

$$T_c, T_d, T_v \xrightarrow{\epsilon} S \text{ emit } \epsilon$$

where $S \in I$ and $T_c, T_d, T_v \in F$. The four transitions in this example represent the following script writing grammar rules: (i) The modifier modifier is made of the ascender ascender connected to a consonant and a vertical bar, (ii) Only one vowel modifier can be attached to a consonant. The resulting FSA is showing in figure 6.3. Emissions on the third and fourth transitions allow us to represent the combination of components ascender and vertical bar to form modifier . The sequence of emissions for this example would be: प, छ, त, ती.

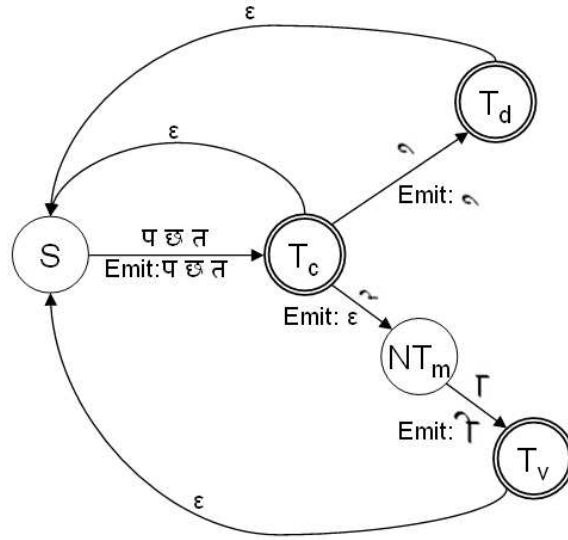


Figure 6.3: Finite State Automata that can accept the classifier labels in Figure 5.4 and emit the correct word पूछती

A finite set of transitions can represent the Devanagari script writing grammar defined by equation 6.1. While this completes the SFSA structure, we are yet to obtain transition probabilities for the model. In some SFSA models (For example handwriting recognition [106]), the relation between states (alphabets) and observations (features) is not apparent, and the Baum-Welch algorithm is used to learn transition probabilities. In our case, the observations are class labels from com-

ponent classifiers and the states are alphabets or characters. We use the script writing grammar defined in equation 6.1 to parse through Devanagari text, and obtain the sequences of states and the classifier symbols for each word. Transition probabilities for the SFSA are obtained from text corpus using equation 6.5.

$$a_{ij}(o) = \frac{\text{count of observing } o \text{ while going from state } i \text{ to } j}{\text{Number of transitions from state } i \text{ to } j} \quad (6.5)$$

Given a sequence of observation symbols, the best word hypothesis is generated using the Viterbi algorithm [33]. The Viterbi probability $\gamma_i(t)$, which is the highest probability of being in state i at time t is given by equation 6.6. The classifier score is also be incorporated into this process (Equation 6.7). By tracing the transitions which generate $\gamma_j(T)$, we get a valid Devanagari word. However, this model does not incorporate n-gram statistics, and generates only the topmost word result. In the next section, we discuss changes to the SFSA model to use alphabets and syllable n-grams and obtain multiple word results.

$$\gamma_j(t) = \begin{cases} 1 & , t = 0 \\ \max(\gamma_i(t-1)a_{ij}(o_t)) & , otherwise \end{cases} \quad (6.6)$$

$$\gamma_j(t) = \begin{cases} 1 & , t = 0 \\ \max(\gamma_i(t-1)a_{ij}(o_t)s_t) & , otherwise \end{cases} \quad (6.7)$$

6.2.3 Implementing Ngram Statistics in the SFSA

The transition and Viterbi probabilities in equations 6.5, 6.6 do not capture alphabet or character n-gram frequencies. This is because while equations 6.5, 6.6 are defined with respect to each state Q of the SFSA, only the start states S can indicate which character has been emitted. Most characters and even some alphabets are composed of multiple components, and several other states will be present between two start states. For example, while accepting classifier symbols for the word

पूछती (Figure 6.3), the state sequences are: $S^1 \rightarrow T_c^2 \rightarrow T_d^3 \rightarrow S^4 \rightarrow T_c^5 \rightarrow S^6 \rightarrow T_c^7 \rightarrow NT_m^8 \rightarrow T_v^9 \rightarrow S^{10}$ (Superscripts added to indicate number of states). To reach the terminal T_c^2 , the FSA accepts \P , at which time the character is \P . At the transition to T_d^3 , it accepts the symbol \P , which changes the character to \P . This is the emission between states S^1 and S^4 . Similarly, characters छ and ती are emissions between $S^4 - S^6$ and $S^6 - S^{10}$. These state transitions can capture n-gram statistics, but are not reflected in Viterbi probability estimation. We define another probability estimation as:

$$\hat{\gamma}_j(t) = \begin{cases} 1, t = 0 \\ \max[\hat{\gamma}_i(t-1)a_{ij}(o_t)], \text{ if } j \text{ is not a start state} \\ \max[\hat{\gamma}_i(t-1)a_{ij}(o_t)p(c_t|h_t)], \text{ if } j \text{ is a start} \\ \text{state, } c_t \text{ is the character between state } j \text{ and the} \\ \text{previous start state, } h_t \text{ is the n-gram history of } c_t \end{cases} \quad (6.8)$$

This modified probability estimation does not conform to a Markov assumption at the start states. Back-tracking is needed in cases where $\hat{\gamma}_j(t)$ evaluates to zero for all states j . Additionally, we need to obtain top-n results from our SFSA. With these constraints, we can no longer use the Viterbi algorithm to obtain word hypotheses. We perform a recursive search through the lattice to get all possible word recognition results and rank them according to $\hat{\gamma}_j(T)$ (Procedures 8, 9, and 10). The time complexity is exponential and would be expensive for a generic graph. In practice, words have limited number of conjunct and descender hypotheses and it is reasonable to perform an exhaustive search through the word recognition lattice.

Equation 6.8 does not incorporate the classifier scores for each observation. To add classifier confidences, the score of each observation s_t is multiplied with $a_{ij}(o_t)$ while calculating $\hat{\gamma}_j(t)$ (Equation 6.9). In the next section, we will discuss a few examples which illustrate the use of SFSA in generating word results.

Procedure 8 GetWords(O, λ)

- 1: Initialize: $t = 1$, Transitions vectors $V = V_1, \dots, V_n, \forall i$, set $V_i = \{ \}$ /* $n = |O|$ */
 - 2: Initialize: Sequence of start states and emissions, $S = \{ \}$, Current emission $E = \{ \}$ and Position of current start state $pos = 0$
 - 3: Using Equation 6.8, find all $\hat{\gamma}_j(t)$ that accept O_1 , add them to V_1
 - 4: $t = t + 1$
 - 5: Parse(V, O, t, S, E, pos)
-

Procedure 9 Parse($V, O, t, S, E, pos, \lambda$)

- 1: **if** V_{t-1} is not empty **then**
 - 2: **if** $t < N$ **then**
 - 3: Find $\hat{\gamma}_i(t-1) = \max[V_{t-1}]$
 - 4: Using Equation 6.8, find all $\hat{\gamma}_j(t)$ that accept O_t , from state i , and add them to V_t
 - 5: $E = E +$ emission corresponding to $\hat{\gamma}_i(t-1)$
 - 6: **if** i is a start state **then**
 - 7: Add emission E to S
 - 8: Set $E = \{ \}$, $pos = pos + 1$
 - 9: **end if**
 - 10: Parse(V, O, t, S, E, Pos)
 - 11: **else**
 - 12: **for all** $\hat{\gamma}_i(t-1) \in V_{t-1}$ and $i \in F$ **do**
 - 13: word = $S + E +$ emission corresponding to $\hat{\gamma}_i(t-1)$
 - 14: score = $\hat{\gamma}_i(t-1)$
 - 15: Print result: (word, score)
 - 16: **end for**
 - 17: Backtrack($V, O, t-1, S, E, pos, \lambda$)
 - 18: **end if**
 - 19: **else**
 - 20: Backtrack($V, O, t, S, E, pos, \lambda$)
 - 21: **end if**
-

Procedure 10 Backtrack($V, O, t, S, E, pos, \lambda$)

- 1: **if** $t < 2$ **then**
 - 2: stop
 - 3: **end if**
 - 4: Find $\hat{\gamma}_i(t-2) = \max[V_{t-2}]$, remove $\hat{\gamma}_i(t-2)$ from V_{t-2}
 - 5: **if** i is a start state **then**
 - 6: $E =$ Last stored emission from S
 - 7: Delete last stored emission from S
 - 8: $pos = pos - 1$
 - 9: **end if**
 - 10: Parse($V, O, t-1, S, E, pos$)
-

$$\tilde{\gamma}_j(t) = \begin{cases} 1, t = 0 \\ \max[\tilde{\gamma}_i(t-1)a_{ij}(o_t)s_t] \\ \max[\tilde{\gamma}_i(t-1)a_{ij}(o_t)s_t p(c_t|h_t)] \end{cases} \quad (6.9)$$

6.3 Generating Words from Components

As seen in section 5.2.5, the result of image based recognition is in the form of sequences of components. Each sequence is not guaranteed to form a valid words. A sequence of components is invalid if they violate the script writing grammar (Section 6.1.1). In addition, several words may be generated from the component sequences. The SFSA validates the component sequence against the script writing grammar and generates multiple word results. Further, each word result is given a score by combining n-gram statistics and classifier confidences. We rank word results by sorting them in descending order of their scores. In this section, we discuss a few examples to enumerate the utility of our SFSA model.

When a sample of the word चैन is presented to our recognizer, three components are generated by the component recognition stage (Figure 5.7(a)). The true word is also composed of three components: Two cores च न, and one ascender ँ . The sequence of three components is valid according to script writing grammar. When the component sequence is processed by the SFSA, we obtain the word results seen in the first row of table 6.3.

In case of the sample of चैन, a single component sequence is obtained during recognition. We examine a sample where multiple component sequences are generated from a word image. The word हरदम contains four core components: ह र द म, but additional components are obtained while recognizing the sample shown in figure 5.9(b). Similarly, figure 5.9(a), and figure 5.9(c) shows words where multiple results are obtained for certain core components. These additional components may occur when a core component is over-segmented, or when multiple class la-

bels are retained for a core component. Multiple components are also obtained while recognizing conjuncts, descenders, or fused characters (Figures 5.7(c), and 5.8). The SFSA model is able to process such multiple components and generate a correctly ranked set of word results.

While recognizing the word कुछ, one of the possible component sequences is: न्, छ (Figure 5.7(c)). According to the script writing grammar, this sequence is valid, and would generate the word नुछ. However, both the alphabet bigram model and the character bigram model provide only one result: कुछ (Table 6.3). In several cases, the character bigram model is able to better rank the word results. For example, in the word ध्यान, the correct word is ranked at the 10th position by the alphabet bigram model, while the correct word is ranked first by the character bigram model. In several other examples, the character model is able assign a better rank to the correct word. We also found examples where the character bigram performs worse than the alphabet bigram, but such cases are relatively few. In the next chapter, we will discuss experiments on a publicly available test set and compare our accuracies with previously reported figures.

S. No.	Image	Truth	Component result	Word result using different language models	
				Alphabet bigram	Character bigram
1	चैन	चैन	Figure 5.7(a)	चैन	चैन
2	जिन	जिन	Figure 5.7(b)	जिन	जिन
3	कुछ	कुछ	Figure 5.7(c)	कुछ	कुछ
4	बासन्ती	बासन्ती	Figure 5.8(a)	बासन्ती बरस	बासन्ती बरस
5	ध्यान	ध्यान	Figure 5.8(b)	न था हानि उधार ... (10th result: ध्यान)	ध्यान था हार घर उधार
6	उनके	उनके	Figure 5.8(c)	हे गले हैं उंच उनके ...(7th result: उनके)	उसके उनके उसने उतरने गले
7	परीक्षा	परीक्षा	Figure 5.9(a)	परीमाड पर परम परनाम परी ... (38th result: परीक्षा)	परीक्षा परीमाड पक्ष परम
8	हरदम	हरदम	Figure 5.9(b)	हार लगातार सारा दान हरदम	हरदम दान हार लगातार रहा
9	बोली	बोली	Figure 5.9(c)	बोल बोली हो बेले हमारे	बोल बारे बोली बनी पत्नी

Table 6.3: Word recognition results after applying language models and dictionary lookup. Word results are sorted according to the score returned by the SFSA.

Chapter 7

Experiments

The ILT data set has word images along with their Unicode truth and can be used to test the final OCR. While words are sufficient to test the final OCR, characters and core components are required to train and test the classifiers outlined in section 5.2.3. We extract characters and core components from word images by using the semi-automated character segmentation technique outlined in section 3.2.4. The run-time operation required to extract these images is described in appendix B. The extracted images have 24,689 samples of consonants and vowels, 686 conjunct characters, 479 descender characters, 6,033 samples of ascender components, and 481 samples of descender components. The image set corresponding to consonants, vowels, ascenders and descenders is divided into 60% training set and 40% testing set. The conjunct and descender characters cover 97 and 56 classes respectively. A dataset of 10,606 words has been used to test the OCR. A summary of the different stages of our OCR and the recognition results at each stage is provided in figure 7.1. In further sections, results from each stage of the OCR are reported separately.

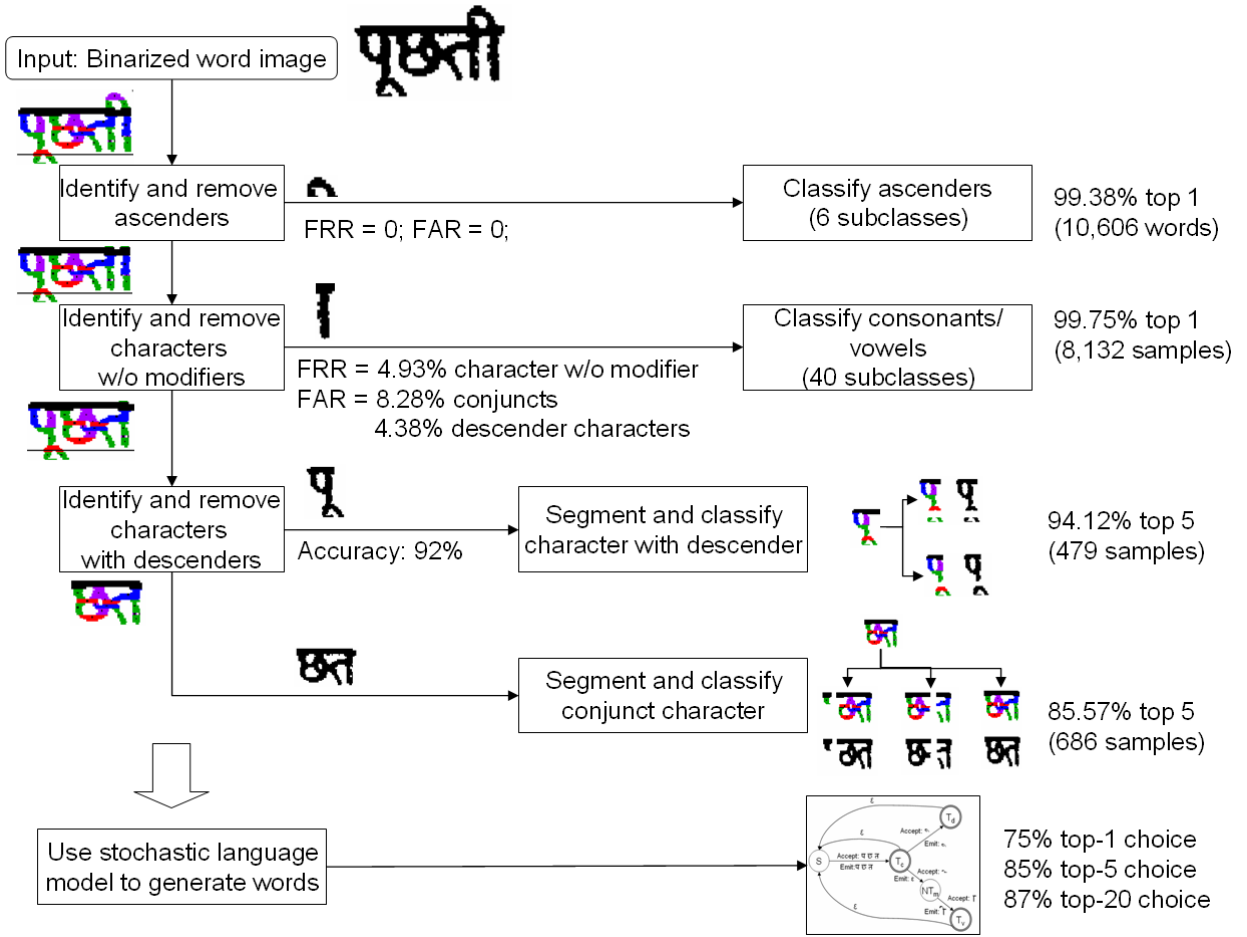


Figure 7.1: Overview of results from each stage of recognition

7.1 Character Recognition

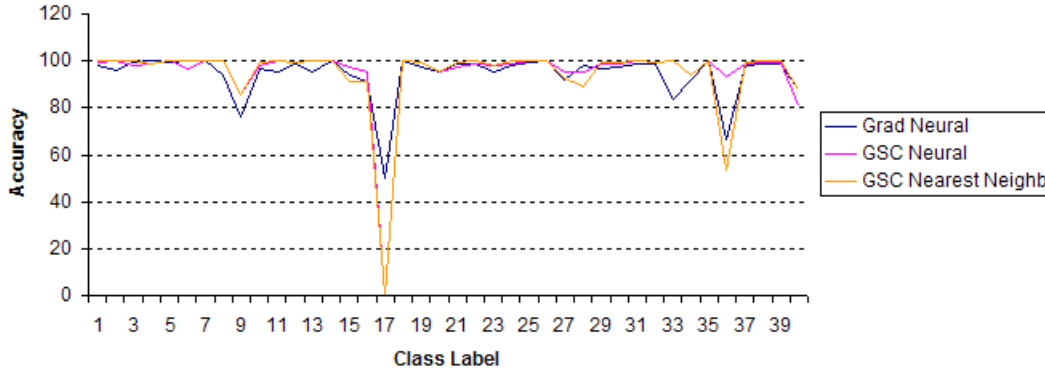
We discuss the accuracy of character recognition in terms of three stages: (i) Accepting a consonant/vowel and rejecting a conjunct/descender by using the component classifier, (ii) Identifying the direction of segmentation for rejected images (conjuncts or descenders) by using baselines, and (iii) Using the BAG to segment a conjunct from left to right or a descender from top to bottom and recognizing the constituent components.

Results for stage (i) include component classifier accuracy, false accept rate (FAR), and false reject rates (FRR). Accuracy of consonant and vowel recognition is in the range of 94-97% (Fig-

Classifier type	FRR	FAR	Breakup of FAR	
			Consonant-descender	Conjunct
Neural Network Gradient features	10.70%	48.85%	48.91%	48.83%
Neural Network GSC features	5.83%	6.88%	7.40%	6.70%
K-NN GSC features	4.93%	7.28%	4.38%	8.28%

Table 7.1: False reject rate (FRR) and false accept rate (FAR) obtained at the equal error rate for consonant/vowel classifiers, Testing is done on 8,029 consonants and vowels, 479 consonant-descenders and 686 conjuncts

ure 7.2). A false accept in our system means that the classifier score for a conjunct/descender is greater than the threshold score. The consonant/descender is wrongly considered as having one core component, and is not segmented. A false reject occurs when the classifier returns a low score for a consonant/vowel, which leads to segmentation of the consonant/vowel. False accepts result in under-segmentation and false rejects lead to over-segmentation. The best FAR and FRR is obtained using GSC neural network classifier: 5.83% FRR and 6.88% FAR (Table 7.1).



(a)

Figure 7.2: Classifier accuracies for 8,029 consonant/vowel images, Average accuracy: Neural network with gradient features: 94.36%, Neural network with GSC features: 95.25%, and K-nearest neighbor with GSC features: 96.97%

For stage (ii), we find that errors occur in words having single characters, E.g.: $\frac{1}{2}$. In our dataset, 846 words have conjuncts and 727 have descenders, and errors were found in only 35 and 94 words respectively.

Character type	Recognized component	Top Choice				
		1	2	3	4	5
Descender	consonant descender	80.17%	86.64%	93.11%	97.08%	98.74%
		84.76%	-	-	-	-
Conjuncts	half-consonant	16.91%	51.46%	67.49%	83.97%	94.17%
	consonant	46.50%	59.91%	71.43%	82.80%	92.27%

(a) K-Nearest neighbor classifier using GSC features

Character type	Recognized component	Top Choice				
		1	2	3	4	5
Descender	consonant descender	38.00%	62.42%	77.45%	92.49%	98.12%
		86.31%	-	-	-	-
Conjuncts	half-consonant	25.66%	52.48%	74.20%	82.94%	85.57%
	consonant	27.41%	47.10%	66.62%	83.10%	91.11%

(b) Neural network classifier using GSC features

Table 7.2: Recognition results from 479 descender and 686 conjunct characters

Results for stage (iii) are shown in Table 7.2. Top choice accuracy in recognizing the conjunct and descender characters is low (26 to 38%) in comparison to the top-5 choice accuracy (91 to 98%). Errors are mainly due to confusions between similar consonants/vowels; e.g: क्य being recognized as फ य or क थ instead of क य.

The character data set used above is identical to the test set used for our segmentation driven recognizer (Section 3.2). A comparison of conjunct recognition using the segmentation driven recognizer and the new BAG based recognition driven segmentation approach is shown in table 7.3. While the former approach gives 66.07% accuracy on 32 frequently occurring conjuncts, the BAG approach achieves 94.17% accuracy in top-5 choices (Table 7.3). The BAG approach is also recognizes an additional 65 classes with 77.09% top-5 choice accuracy.

7.2 Word Recognition

We report word recognition results using five SFSA models. The following SFSA models incorporate classifier confidences: (i) Script writing grammar, with no n-gram statistics, (ii) SFSA with alphabet bigram, (iii) SFSA with character bigrams; the following SFSA models do not incorporate classifier confidences: (iv) SFSA with character bigrams (v) SFSA with character trigrams. Of

Number (#) of classes	Conjunct type	E.g.	# of samples	Percent covered	Technique					
					Segmentation driven technique (Section 3.2)			Recognition driven technique (Section 5.2)		
					# of classes	Percent covered	Accuracy (%)	# of classes	Percent covered	Accuracy (%)
31	High- frequency	क्य	597	33.58	31	33.58	66.09	31	33.58	94.17
65	Low- frequency	ष्ट	335	18.84	-	-	-	65	18.84	77.09
11	<i>Rakaar</i>	प्र	444	29.86	1	20.17	89.28	1	20.17	94.07
17	Others	ज्ज्व क्यु ह	111	6.24	-	-	-	-	-	-
124			1487	100.00	32	60.32		97	87.60	

Table 7.3: Comparison of conjunct recognition using segmentation driven and BAG based techniques

these, (i) uses equation 6.6 to compute a score for each word result, (ii) and (iii) use equation 6.9, and (iv) and (v) use equation 6.8. Component classification is performed by a neural network classifier using GSC features. Dictionary lookup is performed using a lexicon of 4,291 words, and Levenshtien edit distance.

While exhaustive search through a lattice (Procedure 8) is exponential in time complexity, we find that the system takes an average of 3 seconds to recognize each word. The test machine has a single core AMD Athlon 3.8Ghz processor with 512MB RAM. For each decision to segment a character using BAG, an average of 7 consonant-descender or conjunct hypotheses are generated. The number of words actually requiring segmentation is 1,573, but segmentation was performed in 5,173 words. Empirical analysis shows that the correct consonant is always selected despite the tendency to over-segment characters. The best results are obtained using the SFSA with character bigrams and classifier scores: top-choice accuracy of 75%, top-5 choice accuracy of 87%, and top-20 accuracy of 87%.

The performance of recognition driven BAG approach is consistent across different types of documents. Figure 7.3 shows a plot of accuracy across different documents. The top-5 accuracy of recognition driven BAG ranges from 72% to a high of 90%, while its top-1 choice accuracy is

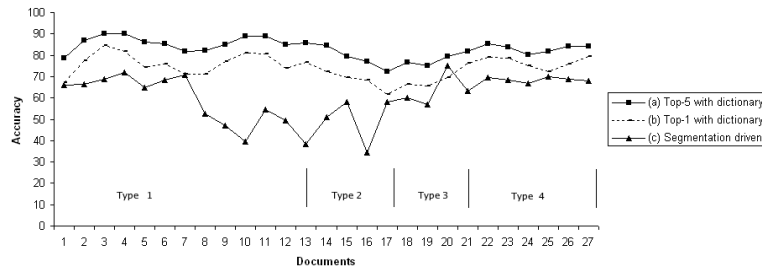


Figure 7.3: Word accuracy on different documents (a,b): Using the recognition driven segmentation approach, (c): Using the segmentation driven approach. The accuracy of segmentation driven approach varies widely across different documents, while the recognition driven segmentation approach provides more consistent accuracy

in the range of 62% to 85%. In comparison, top-1 accuracy of the segmentation driven approach varies widely: 39% to 75%. The segmentation driven approach does not provide multiple word results.

A comparison of our word recognition results with previously reported accuracies is given in table 7.5. Our implementation of a segmentation driven recognizer achieves a word recognition rate of upto 74% on some documents. This is comparable with accuracies in some single-font recognition techniques. However, the technique does not scale when applied to documents with different fonts, and we get a low accuracy of 53% on multi-font documents. The recognition driven BAG approach performs more favorably. While the top-1 choice accuracy is lower at 75%, the top-5 accuracy of 87% beats all methods except one. The performance improvement is significant and we are able to perform consistently across multiple font types.

Top-N	Using classifier confidence (Equation 6.9)						Without classifier confidence(Equation 6.8)			
	Script writing		Alphabet bigram		Character bigram-1		Character bigram-2		Character trigram	
1	13.25	2.75	52.84	1.48	70.17	1.12	68.88	1.18	69.96	1.29
2	36.98	1.82	62.52	1.10	76.45	0.85	75.75	0.88	77.80	0.95
3	51.10	1.48	69.05	0.92	77.90	0.76	77.42	0.78	79.68	0.85
5	65.29	0.99	73.47	0.76	79.60	0.67	79.28	0.69	80.48	0.78
10	74.78	0.71	77.18	0.64	80.38	0.61	80.33	0.61	80.69	0.72
15	77.07	0.63	78.33	0.59	80.55	0.59	80.52	0.59	80.71	0.70
20	78.17	0.59	78.75	0.57	80.64	0.58	80.61	0.58	80.71	0.69

(a) Results without dictionary lookup

Top-N	Using classifier confidence (Equation 6.9)						Without classifier confidence(Equation 6.8)			
	Script writing		Alphabet bigram		Character bigram-1		Character bigram-2		Character trigram	
1	32.99	2.34	61.20	1.46	74.71	1.09	73.33	1.15	74.93	1.11
2	61.42	1.37	73.26	1.01	81.08	0.77	80.24	0.81	81.54	0.80
3	67.74	1.08	77.15	0.84	82.98	0.66	82.41	0.68	83.32	0.71
5	79.80	0.70	81.14	0.67	85.13	0.56	84.74	0.58	84.38	0.64
10	84.85	0.51	84.59	0.53	86.41	0.49	86.32	0.49	85.16	0.59
15	86.67	0.44	85.99	0.47	86.86	0.46	86.78	0.46	85.50	0.57
20	87.32	0.42	86.87	0.43	87.10	0.45	87.06	0.45	85.63	0.56

(b) Results with dictionary lookup

Table 7.4: Accuracy and average string edit distance using different n-gram models. The test set has 10,606 words

OCR Type	Font restriction, Classifier	Accuracy	
		Character	Word
Segmentation driven OCR	Single font, component classifier [6]	93%	67%
	Two font, component classifier [64]	80%	87%
	Two font, character classifier [18]	90%	91%
	Single-font, component classifier (Chapter 3)	84%	74%
	Multi-font, component classifier (Chapter 3)	84%	53%
Recognition driven OCR	Multi-font, top-1 (Chapters 5, 6)	96%	75%
	Multi-font, top-20 (Chapters 5, 6)	96%	87%

Table 7.5: Comparison of Recognition Results

Chapter 8

Conclusions

8.1 Summary

This thesis addresses the problem of multi-font Devanagari OCR. We have designed a novel BAG-based segmentation technique to accommodate characters which may be fused along non-linear boundaries. By using the BAG, we also design algorithms to correct aberrations like noisy character contours or disconnected joints in characters. This is accomplished by taking advantage of the BAG data structure. The BAG is constructed by grouping pixels of a word image into logical blocks, and edges are added to represent connectivity of the blocks. We show that the resulting data structure can provide a higher-level model of characters and can be used to perform image recognition operations that cannot be easily supported by other representations.

Motivated by the success of recognition driven segmentation in Latin and CJK OCR, we have developed a recognition framework that can selectively segment characters in a Devanagari word. We retain multiple segmentation hypotheses in the form of component class labels and scores returned by a classifier. We use script writing grammar to encode rules that govern the formation of Devanagari characters. We design an SFSA framework that can combine component results with script writing grammar and ngram language models to construct a list of word results and

scores. We take advantage of the syllabic-alphabetic nature of Devanagari script where one can use alphabets or orthographic syllables (characters) as the units for ngram modeling.

Font information is not used at any stage in image analysis, and our recognizer requires only the word image as input. This contrasts with prior approaches which use structural information gathered from lines or paragraphs of text to carry out image analysis operations. We incorporate one such segmentation-driven design to construct a baseline system and show that paragraph-based, or even line-based structural features are unreliable for Devanagari OCR.

Word accuracies achieved on multi-font documents are 75% top-choice accuracy and 87% top-10 choice accuracy (Table 7.5). In comparison with segmentation driven techniques, our recognizer performs more consistently across different documents (Figure 7.3). These results are a significant improvement over previous attempts on Devanagari OCR [6, 18, 64]. Publicly accessible datasets and annotation tools have been developed as part of this work. For the first time, a Devanagari recognition engine has been made available online¹.

The system is being integrated with Document Image Analysis (DIA) tools like text-image separation, line/word detection [38] to develop a complete OCR solution. Because our methodology needs only the word image, it can be easily integrated with off-the-shelf DIA tools.

8.2 Contributions

This dissertation contributes to the field of optical character recognition in the following areas:

- *Novel character enhancement and segmentation using Block Adjacency Graphs (BAG):* Characters have natural joints and holes of different sizes, which contribute to their distinct shapes. Traditional image representation and enhancement algorithms are not capable of modeling such structures and often distort the character shape. We adopt a BAG representation to model such natural structures. Previously, it has been shown that BAGs can be used

¹The OCR is located at: <http://phoebe.cedar.buffalo.edu/ocr/>

to extract higher level features from word images and assist in holistic recognition [106]. We design algorithms which can use the BAG structure to enhance character contours while preventing changes to salient character structures. We design a BAG-based method to divide characters at natural breaks or joints and perform character segmentation at non-linear boundaries in at least two directions.

- *Recognition driven segmentation for Devanagari OCR:* By corpus analysis, we determine that 88% of Devanagari text is composed of characters which map to single alphabets, while the remaining 12% contain characters which are formed by fusing one or more alphabets. We design a recognition driven segmentation methodology that can selectively determine the appropriate segmentation for each character. We achieve the following accuracies for character recognition: (i) 97% accuracy in characters having only one component, (ii) 98% top-5 choice accuracy for characters composed of vertically fused components, and (iii) 92% accuracy for characters composed of horizontally fused components. The design uses linguistic models to prune segmentation results and is robust to font variations. Previous approaches in literature have used a dissection based segmentation for Devanagari OCR with limited success in single font systems.
- *Integrating Devanagari script writing grammar and linguistic constraints into a single stochastic framework:* Previous approaches in Devanagari OCR have used script writing grammar, classifier scores and n-grams in separate units. By using a stochastic finite state automata (SFSA), we show that these can be combined into a single framework. We show that by introducing emission states in the transitions of SFSA, one can model Devanagari script writing grammar constraints, particularly cases where multiple components combine into a single linguistic unit. We take advantage of the syllabic alphabetic nature of Devanagari script and design ngram models where the units are orthographic syllables, which may be composed of one or more alphabets. By using an SFSA model, we obtain 75% top choice and 85%

top-5 choice word accuracy on a publicly available multi-font dataset. These results are a significant improvement over previous methods.

- *Devanagari data collection:* While Latin and CJK data resources and baseline systems were readily available from the early days of OCR research, similar data sets have not been introduced for Devanagari OCR. As part of the dissertation work, tools have been designed to enable semi-automated collection of Devanagari image data sets.

The most important facet of this work is that a real-time, web-based OCR has been designed which can accept Devanagari word images and output recognition results in an XML format.

8.3 Future Work

8.3.1 Recognition Driven Segmentation

The system cannot segment certain types of characters that are formed by adding small modifiers to the consonant, E.g.: ऋ ॠ ऌ are formed by adding a small stroke to the consonants ऋ ऋ ऌ. One approach to handle these characters is to treat them as additional classes in the consonant/vowel classifier. However, these characters occur only in 0.23% of the text, and sufficient samples cannot be obtained. Conjunct-descenders which contain fusion of three or more components are also not handled by our segmentation algorithm, E.g.: क्यु स्तू. These type of characters also amount to a very small percentage of Devanagari text (0.24% in ILT, and 0.33% in EMILLE).

8.3.2 Pruning Top-n Word Choices

While our top-1 choice accuracy compares well with prior Devanagari OCR approaches, the top-10 choice results can be enhanced using word or sentence level language models. The design of such language models is an interesting and challenging pattern recognition problem [80]. Our OCR solution provides a useful starting point for such models.

8.3.3 Applying the OCR to Other Languages

While we have used Hindi as the language for our current implementation, the OCR can be directly extended to other syllabic-alphabetic languages like Sanskrit or Marathi which use Devanagari. Extending the OCR to such languages involves extracting the syllable frequencies using a corpus of the respective language. Modifications to the SFSA structure may be needed to accommodate changes if any, in the script writing grammar.

Appendix A

Theoretical Relation of Scripts and Sounds

In this chapter, we present definitions and concepts relevant to our research. These concepts are from the domain of phonology and orthography, and explain the relation between scripts and their sounds [25, 37]. Language is the mode of communication that we refer to as English, French, Spanish, Hindi, Sanskrit, Chinese etc. These languages may be spoken, and the sounds of a language is what we call speech. Languages are written to communicate by vision and the set of written symbols constitute the script. For several languages, speech was the first and most dominant mode of communication. Most scripts have been developed to represent the underlying sounds of a language, and not the other way round [24]. Psychological experiments like the Stroop test (Figure 1.4) show that humans use auditory models to read text [10, 90, 93].

A.1 Articulation: The Generation of Sounds

Speech is caused by a process called articulation: Airflow from the lungs is modified in the throat, nose, and mouth. The primary articulatory organs in the throat are the vocal cords (Also called vocal fold). While producing certain sounds, the vocal cords open, letting air pass through the glottis. Such sounds are called *un-voiced* sounds, e.g. Sound corresponding to the ss in “hiss”.

When the vocal cords are closed or have only a small gap, a vibrating effect is created. Such sounds are called *voiced* sounds, an example being the *zz* in “buzz”. The nasal cavity is separated from the mouth by a small flap called the velum or soft palate. When the velum or soft palate opens during speech, air passes through both the nose and mouth. Sounds generated by such airflow are called *nasal* sounds, e.g.: The *ai* in “pain”. When the velum is closed, sound is generated only from the mouth and is called *non-nasal*.

The terms *voiced*, *un-voiced*, *nasal*, or *non-nasal* indicate how the throat and nose is involved in generating a sound. These terms can be referred to as properties of sound. By specifying the complete set of attributes, one can state the precise generation of a sound. Only two organs which contribute to four attributes have been described so far - there are many auditory organs in the mouth, which give rise to several other attributes.

In the mouth, sound is produced by two types of organs: Passive articulators - organs which do not move (upper lip, teeth, roof of the mouth and pharynx), and Active articulators - organs that move (lower lip, tongue). Different sounds are created by moving one of the active articulators to one of the passive articulators. The organs in mouth allow for two attributes: Manner of articulation, depending on how far the active and passive articulators are when generating a sound, and place of articulation which is based on which passive and active articulators are involved. The place of articulations are listed in table A.1. The manner of articulations are:

- *Stop* The active articulator completely presses against the passive articulator, leaving no gap for air to pass through. Air is then released suddenly by lowering the active articulator, e.g.: *pu* in the “push”.
- *Affricate* Similar to the stop, but the active articulator is lowered very little, giving only a narrow channel for the passage of air e.g.: *ch* in “churn”.
- *Fricative* Articulators are close to each other, but air is not stopped completely, we get fricatives, e.g. the *ze* in “zebra”.

Place of articulation	Active articulator	Passive articulator	Example
bilabial	lower lip	upper lip	<u>p</u> unch
labio-dental	Lower lip	upper teeth	<u>w</u> ater
dental	tip of the tongue	upper teeth	<u>th</u> irsty
alveolar	tip of the tongue	alveolar ridge	<u>l</u> arge
retroflex	tip of the tongue (curled)	area behind alveolar ridge	<u>tr</u> igger
palato-alveolar	blade of the tongue	area behind alveolar ridge	<u>sh</u> ake
palatal	back of the tongue	hard palate	<u>y</u> ak
velar	back of the tongue	velum	<u>g</u> ood
uvular*	tongue back	uvula	(French) <u>P</u> aris
pharyngeal*	tongue root	pharynx wall	
glottal	vocal cords	vocal cords	<u>h</u> are

Table A.1: Place of articulation. *These dont exist in English

- *Glide* Articulators are wide apart and but air flow is completely unhindered, e.g. the ya in “yak”.
- *Liquid* Articulators are in a position that would neither be considered a *fricative* or a *glide*. These are typically the “l” and “r” sounds.
- *Vowel* Glides which form complete syllables. Syllables will be discussed in the next section.

A.2 Phonological Syllable

The basic unit of sound is called a *phone* - it is an audible entity that the human brain can process. Each phone is generated from one of the organs discussed in the previous section. Natural speech is not a disconnected sequence of phones, but a continuous process called co-articulation: The sound of one phone affects the sound of its neighboring phone. However, the co-articulation effect is often seen to begin and end at “syllables”, and groups of sounds/phones can be identified as separate units.

The syllable may be treated as consisting of an onset, a nucleus and a coda, arranged in the representation shown in figure A.1. The onset and coda can be a combination of one or more

consonants, and the nucleus is composed of a vowel sound. The Onset and coda are optional parts of a syllable, whereas the nucleus is a required aspect of each syllable.

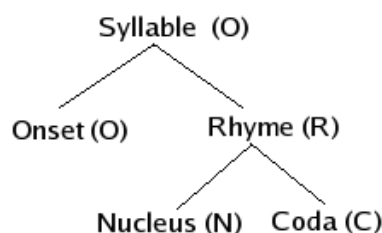


Figure A.1: Syllable structure

A.3 Orthography

As originally shown in [24], scripts used to write languages are broadly classified into two groups: phonological scripts, and non-phonological scripts. In non-phonological scripts, each symbol represents a physical real world concept, For e.g: Egyptian hieroglyphics and Chinese. In phonological scripts, each symbol represents a sound unit. When multiple symbols are placed in sequence, it indicate a sound sequence formed by combining individual sounds of each symbol. Each written symbol sequence and corresponding sound sequence represent real life concept.

Phonological scripts are further classified using the concept of Orthographic Depth (OD), a term used to explain how each written symbol corresponds to the underlying sound [91]. In some scripts, each symbol maps to a different sound based on the surrounding symbols, an example is English: The symbol *w* has different pronunciation in *water* and *owl*. These scripts are said to have deep orthography. In other scripts, the mapping between a written symbol and a sound unit does not vary with context. Such scripts are referred to have Shallow Orthography or Shallow Orthographic Depth. However, OD is not dependant only on the script, it is also dependant on the language. When Latin is used to write English, pronunciation of characters change and it has deep

OD. However, when used to write French, it has shallow OD.

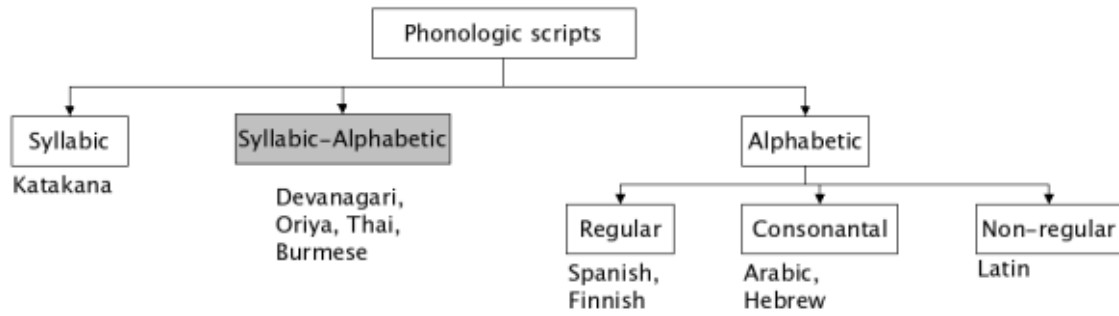


Figure A.2: Different types of phonological scripts

Scripts are further classified by combining orthographic depth and the type of sound units represented by each written symbol [24] (Figure A.2). At one end of the classification are syllabic scripts where each written symbol represents a syllable (Japanese). At the other end are non-regular alphabetic scripts like Latin, where the pronunciation can change with context. Devanagari falls under a class of scripts which are called syllabic-alphabetic.

A.3.1 The Orthographic Syllable in Devanagari

Devanagari has a set of symbols, which are called alphabets or glyphs in orthographic terminology (Figure 1.1). Each alphabet has a specific pronunciation. When used to write Sanskrit, the pronunciation never changes, giving it shallow OD. When Devanagari is used to write Hindi, the pronunciation of symbols can change at certain locations like the *anuswara* (Dot:').

The characteristic feature of Devanagari (and other syllabic-alphabetic scripts) is the formation of an *Orthographic Syllable* by combining one or more alphabets (Chapter 2 [91]). Given a Devanagari word, orthographic syllables can be identified by placing boundaries at the end of symbols which indicate a vowel sound. For example, in the word वरमा (varmaa), there is a vowel sound immediately after व (va) and before the रमा (rmaa). There are two orthographic syllables in

the word: व and माँ.

As observed in [91], “The Devanagari orthographic syllable is not isomorphic with the phonologic syllable”, for e.g.: In the Devanagari word वमाँ the pronunciation of the orthographic syllables (व माँ) is different from the phonological syllable (var and maa). The difference between orthographic syllable and the phonological syllable in Devanagari often occurs at the coda - The coda of a phonological syllable is sometimes written as a part of the following orthographic syllable.

The relation of orthographic and phonological syllables can be summarized as follows:

- If the orthographic syllables and phonological syllables of a word are listed in sequence, the nucleus of each phonological syllable would correspond to the vowel sound in the corresponding orthographic syllable.
- The difference between the orthographic syllable and the phonological syllable is often only at the coda.

A.4 Phonological Studies in Hindi

This section discusses constraints on sequences of consonants as presented by Ohala [70]. The phonological study includes experiments which examined if a native Hindi speaker would accept or reject words corresponding to different consonant sequences. The sequential constraints verified in this work provide an example of how certain consonants can be limited by using the phonetic properties of their neighbors. Some of the constraints outlined by Ohala have been summarized below:

1. The consonant ढ does not occur as a conjunct.
2. Two stops of the same kind will not co-occur, until and unless the conjunct is a geminate,
E.g.: क्क त्त.

3. Breathy voiced consonants ख घ छ झ ठ ढ थ ध फ भ do not occur in a sequence. Affricates, aspirates or breathy voiced stops will not occur before or after an affricate.
4. In a conjunct having two stops, if the second stop is breathy-voiced, the first one should be voiced.
5. If a consonant is [+long] (the same consonant voice repeated, like rra), the preceding segment must be a [-long] (single consonant sound) and the following segment must be a vowel.
6. In conjuncts that occur at the start or end of a word, if the first part is a glide then stops, nasals or fricatives cannot occur as the second part.
7. In conjuncts that occur at the start or end of a word, if both the first and second consonant in the conjunct is a stop or a fricative, then these must agree in voicing.
8. In the middle or end of a word, fricatives do not occur before breathy voiced consonants.
9. If a conjunct having a nasal followed by a stop occurs at the middle or end of a word, the nasal will be homoorganic with the stop. Two sounds are homoorganic when the same active organ is used for both sounds. In our example, this would mean that if the stop is a velar, the preceding nasal would also be a velar.
10. After a palatal nasal, the following segment should be a palatal consonant.
11. After ङ or ण, the following must be a velar or retroflex stop.
12. If the segment before a palatal consonant is a nasal, it must be a bilabial or palatal nasal.
13. Hindi does not contain stop followed by the the consonant sound ya followed by a front vowel sound किय, कयी etc.
14. ऋ does not occur initially and finally in a conjunct.

15. At the start and end of a word, after a य, only vowels occur.
16. If the first part of a consonant is a stop, the second part cannot be a stop or a nasal.
17. If the second part of a two-consonant cluster is a fricative, it must be a ष and it must be preceded by क.
18. If the first part of a conjunct is a nasal, then the second part cannot be a stop, fricative or nasal.
19. ऋ does not occur at the start of a word.
20. Only a vowel can occur before ह.
21. Only a vowel can occur before or after फ़.
22. In a conjunct, if the first part is a fricative and the second one is a stop, the fricative has to be स.
23. If in a conjunct, the second part is a ज्ञ, the first part can only be य.
24. If the second part of a conjunct is a liquid, the first part cannot be an affricate, or an aspirate.
25. If the first part of a conjunct is ह, the second part must be a र.
26. If the first part of a conjunct is न, the second part cannot be ल.
27. Glides do not occur before aspirates or breathy voiced stops.
28. Fricatives do not occur before breathy voiced stops, aspirates or affricates.
29. Nasals do not occur after breathy voiced stops, and stops that do occur after them must be voiceless.
30. Glides do not occur after affricates.

31. If in a two consonant cluster, if the first part is a ड, the second part cannot be a fricative ष श स ह or a stop.
32. Retroflex segments or affricates do not occur after ह.
33. If the second part of a consonant cluster is a ह, the second part cannot be aspirated or breathy-voiced.
34. If the first part of a consonant cluster is a य or व, the second consonant cannot be a glide.
35. य does not occur before or after retroflex glides.
36. Voiced fricative ज्ञ occurs only after voiced segments.
37. Retroflex segments, voiced stops, liquids, and nasals do not occur after ऋ.
38. Retroflexes, fricatives, or liquids do not occur before ऋ.
39. व does not occur before bilabials.
40. Fricatives and retroflexes do not occur after ज्ञ, and if the consonant following ज्ञ is a glide, it must be ह.
41. Bilabials do not occur before फ़, and if the segment before फ़ is a stop, it must be voiceless.
42. Retroflexes or य should not occur after न.
43. The breathy voiced affricates (छ and झ) do not occur as the first part of a consonant cluster. They will occur as the second part of a cluster only if the first part is a homorganic nasal.
44. फ़ does not occur as the first part of a conjunct.
45. In a conjunct having two stops one after the other, velars and bilabials do not occur as the second consonant.

46. Fricatives and stops do not occur after breathy voiced consonants.
47. In a conjunct, if the first part is a non-nasal retroflex, the second part can only be य.
48. After फ, only glides occur.
49. ह does not occur as the first part of a conjunct.
50. Only liquids or nasals occur before ह.
51. Liquids do not occur after न.

The constraints outlined in this section are applicable at specific locations within Hindi words: Constraints 1 through 15 are applicable to all consonant sequences, 16 through 26 are applicable to consonant sequences at the start of a word, 27 through 42 are applicable to consonant sequences in the middle of a word, 43 through 51 are for consonant sequences at the end of a word.

Appendix B

Manual for “ILT Indic Digital Library Toolkit”

The ILT¹ Indic Digital Library (IDL) toolkit, also called *Aakalan*, provides multi-lingual text entry and document image analysis (DIA) functions. The current version supports Devanagari Truth entry and script analysis.

The tool is distributed for windows and Linux. When presented on a CD, the tool, and all sample files are present in a directory named *truthingtool*. In case of a download, you will receive a single zip file, which contains the *truthingtool* directory. The first level of directory structure is shown below:

- bin/* Executable scripts.
- docs/* Documentation for the tool
- libs/* “jar” files needed for the system
- samples/* Sample files, and the truth data set
- TrainingData/* Images of query results are stored here

The tool may be launched by executable scripts located in “bin/jarRun.bat” for Windows and

¹URL: <http://www.cedar.buffalo.edu/ilt/>

“bin/jarRun” for Linux/Unix based systems. Once launched, you will see a window with two options (Figure B.1). The tool has a single window in which multiple interfaces can be opened at the same time. Note: All windows will be closed if you close the main window.

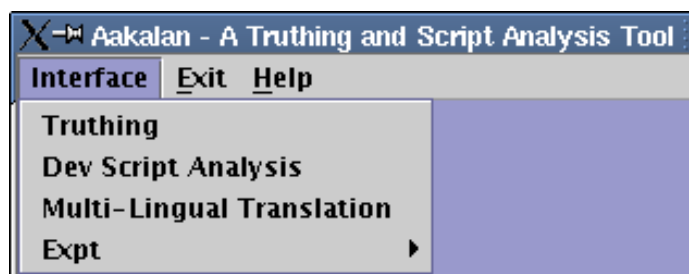


Figure B.1: Main window options

B.1 Truthing Interface

This version is distributed with the Devanagari truthing interface (Figure B.4), which may be opened by using the *Interface: Truthing* option from the Main window (Figure B.1).

B.1.1 Menu Options

Following options are provided by the truthing interface menu (Figure B.2):

File: Open: Image allows users to open an image. The image should be 300dpi gray-scale in tiff (tif), jpeg (jpg), gif, or png format. The tool performs text block identification and word

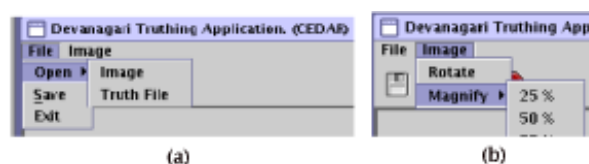


Figure B.2: Truthing Interface Menu: (a) File Options and (b) Image Options

separation on the document. Word images are displayed in the image window (Figure B.4a). Users enter word truth by using transliteration (Figure B.4b) or GUI keyboard (Figure B.4c).

File: Open: Truth File displays a dialog to open a truth file. The truth file has an extension of **imageName.truth.xml**, where **imageName** is the image corresponding to the truth file. DTD for the truth file is located in the “samples” directory.

File: Save displays a dialog to select a file to store the truth. The extension of the selected file should be *truth.xml*.

File: Exit closes the current window. If other windows are open in the main interface, they will not be affected. Note that data in the window is not saved automatically during exiting.

Image: Rotate opens a dialog to rotate the image in a clockwise direction. The angle mentioned in the interface is in degrees. Word boundaries are adjusted after rotating the image.

Image: Magnify magnifies the image to different levels.

B.1.2 Truthing Interface Tool bar Options

The tool bar (Figure B.4e) provides three utilities: The first (save) button is functionally similar to the *File: Save* menu. The second button is used to delete spurious characters and artifacts that are erroneously segmented as words. The third (exit) button is similar to the *File: Exit* menu option.

B.1.3 Image Panel

This holds the display for document images. Word boundaries are usually indicated by blue-colored demarcators. The word being truthed has a red colored boundary. Word bounds may be altered by using the left mouse button to select and drag word boundaries. Hovering the mouse over a word boundary displays word number and other optional information.

To create new words, users draw a rectangular region in the image panel with the left mouse

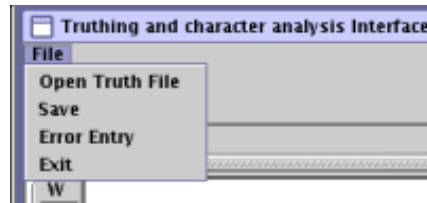


Figure B.3: Menu Options in the Devanagari Script Analysis Interface

button. The system computes a word number for the selected bounds, but in case it fails to determine reading order, the user is prompted to enter a word number.

B.1.4 Transliteration Panel

The transliteration panel (Figure B.4b) is a text editor with features to enable phonetic entry of Indic scripts. Indic script text may be entered using Latin phonetic equivalents. Clipboard features of cut (*control-x*), copy (*control-c*), and paste (*control-v*) are implemented into the system. As one browses through words in the transliteration window, the image and truth panels are updated to show corresponding words.

B.2 Devanagari Script Analysis Interface

The option *Interface: Dev Script Analysis* from the Main Window (Figure B.1) opens a Devanagari script analysis interface. It has controls to examine truth files and extract component images.

B.2.1 Menu Options

The *File: Open Truth File* shows a dialog to select a truth file. Truth files have the extension “truth.xml”.

The *File: Save* menu is similar in operation to the *File: Save* option in the Truthing Interface (Section B.1.1).

File: Error Entry opens a dialog to submit an error. Fields are available to enter line and word location of the error within a document. Text panes allow users to enter a brief explanation and the java stack trace. Errors are usually indicated by exceptions on the command prompt in Unix/Linux systems, or the debug window of jre in Windows. All fields of the error report creator (Figure B.9) are optional, though detailed information is desired to effectively deal with an error.

File: Exit option closes the interface.

B.2.2 Truthing Control

The truthing controls in this interface (Figure B.6) are functionally similar to elements of the Truthing Interface except that images cannot be opened: only existing truth files can be modified. The interface was designed primarily to view word truth and correct it if necessary.

B.2.3 Devanagari Script Analysis

Tab C (Figure B.7) displays Devanagari Script Analysis controls, which provide a tree view of the entire document. When users click on a word node, the underlying truth is parsed to generate different components and their truth (Figure B.7a, b, c). The tree view shows a word divided first into characters, then alphabets and components. This hierarchical division assists in empirical analysis of algorithms that segment words into characters and components.

B.2.4 Query Panel

The query panel allows users to extract different components from truth files. Image processing techniques are used to segment a word into its' semantic parts (components). The rules of Devanagari orthography allow us to predict where these components occur and how they fuse to build a character. Combining image processing techniques with rules of Devanagari orthography, words and word parts are extracted with varying levels of detail.

Option	Description	Examples
Combo Box	To select the text element level	Component <i>ka</i> क : क क क Alphabet <i>ka</i> क : क क क Character <i>kaa</i> का : का Word <i>kyaari</i> क्यारि : क्यारि
Substr	Option used to select substrings	Querying for words with a <i>ka</i> क gives: कारन कर्ता कर्म क्यारि
Half Char	Option to select character half-forms	Querying for characters with <i>half-ka</i> क gives: क्या
Compnd Qry	This is the button used to query for compound characters	

Table B.1: Query options: Examples are results of queries executed over the domain कारन कर्ता कर्म क्यारि

Possible query options and their utility is described in Table B.1. For special vowels like ँ, ए, and ऋ, user needs to enter respective transliteration and select the half-character option. Queries for vowel modifiers ँ(aa) ि(i) ि(ii) ` (e) ै (ai) े(o) ै(au) ू(u) ू(uu) are entered using the GUI keyboard. The upper parts of ि(i) and ि(ii) are obtained by selecting the option “half-char” from the keyboard.

B.3 Multi-Lingual Translation Interface

The Multi-Lingual Translation Interface is opened using the *Interface: Multi-Lingual Translation* option. The interface provides controls to view and annotate documents in Arabic, Devanagari, English and Telugu.

B.3.1 Menu Options

The File menu has following menu options:

File: New Translation blanks out the three text windows and sets up the Transliteration window (Figure B.12(I)) for user input.

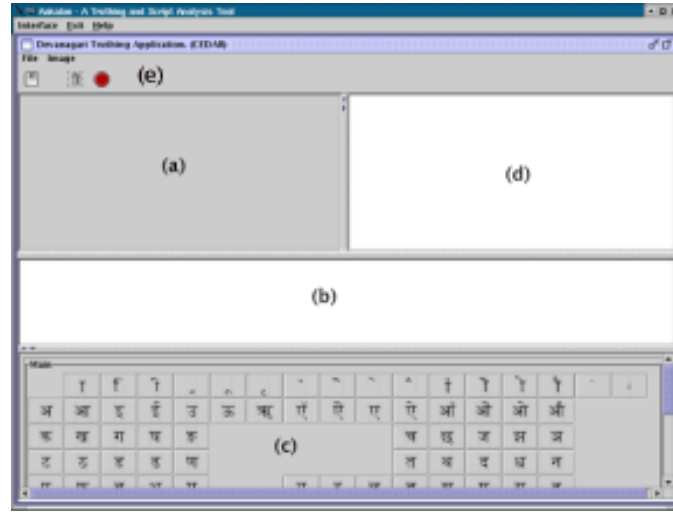


Figure B.4: Empty truthing tool: (a) Image Panel, (b) Transliteration Panel, (c) GUI Keyboard, (d) Truth Panel, (e) Tool bar and Menu options

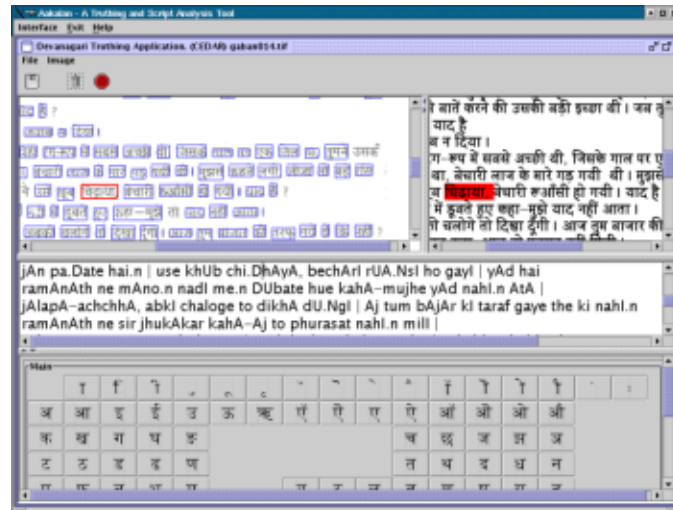


Figure B.5: Truthing interface: Please refer to Figure B.4 for legend. Red outline is for the word being truthed.

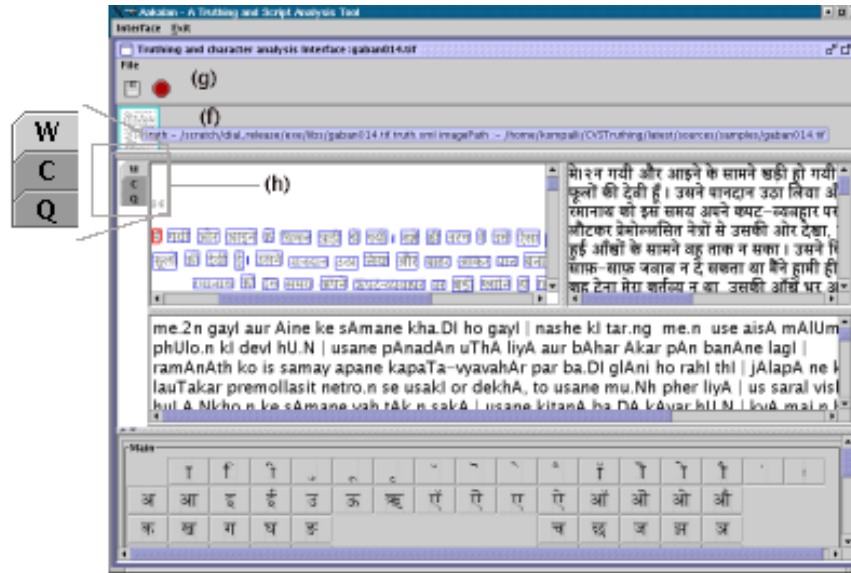


Figure B.6: Devanagari Script Analysis, Truthing Interface: (f) Thumbnail panel, (g) Tool bar, and (h) Tab pane; W: Word level truth, C: Character level truth, and Q: Query Panel.

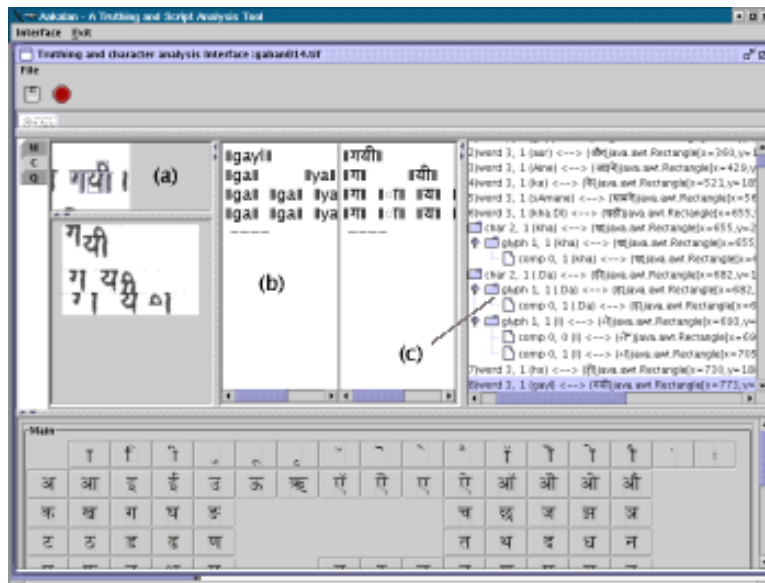


Figure B.7: Devanagari Script Analysis, Segmentation Analysis: (a) Images of the word and its' components, (c) Word and component truth, (d) Tree layout of the truth.

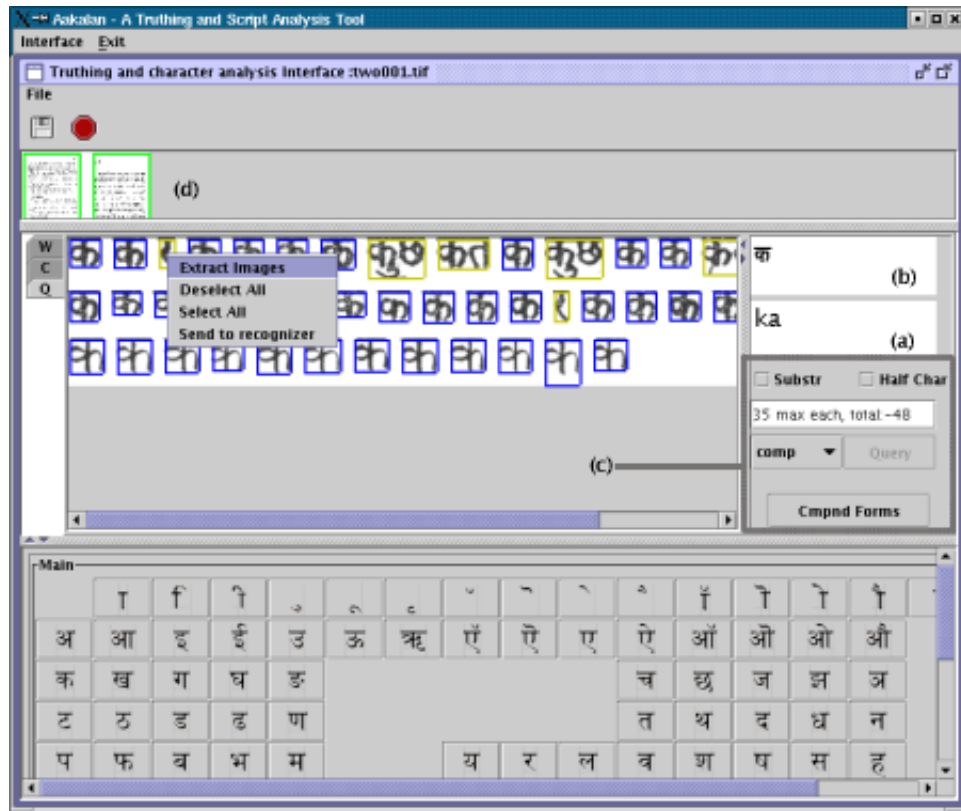


Figure B.8: Devanagari Script Analysis, Query Panel: The display shows a query for the component “ka” क : (a) Query transliteration input, (b) Query display, and (c) Query options

File: Open: Translation File opens a file selection dialog to select a readable translation file. Translation files have an extension of translation.xml. On opening the translation file, the Transliteration input, Annotation window and Script window are initialized with data read from the file, (Figure B.12(I, II, III)) and readied for user input.

File: Open: Truth File opens a file selection dialog to select a readable truth file. Truth files have an extension of truth.xml, and are associated with image information. Such files may be opened in read-only mode². On opening the truth file, the Transliteration input, Annotation window and Script window are initialized with data read from the file, (Figure B.12(I, II, and III)) but user input is disabled.

File: Save As: Translation File Opens a file selection dialog to save the displayed text as a translation file.

File: Save As: Html File Opens a file selection dialog to save the displayed text as a HTML file. Text in different scripts is represented using Unicode-encoded HTML, and Javascript is used to display annotations.

B.3.2 Text Controls

Text controls for the translation interface are shown in Figure B.12. The language of selected text in the transliteration window (Figure B.12(I)) may be altered using the language control (Figure B.11(a)). For ex: When “Devanagari” is selected, highlighted transliteration is parsed according to Devanagari transliteration rules and corresponding locations in the annotation and script windows (Figure B.12(II, III)) are updated.

The script type control (Figure B.11(b)) allows users to select regions for annotation - *Script*: This text is associated with translation/annotation and is highlighted in pink color (Figure B.12(a, f, i)). *Annotation*: When this option is selected, highlighted text is assigned as annotation for the

²Note for developers: This behavior should probably be modified since truth file data might need to be modified and written into translation files.

preceding *Script* (text with pink background). Translations/annotations are identified by a light-blue background (Figure B.12(b, g)). *Text*: Section of the script that is not associated with any annotation/translation. This type of text has a white/transparent background (Figure B.12(e, h)).

The translation file `howto3.translation.xml` in the translation samples provides an example to examine the features of this interface.

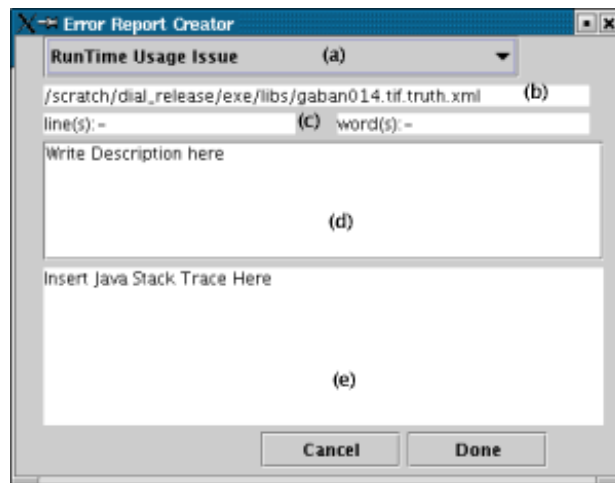


Figure B.9: Error Report Creator: (a) Error Type, (b) Truth File Name, (c) Line and Word Number, (d) Error Description, and (f) Stack Trace.

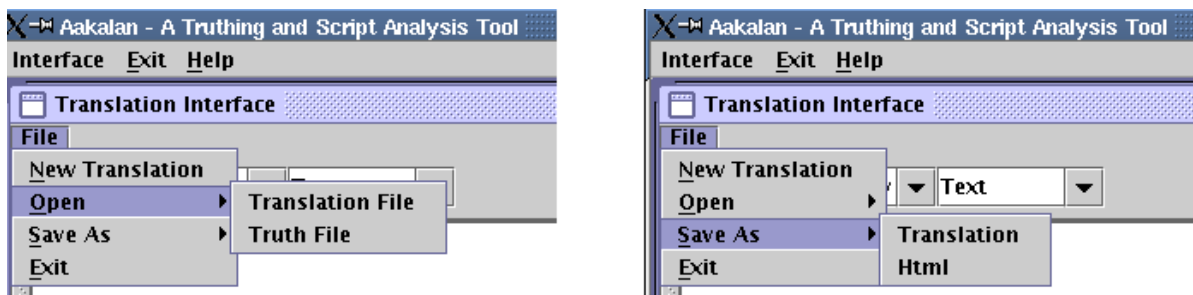


Figure B.10: Menu Options for the Multi-Lingual Translation Interface

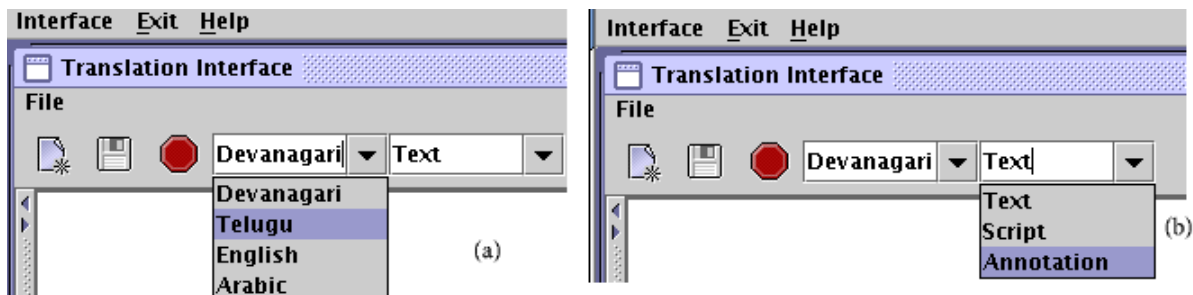


Figure B.11: Text Controls for the Multi-Lingual Translation Interface

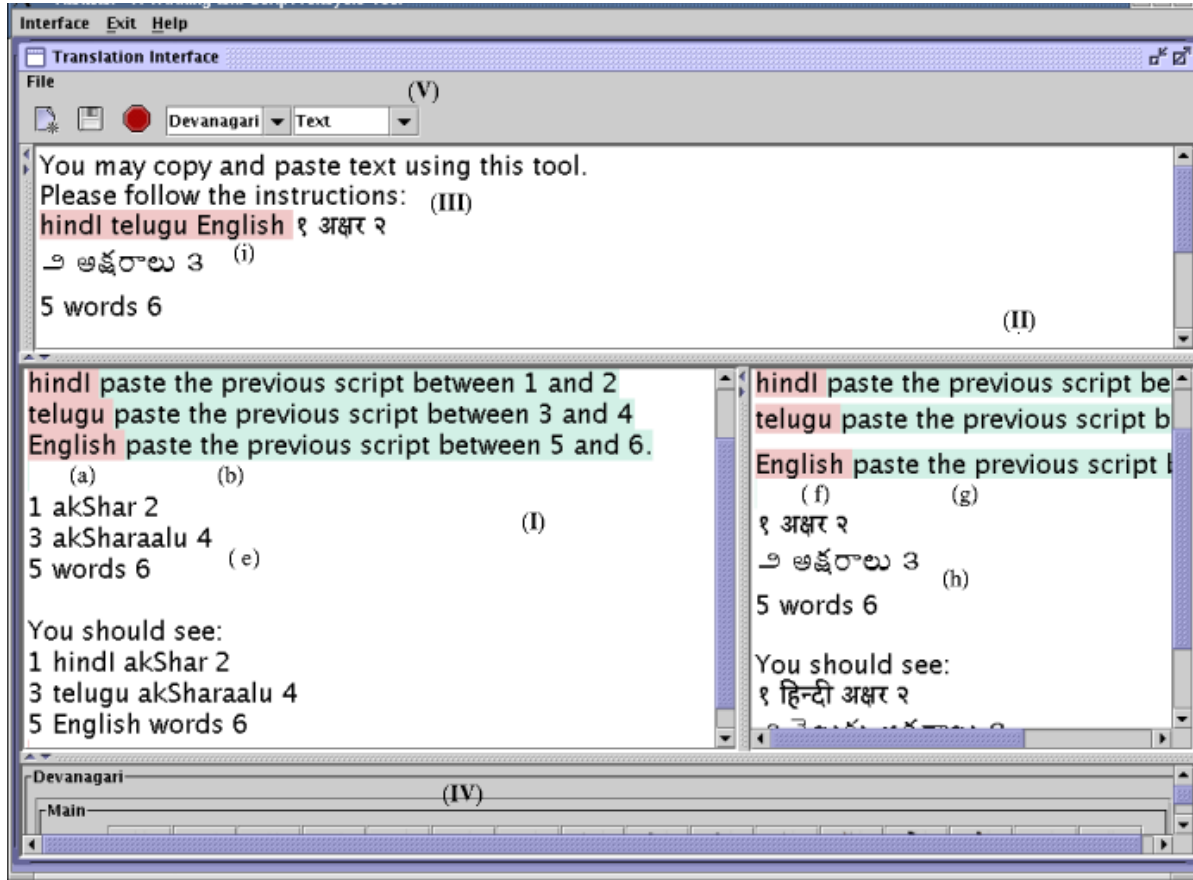


Figure B.12: Multi-Lingual Translation Interface: Multi-Lingual Translation Interface: (I) *Transliteration window* to enter transliterated English text for scripts, (II) *Annotation window* to view parsed text along with annotation, (III) *Script window* displays transliterated text without the associated annotations. (IV) GUI keyboard for text entry, (V) Text controls and tool bar. (a) Text to which annotation has been added, (b) Annotation for pink colored region, (c) Text without any annotation, (f,i,g,h) Corresponding parsed text.

Bibliography

- [1] Amr R. Abdel-Dayem, Ali K. Hamou, and Mahmoud R. El-Sakk. Novel adaptive filtering for salt-and-pepper noise removal from binary document images. *Lecture Notes in Computer Science*, 3212:191–199, 2004.
- [2] L. Alessandro, Robert Sabourin Koerich, and Y. Suen Ching. Lexicon-driven hmm decoding for large vocabulary handwriting recognition with multiple character models. *International Journal on Document Analysis and Recognition*, 6, 2003.
- [3] Robert B. Allen and John Schalow. Metadata and data structures for the historical newspaper digital library. In *Proceedings of the 8th international conference on Information and knowledge management*, pages 147–153, 1999.
- [4] Henry Baird and Tim Kam Ho. Large-scale simulation studies in image pattern recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(10):1067–1079, 1997.
- [5] P. Baker, A. Hardie, T. McEnery, R. Xiao, K. Bontcheva, H. Cunningham, R. Gaizauskas, O. Hamza, D. Maynard, V. Tablan, C. Ursu, B. Jayaram, and M. Leisher. Corpus linguistics and south asian languages: Corpus creation and tool development. *Literary and Linguistic Computing*, 19(4):509–524, 2004.
- [6] Veena Bansal. Integrating knowledge sources in Devanagari text recognition. *IEEE Transactions on Systems, Man and Cybernetics Part A*, 30(4):500–505, 2000.

- [7] Veena Bansal and R.M.K. Sinha. Partitioning and searching dictionary for correction of optically-read devanagari character strings. *International Journal on Document Analysis and Recognition*, 4(4):269–280, July 2002.
- [8] Veena Bansal and R.M.K Sinha. Segmentation of touching and fused devanagari characters. *Pattern Recognition*, 35:875–893, 2002.
- [9] D. Besagni and A. Belaid. Citation recognition for scientific publication in digital libraries. In *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries. (DIAL 2004)*, pages 244–252, 2004.
- [10] Derek Besner and Jennifer A Stolz. Unintentional reading: can phonological computation be controlled? *Canadian Journal of Experimental Psychology*, 52:35–42, 1998.
- [11] Christopher M. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1996.
- [12] Djamel Bouchaffra, Venu Govindaraju, and Sargur N. Srihari. Postprocessing of recognized strings using nonstationary markovian models. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 21(10):990–999, 1999.
- [13] Thomas M. Breuel, William C. Janssen, Kris Popat, and Henry Baird. Paper to pda. In *Proceedings of the 16th International Conference on Pattern Recognition (ICPR 2002)*, volume 1, pages 476–479, 2002.
- [14] Claudia Carello, M. T. Turvey, and Georgije Lukatela. Can theories of word recognition remain stubbornly nonphonological? In R. Frost and L. Katz, editors, *Orthography, Phonology, Morphology, and Meaning: An Overview*, pages 211–226. Elsevier publications, 1992. ISBN: 0444891400.

- [15] R.G. Casey and E. Lecolinet. A survey of methods and strategies in character segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18:690–706, July 1996.
- [16] R. Cattoni, R. Coianiz, S. Messelodi, , and C.M. Modena. Geometric layout analysis techniques for document image understanding: A review. Technical Report 9703-09, ITC-IRST, 1998.
- [17] Bidyut Baran Chaudhari and Umapada Pal. Skew angle detection of digitized indian script documents. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(2):182–186, February 1997.
- [18] B.B. Chaudhuri and U. Pal. An ocr system to read two Indian language scripts: Bangla and Devanagari. In *Proceedings of the 4th International Conference on Document Analysis and Recognition*, pages 1011–1015, 1997.
- [19] Ming Chen, Xiaoqing Ding, and Jian Liang. Analysis, understanding and representation of chinese newspaper with complex layout. In *Proceedings of the International Conference on Image Processing*, volume 2, pages 590–593, 2000.
- [20] Manning Christopher and Schutze Hinrich. *Foundations of Statistical Natural Language Processing*. MIT Press, 1999.
- [21] Scott Connell. A comparison of hidden markov model features for the recognition of cursive handwriting. Technical Report MSUCPS: TR96-31, Michigan State University, 1996.
- [22] T.H. Cormen, C.E., Leiserson, and R. L. Rivest. *Introduction to Algorithms*. The MIT Press, 1990.
- [23] Bamber D. The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *Journal of Mathematical Psychology*, 12:387–415, 1975.

- [24] Peter T. Daniels and William Bright. *The World's Writing Systems*. Oxford University Press, March 1996.
- [25] Mike Davenport and Stephen J. Hannahs. *Introducing Phonetics and Phonology*. Hodder Arnold, 2005.
- [26] Richard O. Duda, Peter E. Hart, , and David G. Stork. *Pattern Classification*. Wiley Interscience, 2 edition, 2000.
- [27] Emmanuel Dupoux, Christophe Pallier, Kazuhiko Kakehi, and Jacques Mehler. New evidence for prelexical phonological processing in word recognition. In *Workshop On Spoken Word Access Processes*, volume 16, pages 491–505, 2001.
- [28] James P. Egan. *Signal detection theory and ROC-analysis*. New York: Academic Press, 1975.
- [29] F. Esposito, D. Malerba, G. Semeraro, S. Ferilli, O. Altamura, T. M. A. Basile, M. Berardi, M. Ceci, and N. D. Mauro. Machine learning methods for automatically processing historical documents: From paper acquisition to xml transformation. In *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries. (DIAL 2004)*, pages 328–335, 2004.
- [30] J 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.
- [31] J.T. Favata. Off-line general handwritten word recognition using an approximate beam matching algorithm. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23:1009–1021, September 2001.

- [32] M.L. Forcada. Corpus-based stochastic finite-state predictive text entry for reduced keyboards: application to catalan. In *Procesamiento del Lenguaje Natural*, pages 65–70, 2001.
- [33] G Forney. The viterbi algorithm. *Proceedings of the IEEE*, 61:268–277, March 1977.
- [34] Herbert Freeman. Computer processing of line-drawing images. *Computing Surveys*, 6(1):57–97, March 1974.
- [35] K Fukushima and T Imagawa. Recognition and segmentation of connected characters with selective attention. *Neural Networks*, 6(1):33–41, 1993.
- [36] U. Garain and B.B. Chaudhuri. Segmentation of touching characters in printed devnagari and bangla scripts using fuzzy multifactorial analysis. *IEEE Transactions on Systems, Man and Cybernetics, Part C*, 32(4):449–459, November 2002.
- [37] John Goldsmith, editor. *The Handbook of Phonological Theory*. Blackwell, 1996.
- [38] Venu Govindaraju, Swapnil Khedekar, Suryaprakash Kompalli, Faisal Farooq, Srirangaraj Setlur, and Ramanaprasad Vemulapati. Tools for enabling digital access to multilingual indic documents. In *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries. (DIAL 2004)*, pages 122–133, 2004.
- [39] I. Guyon, R. Haralick, J. Hull, and I. Phillips. Data sets for ocr and document image understanding research. *Handbook of Character Recognition and Document Image Analysis*, pages 779–799, 1997.
- [40] Venu Govindaraju Hanhong Xue. On the dependence of handwritten word recognizers on lexicons. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(12):1553–1564, 2002.

- [41] M.W. Harm and M.S. Seidenberg. Computing the meanings of words in reading: Division of labor between visual and phonological processes. *Psychological review*, 111:662–720, 2004.
- [42] Teemu Hirsimäki, Mathias Creutz, Vesa Siivola, and Kurimo Mikko. Morphologically motivated language models in speech recognition. In *Proceedings of International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning*, pages 121–126, 2005.
- [43] Jonathan J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, 1994.
- [44] Jonathan J Hull and Sargur N. Srihari. Experiments in text recognition with binary n-grams and viterbi algorithms. *IEEE Transactions on Pattern Analysis Machine Intelligence*, 4(5):520–530, 1982.
- [45] Daisy L. Hung, Ovid J. L. Tzeng, and Angels K. Y. Tzeng. Automatic activation of linguistic information in chinese character recognition. In R. Frost and L. Katz, editors, *Orthography, Phonology, Morphology, and Meaning: An Overview*, pages 119–130. Elsevier publications, 1992. ISBN: 0444891400.
- [46] R. Hunter and A. H. Robinson. International digital facsimile coding standards. *Proceedings of the IEEE*, 68(7):854–867, July 1980.
- [47] L. Hutchison and W. A. Barrett. Fast registration of tabular document images using fourier analysis. In *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries. (DIAL 2004)*, pages 253–267, 2004.
- [48] International Phonetic Association, 1999.
- [49] Anil Jain. *Fundamentals of Digital Image Processing*. Prentice-Hall, 1986.

- [50] Amengual Juan C and Vidal Enrique. Efficient error-correcting viterbi parsing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:1109–1116, October 1998.
- [51] Perez-Cortes Juan C, Amengual Juan C, Arlandis Joaquim, and Llobet Rafael. Stochastic error-correcting parsing for ocr post-processing. In *Proceedings of the 15th International Conference on Pattern Recognition*, volume 4, pages 405–408, 2000.
- [52] Tapas Kanungo, Robert Haralick, and Henry Baird. Validation and estimation of document degradation models. *Proceedings of 4th UNLV Symp. on Document Analysis & Information Retrieval*, pages 217–225, 1995.
- [53] Tapas Kanungo and Song Mao. Stochastic language models for style-directed layout analysis of document images. *IEEE Transactions on Image Processing*, 12(5):583–596, 2003.
- [54] E. Kavallieratou, N. Fakotakis, and G. Kokkinakis. Skew angle estimation in document processing using cohen’s class distributions. *Pattern Recognition Letters*, 20(11-13):1305–1311, November 1999.
- [55] Swapnil Khedekar, Ramanaprasad Vemulapati, Srirangaraj Setlur, and Venu Govindaraju. Text - image separation in devanagari documents. In *Proceedings of the 7th International Conference on Document Analysis and Recognition*, pages 1265–1269, 2003.
- [56] D.H. Kim, Y.S. Hwang, S.T. Park, E. Kim, S. Paek, and S. Bang. Handwritten korean character image database. In *Proceedings of the 2nd International Conference on Document Analysis and Recognition*, pages 470–473, 1993.
- [57] G Kim and Venu Govindaraju. A lexicon driven approach to handwritten word recognition for real-time applications. *IEEE Pattern Analysis and Machine Intelligence*, 19:366 – 379, 1997.

- [58] Suryaprakash Kompalli, Srirangaraj Setlur, Venugopal Govindaraju, and Ramanaprasad Vemulapati. Creation of data resources and design of an evaluation test bed for devanagari script recognition. In *Proceedings of the 13th International Workshop on Research Issues on Data Engineering: Multi-lingual Information Management*, pages 55–61, 2003.
- [59] Karen Kukich. Techniques for automatically correcting words in text. *ACM Computing Surveys*, 24(4):377–439, 1992.
- [60] N. Lachiche and P.A. Flach. Improving accuracy and cost of two-class and multi-class probabilistic classifiers using roc curves. In *Proceedings of the 20th International Conference on Machine Learning*, pages 416–423, January 2003.
- [61] V. Lavrenko, T. Rath, and R. Manmatha. Holistic word recognition for handwritten historical documents. In *Proceedings of the 1st International Workshop on Document Image Analysis for Libraries. (DIAL 2004)*, pages 278–287, 2004.
- [62] Seong-Whan Lee, Dong-June Lee, and Hee-Seon Park. A new methodology for gray-scale character segmentation and recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18:1045–1050, October 1996.
- [63] Vladimir I. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10(8):707–710, 1966.
- [64] Huanfeng Ma and David Doermann. Adaptive hindi ocr using generalized hausdorff image comparison. *ACM Transactions on Asian Language Information Processing*, 26(2):198–213, 2003.
- [65] S. Madhvanath, G. Kim, and V. Govindaraju. Chain code processing for handwritten word recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 21(9):928–932, 1999.

- [66] Sriganesh Madhvanath and Venu Govindaraju. The role of holistic paradigms in handwritten word recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(2):149–164, February 2001.
- [67] Muhammed Melhi, Stanley S. Ipson, and William Booth. A novel triangulation procedure for thinning hand-written text. *Pattern Recognition Letters*, 22(10):1059–1071, August 2001.
- [68] Tom M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [69] Shunji Mori, Ching Y. Suen, and Kazuhiko Yamamoto. Historical review of ocr research and development. *Proceedings of the IEEE*, 80:1029–1058, July 1992.
- [70] Manjari Ohala. *Aspects of Hindi Phonology*. Motilal Banarasidas, 1983. ISBN: 0895811162.
- [71] U Pal and Chaudhuri B. B. Indian script character recognition: A survey. *Pattern Recognition*, 37:1887–1899, September 2004.
- [72] Umapada Pal, P.K. Kundu, and B. B. Chaudhuri. Ocr error correction of an inflectional indian language using morphological parsing. *Journal of Information Science and Engineering*, 16:903–922, November 2000.
- [73] Theo Pavlidis. Contour filling in raster graphics. In *Proceedings of the 8th annual conference on Computer graphics and interactive techniques*, pages 29–36, 1981.
- [74] Kaushalendra Prapann. *India Today Hindi Edition*. The India Toda Group, april 2002.
- [75] Munshi Premchand. *Gaban*. Hansa Prakaashan, 1960.
- [76] Frost Ram and Bentin Shlomo. Reading consonants and guessing vowels: Visual word recognition in hebrew orthography. In R. Frost and L. Katz, editors, *Orthography, Phonol-*

- ogy, *Morphology, and Meaning: An Overview*, pages 27–44. Elsevier publications, 1992. ISBN: 0444891400.
- [77] Mingwu Ren, Wankou Yang, and Jingyu Yang. A new and fast contour-filling algorithm. *Pattern Recognition*, 38:2564–2577, December 2005.
- [78] Jairo Rocha and Theo Pavlidis. Character recognition without segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17:903 – 909, September 1995.
- [79] Azriel Rosenfeld and Avinash C. Kak. *Digital Picture Processing*. Academic Press, 2 edition, 1982.
- [80] Ronald Rosenfeld. A maximum entropy approach to adaptive statistical language modeling. *Computer, Speech and Language*, 10:187–228, 1996.
- [81] Lehal G. S., Chandan Singh, and Ritu Lehal. A shape based post processor for gurmukhi ocr. In *Proceedings of the 6th International Conference on Document Analysis and Recognition*, pages 1105–1109, 2001.
- [82] Daming Shi, Robert I. Damper, and Steve R. Gunn. Offline handwritten chinese character recognition by radical decomposition. *ACM Transactions on Asian Language Information Processing*, 2:27 – 48, 2003.
- [83] Steven J. Simske and Margaret Sturgill. A ground-truthing engine for proofsetting, publishing, re-purposing and quality assurance. In *Proceedings of the 2003 ACM symposium on Document engineering*, pages 150 – 152, 2003.
- [84] R.M.K. Sinha. Rule based contextual post-processing for devanagari text recognition. *Pattern Recognition*, 20:475–485, 1987.
- [85] R.M.K. Sinha and H. Mahabala. Machine recognition of devnagari script. *IEEE Transactions on Systems, Man, and Cybernetics*, 9:435–441, 1979.

- [86] R.M.K. Sinha, B. Prasada, G.F. Houle, and M. Sabourin. Hybrid contextual text recognition with string matching. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15:915–925, September 1993.
- [87] Peter Slavik and Venu Govindaraju. An overview of run-length encoding of handwritten word images. Technical Report 9, Department of Computer Science and Engineering, University at Buffalo, 2000.
- [88] Jiqiang Song, Zuo Li, M.R. Lyu, and Shijie Cai. Recognition of merged characters based on forepart prediction, necessity-sufficiency matching, and character-adaptive masking. *IEEE Transactions on Systems, Man and Cybernetics, Part B*, 35:2–11, February 2005.
- [89] Milan Sonka, Vaclav Hlavac, and Roger Boyle. *Image Processing, Analysis and Machine Vision*. Brooks/Cole Publishing, 2 edition, 1999.
- [90] John A. Spinks, Ying Liu, Charles A Perfetti, and Li Hai Tan. Reading chinese characters for meaning: the role of phonological information. *Cognition*, 76:B1–B11, 2000.
- [91] Richard Sproat. *A Computational Theory of Writing Systems*. Cambridge University Press, September 2000.
- [92] G Srikantan, S. Lam, and Srihari Sargur. Gradient based contour encoding for character recognition. *Pattern Recognition*, 29(7):1147–1160, 1996.
- [93] John Ridley Stroop. Studies of interference in serial verbal reactions. *Journal of Experimental Psychology*, 18:643–662, 1935.
- [94] C.Y. Suen, S. Mori, S.H. Kim, and C.H. Leung. Analysis and recognition of asian scripts-the state of the art. In *Proceedings of the 7th International Conference on Document Analysis and Recognition*, pages 866–878, 2003.

- [95] S Suzuki, N Uedan, and J Sklansky. Graph-based thinning for binary images. *International journal of pattern recognition and artificial intelligence*, 7(5):1009–1030, 1993.
- [96] Yuan Y. Tang, Lo-Ting Tu, Jiming Liu, Seong-Whan Lee, and Win-Win Lin. Off-line recognition of chinese handwriting by multifeature and multilevel classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20:556–561, 1998.
- [97] Insup Taylor and M. Martin Taylor. *The Psychology of Reading*. Academic Press, 1983. ISBN: 0126840806.
- [98] Oivind Due Trier and Anil K. Jain. Goal-directed evaluation of binarization methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 17(12), 1995.
- [99] Oivind Due Trier, Anil K Jain, and Torfinn Taxt. Feature extraction methods for character recognition a survey. *Pattern recognition*, 29(4):641–662, 1996.
- [100] Jyotsna Vaid and Ashum Gupta. Exploring word recognition in a semi-alphabetic script: The case of devanagari. *Brain and Language*, 81:679–690, 2002.
- [101] Hanhong Xue Venu Govindaraju, Petr Slavk. Use of lexicon density in evaluating word recognizers. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 24(6):789–800, 2002.
- [102] An-Bang Wang, Jun S. Huang, and Kuo-Chin Fan. Optical recognition of handwritten chinese characters by partial matching. In *Proceedings of the 2nd International Conference on Document Analysis and Recognition*, 1993.
- [103] C. J. Wells, L. J. Evett, P. E. Whitby, and R. J. Whitrow. Fast dictionary look-up for contextual word recognition. *Pattern Recognition*, 23:501 – 508, 1990.

- [104] Kim Jong Woo and Thoma R. George. Automated labeling in document images. In *Proceedings of the SPIE, Document Recognition and Retrieval VIII*, volume 4307, pages 111–122, January 2001.
- [105] Youshou Wu Xiaofei Huang, Jun Gu. A constrained approach to multifont chinese character recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(8):838–843, 1993.
- [106] Hanhong Xue. *Stochastic Modeling of High-Level Structures in Handwriting Recognition*. PhD thesis, University at Buffalo, The State University of New York, 2002.
- [107] Yaguang Yang, Kristen Summers, and Mark Turner. A text image enhancement system based on segmentation and classification methods. In *Proceedings of the 1st ACM workshop on Hardcopy document processing*, pages 33–40, 2004.
- [108] Zheng Yefeng, Li Huiping, and David Doermann. Machine printed text and handwriting identification in noisy document images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26:337–353, March 2004.
- [109] Daekeun You and Gyeonghwan Kim. An approach for locating segmentation points of handwritten digit strings using a neural network. In *Proceedings of the 7th International Conference on Document Analysis and Recognition*, volume 1, pages 142–146, 2003.
- [110] Bin Yu and A.K. Jain. A generic system for form dropout. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18:1127–1134, November 1996.
- [111] Bin Zhang and Sargur N. Srihari. Binary vector dissimilarities for handwriting identification. In *Proceedings of the 10th International SPIE Conference on Document Recognition and Retrieval*, pages 155–166, January 2003.

- [112] Jing Zheng, Xiaoqing Ding, and Youshou Wu. Recognizing on-line handwritten chinese character via farg matching. In *Proceedings of the 4th International Conference on Document Analysis and Recognition*, volume 2, pages 621–624, August 1997.