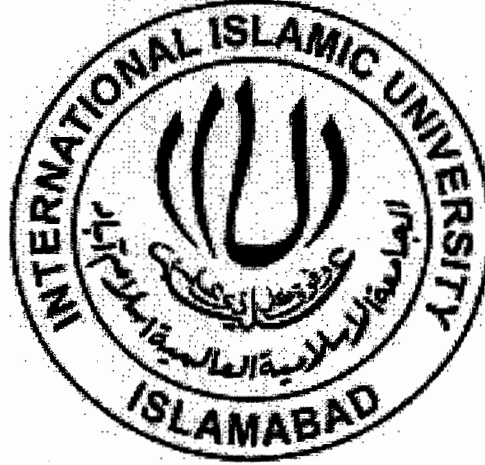


**ONLINE URDU CHARACTER RECOGNITION IN
UNCONSTRAINED ENVIRONMENT**

T08161



DATA ENTERED

Researcher:

Muhammad Imran Razzak

Reg. No. 36-FBAS/PHDCS/F07

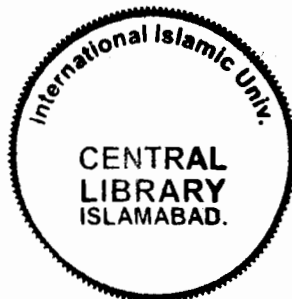
Supervisor:

Prof. Dr. Syed Afaq Hussain

Prof. Dr. Muhammad Sher

**Department of Computer Science,
Faculty of Basic and Applied Sciences,
INTERNATIONAL ISLAMIC UNIVERSITY,
ISLAMABAD,**

2011



DATA ENTERED

Accession No

TH861

M. ^{sil}
Md

PhD

006. 424

RAO

1. optical character recognition

ONLINE URDU CHARACTER RECOGNITION IN UNCONSTRAINED ENVIRONMENT



Muhammad Imran Razzak

36-FBAS/PHDCS/F07

Submitted in partial fulfillment of the requirement for the degree of Degree of Philosophy
in Computer Science at Faculty of Basic and Applied Sciences,
International Islamic University,
Islamabad

Prof. Dr Syed Afaq Hussain

Prof. Dr. Muhammad Sher

June, 2011

APPROVAL

Title of Thesis: Online Urdu Character Recognition in Unconstrained Environment

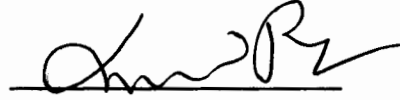
Name of Student: Muhammad Imran Razzak

Registration No: 36-FBAS/PHDCS/F07

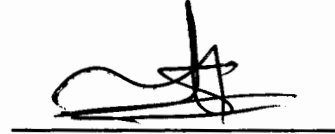
Accepted by the Department of Computer Science, INTERNATIONAL ISLAMIC UNIVERSITY, ISLAMABAD, in partial fulfillment of the requirements for the degree of Doctor of Philosophy in Computer Science.

Viva Voce Committee

Prof. Dr. Muhammad Riaz
Dean, Faculty of Basic and Applied Sciences
International Islamic University, Islamabad



Prof. Dr. Muhammad Sher (Co-Supervisor)
Chairman, Department of Computer Science
International Islamic University, Islamabad



Prof. Dr. Syed Afaq Hussain (Supervisor)
Faculty of Computing,
Riphah International University, Islamabad



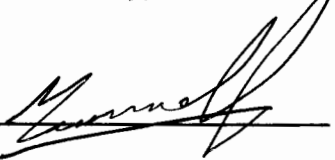
Dr. Ayaz Hussain (Internal Examiner)
Department of Computer Science
International Islamic University, Islamabad



Dr. M. Sikandar Hayat Khayal (External Examiner-I)
Chairperson, Department of Computer Science & SE
Fatima Jinnah Women University, Rawalpindi



Dr. Hammad Qureshi (External Examiner-II)
School of Electrical Engineering and Computer Science
National University of Sciences and Technology, Islamabad



Thursday, 9th June, 2011

ABSTRACT

Computer, the humongous giant of technology, has brought innovative changes in every aspect of life, especially in applications imitating humans. Currently, it is used in every field of life to facilitate human endeavor. One such application is character recognition. Character recognition is an important offshoot of pattern recognition problems. It imitates a human's ability to read, using a machine. It has been a field of intensive, if exotic, research since the early days of the computer. This task becomes more complex and demanding in case of handwritten and cursive text. Arabic script-based languages, which are used by almost a quarter of the world's population [Belaid et.al, 2010], are cursive, rich in diacritical marks and variety of writing styles present a challenging task for the researchers. Urdu is an Arabic script based languages however the Urdu character set is the superset of all Arabic script-based languages. Character recognition has been performed either through segmentation free or segmentation based approaches. There are numerous issues with a segmentation free approach, and it is very difficult to train using a large dataset. On the other hand in Urdu, a segmentation based approach has a large overhead and has less accuracy for cursive script as compared to segmentation free methods. In terms of classification, this thesis presents two approaches for Urdu character recognition: segmentation free method based on a hybrid approach (HMM and fuzzy logic), and bio-inspired character recognition system that uses fuzzy logics. Fuzzy is used as inner and outer shells for preprocessing and post processing of HMM. Biologically inspired multilayered fuzzy rules based system has been presented. Using the human visual concept, a layered approach has been suggested where the diacritical marks are separated from the ghost characters and mapped onto the primary ligature in the final layer. The proposed technique also caters to Multilanguage character recognition system for all Arabic script-based languages like Arabic, Persian, Urdu, Punjabi etc. The presented multilayered bio-inspired approach recognizes the ligature by extracting the features and combining them to find new premises in a bottom up fashion and it provided accuracy of 87.4%.

DECLARATION

I, **Muhammad Imran Razzak**, Registration No: **36-FBAS/PHDCS/F07**, hereby declare that this thesis titled “**Online Urdu Character Recognition in Unconstrained Environment**” for the fulfill of requirements of Doctor of Philosophy in Computer Science submitted to Department of Computer Science, International Islamic University, Islamabad, Pakistan is my own work and has not copied from any other resource. It is further declared that I have conducted this research and have accomplished this thesis entirely on the basis of my personal efforts and under the sincere guidance of my supervisors.

Muhammad Imran Razzak
36-FBAS/PHDCS/F07

To My Beloved Parents.

ACKNOWLEDGEMENTS

First of all, I would like to extend special thanks to my supervisors: Prof. Syed Afaq Husain, Professor, Faculty of Computing, Riphah International University and Prof. Muhammad Sher, Chairman, Depart of Computer Science, International Islamic University for giving me their golden time and helping in research. Secondly, I would like to thanks Professor Fateh Muhammad Malik, Rector International Islamic University and Dr Attash Durani, Director, Center of Excellence for Urdu Informatics, National Language Authority, for guiding me to develop multi-language character recognition system. I am also grateful to Prof. Abdel Belaid, READ, LORIA Lab, France for his kind guidance in writing research. I would also like to thank Prof. Rubiyah Yosuf, Director of Center of Artificial Intelligence and Robotics, University of technology, Malaysia for helping me in research during my stay at CAIRO University of Technology, Malaysia. I would like to pay special thanks to Ghulam Rasool Tahir for his encouragement and moral support during my PhD.

I would also like to thank my PhD colleagues, Mr. Aneel Rahim, Mr. Shoaib Ishaq, Mr. Ifikhar Watto for their encouragement and assistance in research work.

Last but not least, my younger brother Usman Razzak for his loving support during my PhD, he had to deal with stressed PhD students. Anyhow, he did well.

TABLE OF CONTENTS

ACKNOWLEDGEMENT	VI
LIST OF FIGURES	X
LIST OF TABLES	XIV
LIST OF ABBREVIATIONS	XV
1. INTRODUCTION	1
1.1 Motivation ..	4
1.2 Classification of Character Recognition	5
1.3 Constrained Vs Unconstrained	8
1.4 Urdu/Arabic Character Recognition Process	9
1.5 Challenges to Urdu Script Character Recognition	10
1.6 Objectives	13
1.7 Knowledge Based Urdu Character Recognition System	14
1.8 Contributions	15
1.9 Thesis Organization	16
2. HANDWRITTEN CHARACTER RECOGNITION:LITERATURE SURVEY	21
2.1 Preprocessing	24

2.2 Feature Extraction	33
2.3 Classifier	37
2.4 Final Remarks	43
3. FUZZY LOGIC AND HUMAN BIOLOGICAL VISUAL PERCEPTION.....	45
3.1 Fuzzy Logics	47
3.2.1 Linguistic Variables.....	49
3.3 Human Visual Perception	51
3.2.1 Autonomy of Human Visual System	53
3.3 Bio-Inspired Image Processing	55
3.4 Bio-Inspired Fuzzy Logics System	59
3.5 Final Remarks	60
4. GHOST CHARACTER RECOGNITION THEORY.....	62
3.1 Urdu Script Based Languages	63
3.2 Arabic Script Based Languages Character Recognition.	67
4.3 Ghost Character Theory	67
4.4 Ghost Character Recognition Theory	69
4.5 Effect of Ghost Character Theory	73
4.6 Merits	75
4.7 Demerits	75
4.8 Final Remarks	77
5. PREPROCESSING AND FEATURE EXTRACTION	78
5.1 Preprocessing Methodology	79

5.1.1 De-Hooking.....	81
5.1.2 Smoothing.....	82
5.1.3 Interpolation.....	83
5.1.4 Secondary Stroke Separation.....	84
5.1.5 Slant Correction.....	93
5.1.6 Stroke Mapping.....	96
5.1.7 Baseline Estimation.....	98
5.2 Feature Extraction ..	103
5.2.1 Directional Features ..	107
5.2.1 Structural Features ..	109
5.3 Final Remarks ..	113
6. BIOLOGICAL INSPIRED URDU CHARACTER RECOGNITION	117
6.1 Multilayered Human Visual System ..	120
6.2 Multilayered Character Recognition ..	123
6.2.1 Level-0: Stroke Acquisition ..	128
6.2.2 Level-1: Low level processing ..	128
6.2.3 Level-2: Geometrical processing ..	129
6.2.4 Level-3: Simple Feature Extraction ..	131
6.2.5 Level-4: Complex Feature Extraction ..	134
6.2.6 Level-5: Character Formation ..	134
6.2.7 Level-6: Word Formation ..	137
6.3 Discussion ..	138
6.4 Final Remarks ..	139

7. HOLISTIC APPROACH FOR URDU CHARACTER RECOGNITION USING	
MODIFIED HMM.....	141
7.1 Modified HMM Based Urdu Classification ..	142
7.2 HMM And Fuzzy Logics: A Hybrid Approach ..	149
7.2.1 Phase-I Clustering ..	151
7.2.2 Phase-II Classifier- I HMM ..	153
7.2.3 Phase-III Classifier- II Fuzzy Rules ..	155
7.3 Post Processing ..	156
7.4 Final Remarks ..	158
8. COMPRATIVE ANALYSIS AND CONCLUSIONS	162
8.1 Ghost Character Recognition Theory ..	163
8.2 Feature Extraction and Feature Fusion ..	166
8.3 Segmentation Free and Segmentation Approach ..	167
8.4 Conclusions ..	171
8.5 Future Work ..	174
REFERENCES.....	176
APPENDIX A PUBLICATION LIST	190

LIST OF FIGURES

<i>Figure 1.1 World vision in order to facilitate students using iPod.....</i>	<i>5</i>
<i>Figure 1.2 Character recognition methodologies</i>	<i>7</i>
<i>Figure 1.3 Digital pen using Anita functionality.</i>	<i>9</i>
<i>Figure 1.4.a Discrete writing</i>	<i>10</i>
<i>Figure 1.4.b Cursive writing</i>	<i>10</i>
<i>Figure 1.5 Character recognition process</i>	<i>11</i>
<i>Figure 1.6.a Different shapes of Nasta'liq script.....</i>	<i>12</i>
<i>Figure 1.6.b Diacritical marks complexity</i>	<i>12</i>
<i>Figure 2.1 Baseline estimation using histogram</i>	<i>25</i>
<i>Figure 2.2 Diacritical marks handling</i>	<i>31</i>
<i>Figure 2.3 Vertical Projection of Diacritical marks</i>	<i>31</i>
<i>Figure 2.4 Issues diacritical marks separation</i>	<i>32</i>
<i>Figure 2.5 Arabic ligature and their constituent character</i>	<i>33</i>
<i>Figure 2.6 Segmentation point estimation based on histogram</i>	<i>35</i>
<i>Figure 2.7 HMM modal for digits 4.....</i>	<i>39</i>
<i>Figure 3.1 Fuzzy Membership function</i>	<i>49</i>
<i>Figure 3.2.a Stimulus diagram for frog. prey acquisition behavior</i>	<i>53</i>
<i>Figure 3.2.b Detouring around the barrier in reaching stimuli</i>	<i>53</i>
<i>Figure 3.3 The autonomy of human visual object recognition</i>	<i>54</i>

List of Figures

Figure 3.4 Biological perception mechanism of human mind	55
Figure 3.5 Visual ventral stream based object recognition	56
Figure 3.6 Bio-inspired Numeral recognition	57
Figure 4.1 Different shapes wrt position from left to right isolated, start, mid, end	62
Figure 4.2.a Arabic Alphabets	62
Figure 4.2.b Persian Alphabets	62
Figure 4.3.a Urdu Alphabets	64
Figure 4.3.b Punjabi Alphabets Shahmukhi	64
Figure 4.4.a Sindhi Alphabets	66
Figure 4.4.b Pashto Alphabets	66
Figure 4.5.a Convergence of four dots to "Tota"	69
Figure 4.5.b Additional shapes in Urdu and Persian	69
Figure 4.6 Ghost characters used in Arabic script based language	67
Figure 4.7 Recognition of 2 nd ghost character letter with associated dot	71
Figure 4.8 Urdu Samples in three different styles. Urdu Nasta'liq, Nasq, Naskh	72
Figure 4.9 Different shapes of "ب" in Nasta'liq wrt neighbor character	70
Figure 4.10 Feature comparison of Nasta'liq and Naskh	71
Figure 4.11.a Mapping of diacritical marks and ghost character	73
Figure 4.11.b Word level recognition HMM for same word in Urdu and Arabic	73
Figure 4.12 Combination of diacritical marks with respect to languages	73
Figure 5.1 Urdu Samples in three different styles. Nasta'liq, Nasq, Naskh	79
Figure 5.2 Before Preprocessing issue (a) ل (b) ل 82	82
Figure 5.3 Missing points due to fast speed of writing	83

<i>Figure 5.4 Smoothed and Interpolated stroke.....</i>	<i>84</i>
<i>Figure 5.5 Secondary stroke separation</i>	<i>85</i>
<i>Figure 5.6 Membership function for Table 5.1 and 5.2</i>	<i>85</i>
<i>Figure 5.7 Smoothed and Interpolated stroke.....</i>	<i>86</i>
<i>Figure 5.8 Secondary stroke separation</i>	<i>87</i>
<i>Figure 5.9 Diacritical marks issue with three dots in Naskh</i>	<i>90</i>
<i>Figure 5.10 Fuzzy Membership function for Table 5.1 and 5.2</i>	<i>91</i>
<i>Figure 5.11 Recognition rate of diacritical marks</i>	<i>92</i>
<i>Figure 5.12 Left slanted, right slanted and normal words.....</i>	<i>95</i>
<i>Figure 5.13 Locally angle computation for Slant normalization</i>	<i>96</i>
<i>Figure 5.14 Base storks that are difficult to write in single stroke</i>	<i>97</i>
<i>Figure 5.15.a Resolved problem shown in figure 5.14.</i>	<i>97</i>
<i>Figure 5.15.b Error in combining</i>	<i>97</i>
<i>Figure 5.16 Baseline and Descender lines for Nasta'liq and Naskh for Urdu</i>	<i>90</i>
<i>Figure 5.17 Baseline for Nasta'liq and Naskh and blue line show baseline issues ...</i>	<i>90</i>
<i>Figure 5.18.a Raw input strokes.</i>	<i>101</i>
<i>Figure 5.18.b Ghost shapes after separation of secondary strokes.</i>	<i>101</i>
<i>Figure 5.18.c Primary baseline estimation based on projection.</i>	<i>101</i>
<i>Figure 5.18.c Locally baseline estimation based on features.</i>	<i>101</i>
<i>Figure 5.19 Features for baseline estimation</i>	<i>102</i>
<i>Figure 5.20 Fuzzy Membership Function for Angle fusion.....</i>	<i>102</i>
<i>Figure 5.21 Recognition result of baseline estimation.....</i>	<i>103</i>
<i>Figure 5.22 Biologically Inspired feature extraction</i>	<i>105</i>

Figure 5.23 Loop confliction and detection	106
Figure 5.24 Chain code extraction.....	108
Figure 5.25.a Layered based low level to complex level feature extraction.....	108
Figure 5.25.b Level-3 and level-4: Fuzzy based small pattern to larger pattern	108
Figure 5.26 Structural feature (inflection point, cusp, end point, branch point).....	109
Figure 5.27 Level-4 Fusion of structural features and directional features.....	111
Figure 5.28 Fusion of structural and direction features.....	113
Figure 6.1 Framework for the study of living organisms.....	119
Figure 6.2 Human visual system.....	121
Figure 6.3 Human behavior diagram for character recognition.....	122
Figure 6.4 Bio-Inspired Machine based character recognition.....	124
Figure 6.5 Machine based character recognition behavior diagram	126
Figure 6.6 Multilevel bio-inspired rule model.....	127
Figure 6.7 Diacritical marks issue with three dots in Naskh style.....	127
Figure 6.8 Basic GC for Arabic script based languages written in Nasta'liq.....	129
Figure 6.9 Simple example using multilayer fuzzy based approach	133
Figure 6.10 Level-4 complex features and candidate independent character.....	135
Figure 6.11 Combination of diacritical marks with respect to languages.....	136
Figure 6.12 Bio-Inspired diacritical mapping	137
Figure 6.13 Combination of diacritical marks with respect to languages.....	137
Figure 6.14.a Ligature formation by combination of diacritical marks and GC.....	138
Figure 6.14.b Word level recognition HMM for same word in Urdu and Arabic ...	138
Figure 7.1 Recognition system architecture.....	144

Figure 7.2 Simple right to left HMM.....	147
Figure 7.3 Fuzzy modeling on HMM wrt language modeling	152
Figure 7.4.a Fuzzy member function for classification	154
Figure 7.4.b Proper classification	154
Figure 7.4.c Confused classification.....	154
Figure 7.5 Layered Fuzzy with HMM.....	156
Figure 7.6 Stroke Mapping	158
Figure 8.a.1 Combination of diacritical marks with respect to languages.....	164
Figure 8.a.2 'Bey' character in all languages	164
Figure 8.2 Division in to ligature from word.....	166
Figure 8.3 Combination of diacritical marks with respect to languages.....	166
Figure 8.4 Combination of diacritical marks with respect to languages.....	168
Figure 8.5.a Ligature formation by combination of diacritical marks and GC.....	168
Figure 8.5.b Word level recognition HMM for same word in Urdu and Arabic	168
Figure 8.6 Bio-inspired character recognition process of word Pakistan.....	170
Figure 8.8 Recognition rate	173
Figure 8.9 Computational Complexity.....	173
Figure 8.7 Probabilistic based modeling for future work.....	175

LIST OF TABLES

<i>Table 5.1 Distance Vs Position Fuzzy Rules..</i>	<i>88</i>
<i>Table 5.2 Result Vs Size Fuzzy Rules.</i>	<i>88</i>
<i>Table 5.3 Diacritical Marks Separation Algorithm</i>	<i>89</i>
<i>Table 6.1 Classes division of strokes.....</i>	<i>121</i>
<i>Table 6.2 Code Book for HMM.....</i>	<i>122</i>
<i>Table 6.3 Recognition result in % on two data set A and data set B</i>	<i>124</i>
<i>Table 6.4 Comparison with previous systems</i>	<i>124</i>
<i>Table 6.5 Vector Quantization</i>	<i>127</i>
<i>Table 6.6 Rule for Secondary strokes handling</i>	<i>132</i>
<i>Table 6.7 Comparison of recognition rate WRT to Methodology</i>	<i>134</i>
<i>Table 6.8 Comparison of recognition rate with previous System</i>	<i>135</i>
<i>Table 7.1 Biological Inspired Character recognition</i>	<i>153</i>
<i>Table 7.1 Recognition Rate</i>	<i>160</i>
<i>Table 8.1 Comparison of Recognition Result of different system</i>	<i>170</i>
<i>Table 8.2 Comparison with existing systems.....</i>	<i>173</i>

LIST OF ABBREVIATIONS

BPNN	Back Propagation Neural Network
CAPTCHA	Completely Automated Turing Test to Tell Computers and Humans Apart
CHMM	Continuous Hidden Markov Model
CNS	Central Nervous System
DHMM	Discrete Hidden Markov Model
EM	Expectation Maximization
GCRT	Ghost Character Recognition Theory
GCT	Ghost Character Theory
IFN/ENIT	Arabic Handwritten Database
IT	Inferotemporal Cortex
HMM	Hidden Markov Model
KNN	K Nearest Neighbor
L-0	Level-0: Stroke Acquisition
L-1	Level-1 Low Level Processing

L-2	Level-2 Geometrical Processing
L-3	Level-3 Simple Feature Extraction
L-4	Level-4 Complex Feature Extraction
L-5	Level-5 Ligature Formation
L-6	Word Formation
LDA	Linear Discriminate Analysis
LGN	Lateral Geniculate Nucleus
LVQ	Linear Vector Quantization
MER	Minimum Enclosing Rectangle
MLP	Multilayer Preceptron
MT	Visual Area
V1	Primary Visual Cortex
PCA	Principle Component Analysis
PDA	Personal Digital Assistants
OCR	Optical Character Recognition
SVM	Support Vector Machine
V1	Striat cortex

CHAPTER 1

INTRODUCTION

Computers have brought revolutionary changes in the way we do things and they are even imitating the human's capabilities. Character recognition is to mimic the human ability to read just like a human using a computer although here the machine has not been able to compete with the human visual system. Handwritten text provides difficulty for even human readers who utilize contextual information as well as the domain knowledge to interpret and recognize such text. The concept of handwritten text is very old, for the purpose of expanding people's memory and facilitating communication. Character recognition is an important offshoot of pattern recognition problems that imitate the human reading capabilities in machine. This task becomes complex and demanding if it involves cursive text and writing styles like Arabic. One of the aims of character recognition is to read much faster rate by associating symbolic identities with images of characters. Some of the potential applications of character recognition are number plate recognition, automatic mail sorting, automatic form reader archiving and retrieving text

provision of interface for entering text in PDAs through pen like devices, etc shown in figure 1.1.

The tremendous advances in the computational intelligence algorithms have provided significant improvement in the development of character recognition systems. OCR has become one of the most successful and challenging applications of artificial intelligence and pattern recognition. It is an interdisciplinary research area, involving researchers from pattern recognition, computer vision, image processing, statistical computation and machine learning. Many commercial character recognition applications such as MyScript, Calligrapher, ICHITARO etc. exist for different languages and different applications. However in the case of languages like Persian, Urdu Arabic, etc which share the same basic Arabic character set and have been coined as Arabic character script based languages very little work has been reported [Husain et.al, 2007]. Moreover, we could not find any no public database for Arabic script based languages [Biadisy et.al, 2006], [Halavati et.al, 2005] especially for Nasta'liq writing style. It is after long time that Middle East Asian script has gotten attention in research, which may be due to the complexities of Arabic script [Sternby et.al, 2009].

A computer can perform most problems more precisely and efficiently as compared to a human. However, there are certain types of complex problems like pattern recognition and visual perception where a few years old child can perform much better than the state-of-the-art computer with latest algorithms that are available today. With the growth of child, he starts to acquire reading and writing skills and he is able to read the texts, whether it is written in different fonts and styles. Moreover, even a human may have problem in reading complex text i.e. poorly written and requires context knowledge

and experience for correct recognition. The complex mechanism behind human visual perception and context knowledge makes the human more powerful than the current machines. Due to these factors, lot of research is required in this field.

The computer can perform many complex problems more precisely, efficiently and faster than human however there are still many problems such as pattern recognition, where few months old child can perform much better than the latest algorithms available today and one such problem is character recognition. The real world is dynamic, enormous and very complex. To deal with it, biological systems have been evolving and self-adjusting over millions of years and adapting to the changes in the environment. Imitation of human capabilities in machine is a field of intense research since the early days of computer and among them one of the interesting issues is human vision modeling. In this regard biology has been an important source of inspiration in image processing and pattern recognition applications. The high context knowledge, adapting and learning capability of human is robust intelligence in complex problems.

To imitate the human brain capability into machine is crucial and very difficult in real world application. Human brain can be modeled by deeply investigating the human visual system. To model the biological vision of human the research on human visual system has been in progress since 80, but still it is facing many challenges. The image processing inspired by biological perception is an efficient way to handle the complex problems. The high context knowledge, learning ability and efficient visionary devices are very robust and cannot be achieved perfectly by machine. Thus modeling the human brain exactly into our daily life applications is almost impossible in current age but the integration of biological visionary system of humans helps to recognize the

complex pattern recognition problems efficiently than conventional methods. More detail of biology inspired system is briefly described in chapter 3.

The happy relation of fuzzy logics with biologically inspired processing for character recognition is an ideal way to deal with complex handwritten patterns. Fuzzy logic has proved itself a powerful classifier to recognize irregular and complex patterns. Naturally translation, rotation and position are independent of human visual perception which cannot be attained by computer system but it has substantial effect on reading speed and accuracy even with human at some level. Due to the success of fuzzy logics for complex patterns, high level fuzzy IF-THEN thinking and reasoning rules are modeled with human visionary concept to resolve the issue of Urdu handwritten text. Further detail of fuzzy with biology inspired recognition is briefly described in chapter 3.

Since the inception of computers, a lots of research efforts has gone in to the field of computer human interface. Input devices such as mouse, keyboard, etc. have several limitations. Two quick ways of interaction with computers is speech and handwriting. Speech recognition has several limitations. Handwriting is the easiest and accurate way of interaction between the machine and humans. There are also several limitations of handwriting due to the variations involved in it.

1.1 Motivation

On-line Latin or Asian language has been a research issue for thirty years. However, only few researchers have focused on Arabic language [Kherallah et.al, 2008]. Most of the Arabic script recognition systems do not allow noisy input. Therefore, multi-

cultural and Multilanguage handwriting styles is the current and hot research issue for Arabic script based languages. The Urdu/Arabic characters are used by at least ¼ of the world population and in many countries and by many nations in the form of different languages such as Arabic, Persian, Urdu, Sindhi, Punjabi, Pashto, etc. Thus the automation of Urdu script based languages has wide spread benefits whereas different writing styles are used by these languages; the most common are Nasta'liq and Naskh. The research for Multilanguage character recognition is still open and challenging problem due to different writing styles and language property. Moreover segmentation based approach for handwritten Urdu script based languages are still in demand due to limitation of dataset for segmentation based approach.

سنگاپور، گرلز سکول میں آئی پیڈ نے کتابوں کی جگہ لے لی

طالبات کتابوں کا بوجھ اٹھانے سے بچ گئیں، اساتذہ، یہ طریقہ آسان ہے، شاگرد

سنگاپور (آئی این پی) سنگاپور کے گرلز سکول میں آئی پیڈ نے کتابوں کی جگہ لے لی۔ گرلز سکول میں اساتذہ اور طالبات کے ہاتھ میں کتب نہیں بلکہ آئی پیڈ ہیں، معلم اور شاگرد محض انگلیوں کی حرکت سے مطلوبہ مواد حاصل کر لیتی ہیں یہ طریقہ کار صرف آسان ہی نہیں جدید بھی ہے کیونکہ طالبات کتابوں کا بوجھ اٹھانے سے بچ گئی ہیں۔ اساتذہ کا کہنا ہے طالبات اسے زیادہ آسانی سے استعمال کر رہی ہیں۔ آزمائشی طور پر سکول کی 140 طالبات اور 10 اساتذہ کو آئی پیڈ فراہم کئے گئے جس پر ایک لاکھ 35 ہزار ڈالر خرچ ہوئے۔ سکول انتظامیہ کے مطابق آئی پیڈ کے استعمال کا مقصد درس و تدریس کو جدید ترین خطوط پر استوار کرنا ہے۔

Figure 1.1: World vision in order to facilitate students using iPod.

1.2 Classification of Character Recognition.

The origination of character recognition starts from 1870's with the invention of scanner whereas the first successful attempt was made by Tyurin, Russian scientist in 1900. The modern character recognition system started from 1940s [Mantas 1986]. Now a day it is a most intensive field of research. There is a worldwide interest in the development of handwritten character applications and it is the most important part of research of artificial intelligence and pattern recognition. The aim of the character recognition system is to behave like human in reading the spatial form of handwritten marks. With respected to the user data, character recognition can be writer independent or writer dependent. Writer dependent character recognition is much more easy and more accurate than writer independent [Subrahmonia, 2000] because in writer dependent the system is trained on few user's and only these users can use the system. The writer independent must be expert in commonly used writing style and intelligent to large variations. Character recognition system differs with respect to mode of input i.e. input acquisition (offline or online), writing mode (printed or handwritten) and font (single or omni) [Al-Badr 1995], [Govindan 1990]. This categorization is illustrated in figure 1.2. Depending upon the nature of application character recognition system is classified into two groups.

1.2.1. Offline Character Recognition

Offline character recognition system deals with recognizing the input text after it is written. The input is obtained by digitizing the paper using a camera or scanner. Therefore input for offline character recognition is the two dimensional spatial

information of the written text. The offline character recognition is further classified into two types with respect to the mode of input offline printed character recognition and offline handwritten character recognition. Off-line handwriting character recognition refers to the process of recognizing the handwritten words that have been scanned/captured from paper.

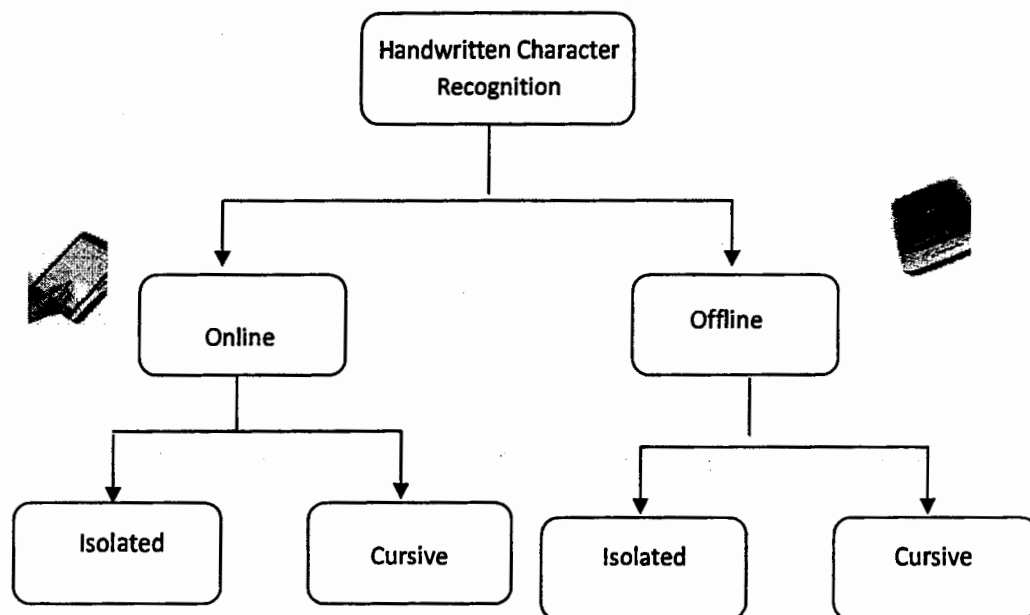


Figure 1.2: Character Recognition methodologies.

1.2.2. Online Character Recognition

Another attractive area of character recognition is a digital pen and handwriting document management. The digital pen allows people to read, write, and memorize, and integrate these acts into information systems [Fujisawa, 2008]. The digital pen can capture handwritten shapes and store it in a very natural way. The system can recognize the text during the time user is writing with a digital pen or pc tablet shown in Figure 1.3. The input through pc tablet or digital pen consists of duration, stroke elements and stroke

order. It has great potential to improve the communication between human and machine. The use of digital pen enables humans to work in digital world just as in real world. Handwritten characters provide the most promising method for interaction with small portable machines. Online character recognition is becoming popular input mode especially in the growing pen applications and improvement in handheld devices such as Tablet PCs, Palm, and Pocket PC based PDAs, etc [Tanaka 2003]. This may be due to the natural way of human interaction with these machines. Online handwritten character recognition is suitable for mobile application whereas it is difficult to use the keyboard i.e. small PDA has a very small keyboard and it is very difficult to use the keyboard. Whenever the term handwriting recognition is used without prefixing it with offline and online terms, people normally think of handwritten input using recognition devices i.e. personal digital assistants (PDAs), pc tablet, etc.

Now a day's online character recognition is gaining more intention due to the latest improvements in handheld device, latest mobile phones, etc. Due to the complexities involved in handwritten text i.e. variability, distortions, etc. the accurate handwritten character recognition is very difficult.

1.3 Constrained Vs Unconstrained

Handwriting character can be constrained or unconstrained [Tappert, 1984]. In nature, constrained handwriting text is the spaced discrete and boxed discrete. It is written inside the special boxes as shown in figure 1.4(a). The character written separately with spaces and don't touch the other characters are called spaced discrete handwriting. If the characters are written separately and they touch each other, the

writing is called run-on discrete handwriting whereas the unconstrained handwriting text is cursive and complex in nature. In the cursive text, the characters are connected and the strokes are more than one in individual character shown in Figure 1.4(b). Normally, the people write in mixed cursive styles that include spaced, run-on discrete and cursive but it is called the cursive handwritten text. It is a very difficult task to recognize the cursive handwriting text due to the complexities involved in it i.e. order of strokes, noise, and speed of writing and large variations in shapes.

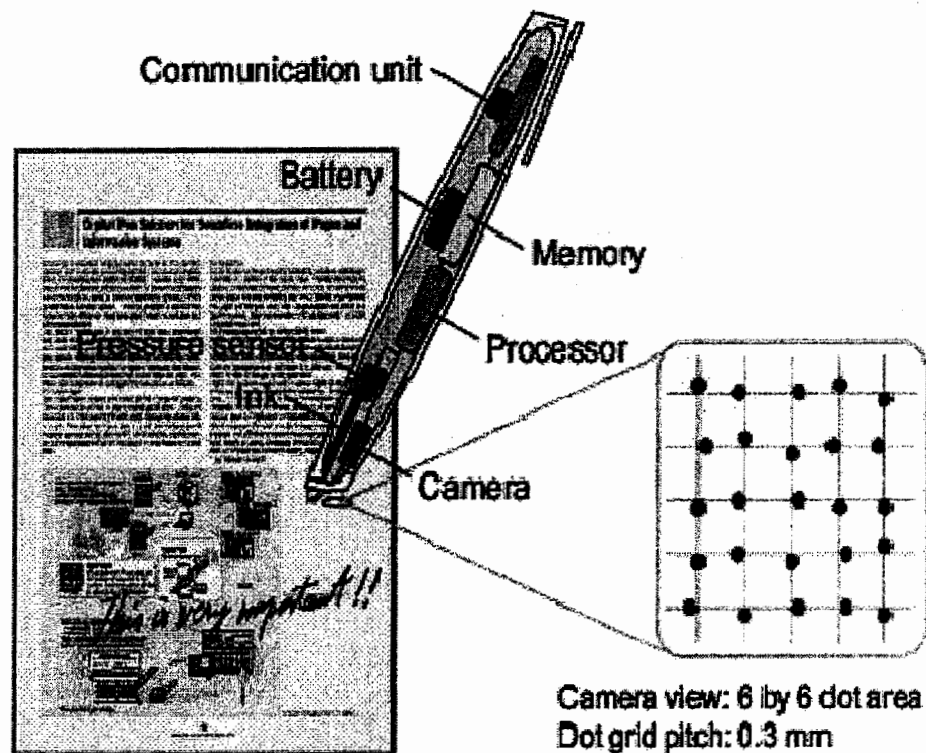


Figure 1.3. Digital pen using Anoto functionality [Fujisawa, 2008]

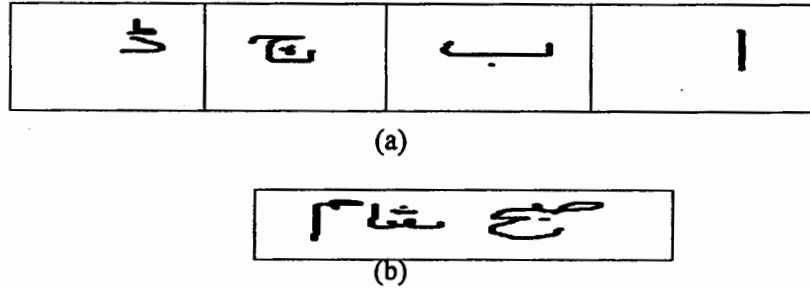


Figure 1.4 (a) discrete writing (b) cursive writing

1.4 Urdu/Arabic Character Recognition Process

Generally, the character recognition process is divided into four phases which are preprocessing, feature extraction, training/classification and post processing, shown in Figure 1.5. Preprocessing phase is responsible to clean the raw input from noisy data and normalize it for further processing. Feature extraction allows extraction of some unique patterns from large data that can better discriminate the shapes. A successful character recognition methodology depends highly on the particular choice of features used by the pattern classifier. Feature selection in pattern recognition problems involves the derivation of salient features from the raw data input in order to reduce the amount of data used by the classifier for classification and simultaneously provide the enhanced discriminatory power.

The classification phase is responsible for recognition of learning based on the discriminate features whereas, post processing phase is further applied on to the recognized text to further investigate the recognized data based on context knowledge/word knowledge i.e. diacritical mapping, dictionary based processing, word formation, etc.

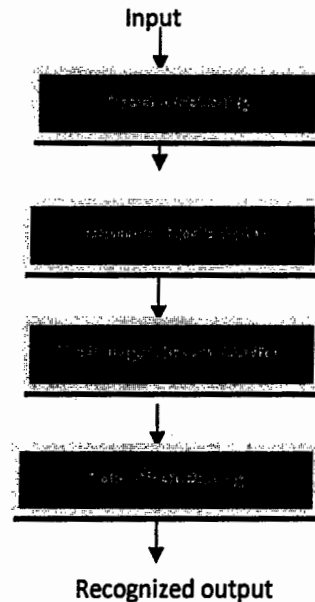


Figure 1.5: Character recognition process

1.5 Challenges in Urdu Script Character Recognition

Arabic civilization has wide culture because many civilizations have been inspired by Arabs culturally and Arabic language is closely associated with Islam and with a highly esteemed body of literature. The Arabic/Urdu scripts consist of many challenges that make it more complex as compared to other scripts. Moreover the Urdu script is written in Nasta'liq that makes it more difficult in the family of Arabic script based languages. Unlike other scripts Arabic script based languages are written from right to left. It is written in cursive form in both printed and handwritten form. The shape of the character depends upon the associated connected characters. Nasta'liq script is more complex as compared to Naskh style. For example the character in Naskh are up to 4 depending upon the position whereas the character shapes may be up to 32 depending upon the associated characters on both sides and secondly the shape of Nasta'liq is more complex as compared to Naskh shown in figure 1.6.

Moreover, Arabic script based languages are rich in diacritical marks especially Urdu, hence it is difficult to differentiate between similar shapes. Historically, diacritical was not the part of languages and hence, it was added later on for nonnative. The extensive amount of dots associated with each word makes it very complex. The presences of diacritical marks make base line detection, line detection and segmentation more difficult.

ب، پ، ت، د، ر، ی، ص، ط، ن، یو، یق، یف، یج، یز، یل، یاء، نی

(a)



(b)

Figure 1.6(a): different shapes of Nasta'liq script (b) diacritical marks complexity

Characters may overlap vertically even without touching each other as shown in Figure 1.5. Moreover the size of each ligature varies due to cursiveness and size of character. Thus, the Arabic characters have not got symmetry in height and width. Words may consist of one or more ligatures whereas the ligature may be composed of one or more alphabets i.e. only Urdu consists of more than 22000 ligatures [Attash Durani, 2010]. It is very difficult to segment the ligatures into characters especially in the case of Nasta'liq script. Some of the Urdu script letters have same form however; they are differentiated from each other by the addition of diacritical marks at different positions. Moreover, some character use special marks which leads to modify the characters accent [Kherallah et.al 2009]. These diacritical marks are positioned at certain distance either below or above from the character. In most of the cases, the Urdu script based languages writing for native reader does not use vowels because the sense of the word can be

determined by the context of the sentence. However they are used for nonnative user. The richness of Urdu script based languages in diacritical marks makes it more difficult and complex. Moreover the Urdu words composed of sub words and sub words may compose of one or more letters and these letters are combined in cursive form. This property of similarity of character, ligature and words makes the Urdu script based language more difficult due to the absence of context knowledge. The very similar contours are very difficult to segment and recognize [Hamad and Zitar, 2010]. The characters such as the shown in Figure 1.6.b complicate segmentation and recognition when it occurs in the middle of a word.

Variations in handwritten text depend upon the alignments and the different forms of handwritten strokes and these variations are geometrical. Most common geometrical properties are position, size, aspect ratio of strokes, slant and number of strokes in a character. Characters may look similar although the number of characters, number of strokes, drawing order and direction of the strokes may vary considerably. Thus recognition and segmentation of Urdu handwritten script based languages is a difficult task because the Urdu script handwritten characters are naturally cursive and unconstrained and it is believed that a good segmentation technique is one of the important reasons for high accuracy character recognition.

1.6 Objectives

Urdu script based languages has received limited attention as compared to the other scripts and therefore no single commercial application exists for the purpose. Moreover, Urdu script based languages are similar to each other due to the same script. The idea of

Multilanguage Urdu script based languages has not been considered. The one of the basic objectives is to provide a framework for Multilanguage Urdu script based upon the language characters. While the main objective of the thesis is to explore the segmentation based techniques to avoid the heavy training required for ligature approach and limitation of dataset i.e. only Urdu script contain more than 22000 ligatures and it is very difficult and complex to train large datasets. The third objective is to explore the character recognition according to the biological perceptive of human vision for Urdu script based languages which are rich in diacritical marks and complex shapes. An efficient solution to recognize the complex Urdu script is to observe them based on biologically vision. Natural images contain ambiguous and incomplete data due to various effects. Due to the large variations in modeling of human visual perception, a fuzzy logic is a significant tool to handle such complex problems. Most of the cases, the Urdu script does not use the special diacritical marks for accent because the sense of the word is determined by the context of the sentence. Thus all the vowels in this work are not considered.

1.7 Knowledge Based Urdu Character Recognition System

The intelligent knowledge-based methods used to imitate the human visual system include expert systems, genetic algorithms, neural networks, evolutionary computing and fuzzy logics. Fuzzy logic provides a basis for representing uncertain and imprecise knowledge and uncertainties can be modeled easily as inspired from the human visual system. This thesis presents biology inspired Urdu script based languages character recognition using fuzzy logic. The ghost character recognition theory provides the framework for Multilanguage character recognition system similar to human visual system. Although the human visual system is translation, rotation invariant, etc, yet it has

some effect on accuracy and speed. Fuzzy based several preprocessing steps e.g. skewness, baseline estimation, slant correction, smoothing, etc. are applied to normalize the handwritten input strokes. These series of steps are applied only on the ghost character to avoid the error introduced due to diacritical marks. The preprocessing results in clean and smoothed input for feature extraction like the input is normalized by eyes so that the data is translation and rotation invariant. The structural and directional features are extracted in a layered manner like the human and fusion is performed to create a new feature matrix. The lower layer extract similar features and these similar features are further refined to form new complex features. The languages rules are added with fuzzy rules to extract unique and discriminant features. Modified HMM is used for the segmentation free approach whereas fuzzy logic based biological inspired method is used and partially implicit segmentation based recognition is performed. Fuzzy provided an excellent way to model the biological concept for Urdu handwritten character recognition. Finally, the diacritical marks are mapped on to the primary strokes with the help of language rules.

1.8 Contributions

The original contributions of this works are

A new method “Biological inspired character recognition for Urdu script based languages” is presented. It is built as human visual system which finds the larger more significant patterns from the smaller patterns.

- I. Multilanguage technique “Ghost Character Recognition Theory” for Arabic script based language is presented. Due to the similarity of all

Arabic script based languages, ghost character recognition theory helps to build the Multilanguage character recognition system. The diacritical marks are separated and ghost shapes are treated separately and finally mapping of diacritical is performed based on the language.

- II. In order to make handwritten communication with machine more natural, several preprocessing steps for handwritten online Urdu script character recognition have been presented to overcome the issue in raw input.
- III. Fusion is performed on structural and directional features to form new feature sets. The combination of structural and directional features is based on time and position information.
- IV. Segmentation free (Modified HMM, Hybrid Approach: Fuzzy logics and HMM) and simultaneous bio-inspired segmentation based approach are presented. HMM classes are reduced using fuzzy logics based on the word structure whereas bottom up approach is used to find the more discriminatory patterns from the small sub patterns using recursive approach.
- V. A fuzzy logics based method for diacritical marks handling is presented. As these languages are rich in diacritical marks, and they have lots of issues in handling of diacritical marks. Fuzzy clustering is performed to read diacritical marks.
- VI. Finally mapping of diacritical marks is performed based on the languages structure.

1.9 Thesis Organization

The focus of the thesis is the applicability of bio-inspired preprocessing, feature extraction and recognition of text in layered manner and analysis of segmentation and segmentation free approach for Urdu script based languages character recognition especially languages written in Nasta'liq script. The rest of the thesis is divided into eight chapters. The organization of the thesis is as follow:

Chapter 2. Handwritten Character Recognition: A Literature Survey

Chapter 2 briefly presents the state of the art on Arabic script based language character recognition. The main phases on Urdu character recognition i.e. preprocessing, feature extraction, classification and post processing is reviewed. Finally, the gaps between the previous works are described.

Chapter 3 Fuzzy Logic and Human Biological Visual Perception

This chapter briefly describes the fuzzy logics and analysis of human visual perception and computational modeling on applications. The first part of this chapter presents fuzzy logic and the second part exposes the terminology of the material studied inspired by biological systems. Object recognition in real world still remains challenging and important in daily life applications such as character recognition, robot navigation, image understanding, security, etc. as compared to application in constrained environment. An efficient solution to recognize the complex pattern is to observe them based on biological vision. Due to the large variations the modeling of human visual perception using fuzzy logics is a significant tool to handle such complex problems. To deal with the complex problem and involve human intelligence, fuzzy logic inspired by biological vision

concept is presented. Finally it presents biological inspired character recognition with the concept of fuzzy logic. The main contribution of this chapter is to discuss the biological inspired fuzzy model.

Chapter 4 Ghost Character Recognition Theory and Arabic Script based Languages Character Recognition

This chapter presents the ghost character recognition theory based on ghost character theory. The proposed ghost character recognition theory is the important step towards the biological inspired and Multilanguage character recognition system. The main benefit of proposed approach is that it works well for all Arabic script based languages by concentrating ghost shapes and diacritical marks separately and developing dictionary for every language that helps in ligature and word formation. Handling all Arabic script based languages has several issues as compared to system for specific languages and specific writing style i.e. Nasta'liq or Naskh.

Chapter 5:

Preprocessing and feature extraction is the most important and complex part of the thesis. The handwritten Urdu script contains lot of variations and noisy pattern. Based on the ghost character theory proposed in chapter 4, the diacritical marks are separated from the base stroke and treat the ghost character separately. The raw input data is preprocessing by applying several preprocessing steps i.e. baseline estimation and correction, smoothing, de-hooking, normalization, etc. Secondly, the unique features are extracted, based on the geometry of handwritten

strokes. The aim of feature extraction from the input strokes reduces the input pattern to avoid complexities while maintaining the high accuracy and the extraction of these distinct patterns that uniquely define the strokes and most important for classification. This phase includes the feature extraction from the primary strokes that provides sufficient information to the inference engine for classification. The biological inspired features are extracted in bottom up fashion. Both directional and structural features are extracted using fuzzy logic.

Chapter 6: Holistic Approach for Urdu Recognition using Modified HMM

Segmentation based approach for recognition of handwritten Urdu script incorporates a considerable overhead and has less accuracy as compared to other script. This chapter presents two approaches for Urdu character recognition. The first phase introduces a ligature based approach using Hidden Markov Model that provides solution for recognition of Urdu script. HMM database is divided into 54 subclasses based on the starting and ending shapes of the ligature. The sub division in classes reduces the time complexity and increases the efficiency. The second part presents a segmentation free approach for recognition of Online Urdu handwritten script using hybrid classifier, HMM and Fuzzy logic. The classification is performed using fuzzy instead of exactly based on feature. Trained data set consists of HMM's for each stroke is further classified into 62 sub pattern based on the primary stroke using fuzzy rule. Fuzzy linguistic variables based on language structure are used to model features and provide suitable result for large variation in handwritten strokes.

Chapter 7: Biology Inspired Character Recognition

Due to the limitation involved in segmentation approach, this chapter presents segmentation based approach for online Urdu script based language character recognition. The concept of ghost character recognition discussed in chapter 4 helps to build biological inspired and Multilanguage character recognition system. This chapter integrates the biological concept with fuzzy logics to achieve robust system in terms of speed, accuracy and efficiency. The visual perception system of human is very powerful in pattern recognition problems. Bio-inspired multi-layered fuzzy rules based expert system is presented for handwritten Multilanguage character recognition i.e. Arabic script based languages inspired by visual ventral stream to improve the recognition accuracy.

Chapter 8: Comparative Analysis and Conclusion

This chapter analyzes the presented approach i.e. ghost character theory, effect of fuzzy preprocessing and feature extraction and finally compare the segmentation free approach with biological inspired segmentation base approach. Finally future work is presented.

CHAPTER 2

HANDWRITTEN CHARACTER RECOGNITION: LITERATURE SURVEY

Character recognition is a fundamental but most challenging in the field of pattern recognition with large number of useful applications. It has been an intensive field of research since the early days of computer science due to it being a natural way of interactions between computers and humans. There is a worldwide interest in the development of handwritten character applications and the tremendous advances in the computational intelligence algorithms have provided new tools for the development of intelligent character recognition. With respect to the mode of input, Handwritten character recognition is classified into two classes namely; online and offline. In offline handwritten character recognition, the input is available in the form of an image obtained through a scanner or a digital camera whereas in online character recognition coordinate information of strokes is available with timing information, this additional timing

information makes online recognition easier as compared to offline. Furthermore online character recognition can be processed from both online and offline perspective while versa is not possible.

From the classifier perspective, character recognition systems are classified into two main category i.e. segmentation free (global) and segmentation based (analytic). The segmentation free also known as the holistic approach to recognize the character without segmenting it into sub units or characters. Each word is represented as a set of global features e.g. ascender, loops, cusp, etc. Whereas segmentation based approach; each words/ligature is segmented into subunits either uniform or non-uniform and subunits are considered independently [Lorigo and Govindaraju, 2006]. The segmentation systems have low recognition rate as compared to segmentation free approach due to error involved in segmentation. On the other hand segmentation based approach are more powerful for large datasets as compared to the segmentation free systems.

Automatic recognition of cursive handwritten script remains a challenging problem even with the promising improvement in the classifier. Online Urdu script based languages character recognition is very challenging task due to its complex structure of graphical marks. Vast variation and inconsistencies makes handwritten Arabic script based languages character recognition much more complex than any other language. Segmentation free methods are very helpful in the recognition of complex patterns, but for Urdu script, the numbers of ligatures are more than 22000 [Attash Durani, 2010]. Thus it's very difficult to build a system based on segmentation free approach for such large amount of ligature in term of both complexity and recognition rate. The segmentation based approach for recognition of handwritten Urdu script incorporates a

considerable overhead and has extremely low accuracy as compared to other script. Moreover presence of complimentary characters in Urdu language makes it complicated as these have to be segmented into secondary strokes. A lot of research works have been done for offline printed Arabic Persian and Urdu documents with reasonable levels of accuracy. However, within the context of online handwritten character recognition, studies dealing with Arabic script based languages characters are scarce [Hussain et.al, 2007] especially for Urdu script. Furthermore the work done for Arabic and Persian cannot be directly applied to Urdu due to more complexities involved in Nasta'liq style as compared to Naskh. Moreover, there is no standard dataset available for Urdu script especially for online character recognition. For handwritten offline only few available databases are IFN/ENIT [Pechwitz et.al, 2002], LMCA[Kherallah et.al, 2008], AHDB [Adeed et.al 2004] and Urdu database and Farsi database by CENPARMI respectively [Sagheer et.al., 2009],[Haghighi et.al, 2009].

Normally handwritten recognition is divided in to four phases which are preprocessing, feature extraction, classification and post processing. The input is obtained from PC tablet and it consists of strokes elements x, y in order of occurrence, time, force, velocity. The additional timing information and order of strokes elements makes online character recognition little easy over offline handwritten character recognition. Force and velocity information are helpful in personality observation, person identification, etc. but not helpful for character recognition due to different force and writing speed. Strokes are acquired from the movement of pen during the pen down and pen up event. These raw inputs strokes contain lot of inconsistencies and large data which cannot be directly used

by classifier, thus preprocessing and pattern extraction are performed on raw input strokes.

2.1 Preprocessing

Preprocessing is the basic phase of character recognition and its crucial for good recognition rate. The main objective of preprocessing steps is to normalize strokes and remove variations that would otherwise complicate recognition and reduce the recognition rate. The input of preprocessing phase is raw data x,y strokes elements and output is normalized strokes. Several basic preprocessing steps are skew detection, noise removal, filtering, slant correction, normalization, etc.

J. Sternby et.al performed noise reduction at the time of recognition and some improper pattern are considered as noise whereas segmentation of strokes into sub patterns that contain shapes of individual character is performed as a preprocessing step. Some robust rules are added with traditional segment point to segment the strokes into independent subunits whereas the missing segmentation and complex curve pattern are not considered [Sternby et.al, 2009]. J. Schenk et.al performed histogram based skew detection and slant normalization and for normalization the handwritten strokes re-sampling is performed to achieve the strokes elements at equidistant [Schenk et.al, 2008]. Ahmad and Maen performed slope estimation and normalization for online Arabic digits [Ahmad and Maen, 2008]. The sign of slope values zero and infinity are extracted to use for feature extraction and breaks points are estimated from slope value. Primitives are extracted from normalized slope and break points and finally primitives are sorted to reduce the input pattern size. Hussam et.al presented several preprocessing steps in order

to normalized the handwritten strokes i.e. noise removal, baseline estimation, diacritical marks extraction, middle region detection, histogram estimation, modified histogram estimation, and stroke width estimation, etc. and the diacritical marks are extracted at first stage [Hussam et.al, 2010]. Lui et.al performed nonlinear moment based normalization on Bangla and Farsi numerals character on both binary and gray scale image [Lui et.al, 2009]. Sternby et.al performed segmentation by segmenting the character at vertical extreme point with respect to the direction of writing of character and the output of segmentation is individual character by using heuristic rules [Sternby et.al, 2009].

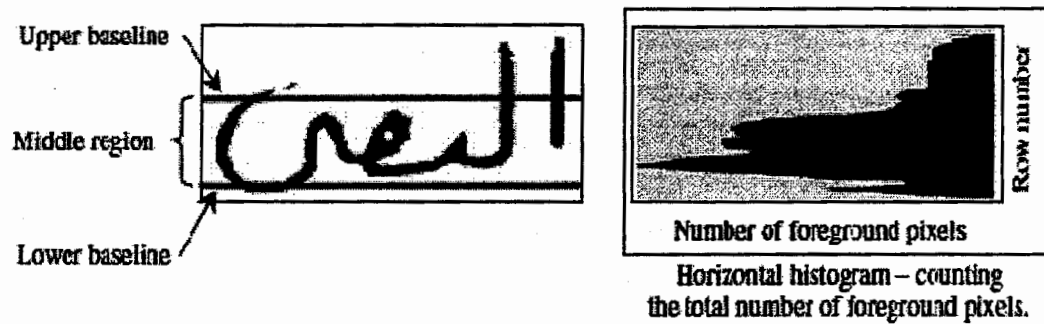


Figure 2.1: Baseline estimation using histogram [Hamad and Zatar, 2009]

Husain et.al performed smoothing; de-hooking on chain code and Malik and Khan performed repetition removal and filtering step to remove the repetition and staircases for Urdu input. El-Anwar et.al presented a method to relate the broken strokes due to pen lift during writing due to fast pen movement and segment the dot from the strokes to reduce the feature set for Arabic language [Randa et.al, 2007]. Al-Ghoneim K. et.al performed translation, scaling, connected line generation, and smoothing for offline input and translated the image using the center of gravity and used bresenham's line algorithms to fill the missing [Al-Ghoneim et.al, 2001]. Biadisy et.al used low pass filter for noise

reduction and smoothing and Douglas and Peucker's algorithm [Douglas D, 1973] is used to eliminate the redundant points [Biadsky et.al, 2006]. Benouareth et.al performed binerization, smoothing, baseline detection and thinning operation for offline Arabic handwritten recognition [Benouareth et.al, 2008]. Razzak et.al proposed several preprocessing steps for online Urdu character recognition from both online and offline perspective [Benouareth et.al, 2008]. Adeed presented normalization to reduce the character to uniform height of strokes based on slope, stroke width, and stroke height [Adeed et.al, 2002]. Ibrahim et.al performed smoothing on chain code. Baseline is the virtual line on which semi cursive or cursive text are aligned/ joined. Generally baseline is kept in mind during both writing and reading. Baseline detection is not only used for automatic character recognition but it is also necessary for human. Without baseline detection it is very difficult to read the text even for human due to improper visibility and error rate increase up to 10% while the context sensitive interpretation is involved in human reading. Whereas in automatic classification no context based interpretation is involved and decision of baseline may be on one word which is very difficult especially estimation of diacritical marks without proper baseline being very difficult. Thus baseline detection is the necessary part of better classification especially for Arabic script based languages. Several baseline detection methods based on horizontal projection have been proposed in literatures but they are for large text lines.

The horizontal projection based approach counts the elements on horizontal line and assumes that maximum number of elements on horizontal line is the baseline. Although it is robust and very easy to implement but it needs long straight line of text but in the case of handwritten text especially for online handwritten text the length of the line

may be very short. Thus the histogram projection mostly fails in estimating the correct baseline for isolated handwritten text and ligatures having greater number of ascender and descender.

Boubaker et.al presented a novel method for both online and offline Arabic handwritten text based on the graphemes segmentation. The segmentation of Arabic character works on the previous detection of baseline using the topologically points i.e. the backs of the valleys adjoining the baseline and the angular points. [Boubaker et.al, 2009]. Maddouri and Abed compared the six methods for Arabic character recognition on IFN/ENIT database. Projection based method fails to estimate the baseline for short word length and words having more diacritical marks, ascender, descender [Maddouri and Abed, 2008]. Min-Max is presented based on the projection based method. Min-Max contour method used critical points from the word contour and two baselines upper and lower are extracted from mean of maxima and minima respectively. The combination of Min-Max and some structural primitives i.e. loops, diacritical marks is used for baseline estimation. These additional primitives are used to differentiate contours from the others.

Faisal et.al modified the RAST algorithms by introducing two descender lower line d_1 and descender upper line d_2 for Urdu images [Faisal et.al 2005]. Farooq et.al and Benoureh et.al presented linear regression on local minima of the word for baseline detection [Farooq et.al 2005], [Benoureh et.al 2009]. Mohammad et.al presented a vertical projection algorithm obtained by summing the value along x-axis and detected two baselines [Mohammad et.al 2009]. The lower baseline is identified by maximum projection profile. The upper baseline is estimated by scanning the image from top to bottom. Alkhateeb presented knowledge based baseline estimation by using the location

information for baseline estimation [Alkhateeb,2009]. The algorithm is improved by estimating the baseline at bottom half of image because of baseline existence at bottom of word. The vertically projection is inefficient for small length text. Benouareth et.al used projection after transforming image into Hough parameter for baseline estimation [Benouareth et.al ,2009]. M.Pechwitz et.al used linear piecewise curves using projections for baseline estimation [M.Pechwitz and V.Maergner, 2003]. Izadi et.al performed re-sampling; smoothing and de-hooking for online writer independent character recognition and the character are normalized at same height and the original aspect ratio of the character remains same [Izadi et.al, 2008]. For training SVM, point re-sampling is performed to make all the instants of same length in order to get feature matrix of same length. Amor and Amara et.al presented several preprocessing step i.e. noise reduction, edge detection, etc. [Amor and Amara, 2006]. Aly et.al presented method for baseline discrimination dealing super script and subscript for mathematical equation and they normalized the text in bounding box by estimating the partial height of character [Aly et.al , 2007]. Deepu et.al performed smoothing, translation and re-sampling for online Tamil handwritten character recognition. The original coordinates are replaced with new coordinates using piece wise linear transformation [Deepu et.al, 2004]. Baghshah et.al performed stroke segmentation, noise reduction and smoothing for online Persian handwritten character recognition. Smoothing is performed using the average point of neighbor pixel and apply filtering step to reduce the close point [Baghshah et.al, 2006].

Arabic script based languages are rich in diacritical marks. The primary shape is called ghost character where as secondary strokes associated with each character is called diacritical marks. The diacritical marks are small secondary strokes such as dot, accents

to distinguish basic shape from other [Sternby et.al, 2009]. The basic rule of Arabic script based languages is that primary stroke contains zero or more diacritical marks shown in figure 1.5.a. Moreover one character may contain up to four diacritical marks and it may be increased to six by adding zeer, zaber, etc. These diacritical marks may appear at inside, top, below, right side of the ghost character. Basically each word in Arabic script is broken into ligature which consist one or more alphabets.

Diacritical marks separation in Nasta'liq is more complex as compared to Naskh style due to the complex nature of Nasta'liq over Naskh i.e. پ has 32 shapes depending upon the position of character and adjoining character on both side as shown in figure 1.6.b. Thus the recognition for handwritten Nasta'liq script is more complex. More than half of the Urdu script based languages character set consists of diacritical marks whereas these characters have as unique main stroke. Normally maximum number of dots for one character is three but for some languages i.e. Sindhi, Pashto it may be four. The huge number of dots sometimes causes word to be read in various forms with completely different meaning due to incorrect association.

Researchers have proposed different methods for Arabic script character recognition but the attention paid to diacritical marks is very low where as it is more complex in the case of handwritten character recognition especially for Nastaliq script. Sternby et.al handled diacritical marks by adding to the segmentation graph as extra edge [Sternby et.al, 2009]. Biadsy et.al estimated the diacritical position by projecting the diacritical marks onto the base stroke and detection is performed by using size, shape of bounding box, location and time of the stroke [Biadsy et.al, 2006]. The last point of diacritical marks is added with the base stroke where the projection falls and new

sequence includes all points of base stroke mixed with secondary strokes. Sari et.al extracted morphological features which are stored in associates list to consider diacritical marks [Sari et.al, 2002]. Saabni and Sana handled the diacritical marks based on sequential order, size and location [Saabni and Sana, 2009]. Benouareth et.al used location information for offline Arabic handwritten character recognition [Benouareth et.al al, 2008]. Khorsheed modeled the order of diacritical marks with other base character features [Khorsheed, 2003]. Husain et.al handled the diacritical marks using position and order information. Husam et.al removed the punctuation marks to find the best vertical histogram. The extraction of punctuation marks is based on the density of the stroke. The stroke is removed if its density is less than constant the c and on the other hand they also face problem in the case of small base characters such as ج [Husam et.al, 2010]. The separation of dot is only for vertical histogram estimation whereas the consideration of diacritical marks on to associated character is based on this vertical histogram. This approach is somehow suitable in the case of Naskh style for Arabic script while it has several issues for Arabic script used Indian region i.e. Urdu, Punjabi, Sindhi.

Kherallah et.al treated some most commonly used diacritical marks for Arabic script and their mapping is based on elliptic representation [Kherallah et.al, 2009]. The above short review on Arabic script base languages shows that very less attention has been paid to deal with diacritical marks. While there are lot of issues are involved in diacritical marks especially in Nasta'liq script. The Arabic script based languages are rich in diacritical marks; it is very difficult to handle the diacritical marks for some character. The position, order and size are not enough to estimate the associated character. Husam et.al separated the punctuation marks at first stage from the character to reduce the complexity in the

handwritten text then they removed the diacritical marks. Vertical histogram on base character technique is used in locating the prospective segmentation point because the punctuation marks cause incorrect segment point due to the position of diacritical marks. The other hand the author did not treat all the diacritical marks i.e. diacritical marks used accent. Based on the density of the diacritical marks, diacritical marks are removed from ghost character but this step is only for the preprocessing phase [Husam et.al 2010]. Sternby et.al extracted diacritical marks from the base shape and treated separately. They did not tread like base shape thus all letters with different diacritical marks and same basic shape shared the same template. Later on the diacritical marks are attached on to the primary stroke with respect to their relative position shown in figure 2.2 [Sternby et.al, 2009].

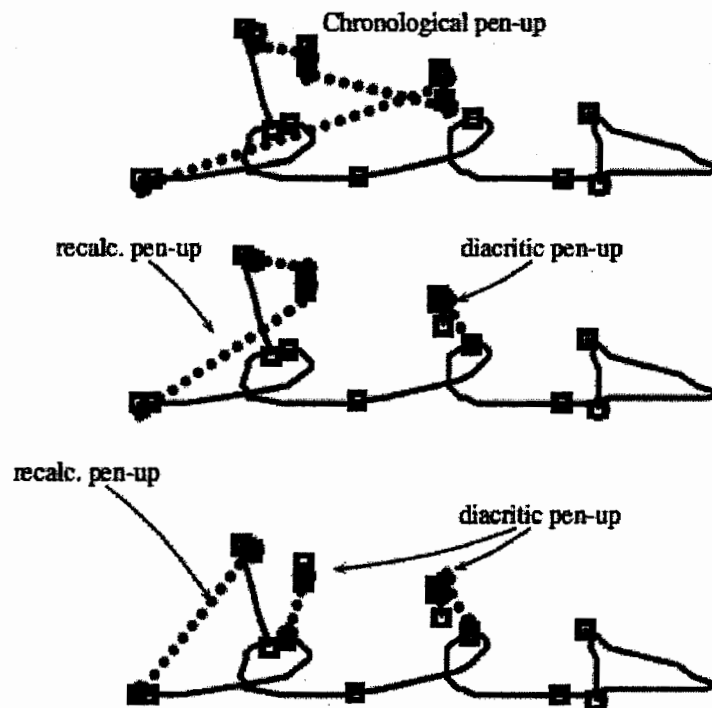


Figure 2.2: Diacritical marks handling [Sternby et.al, 2009].

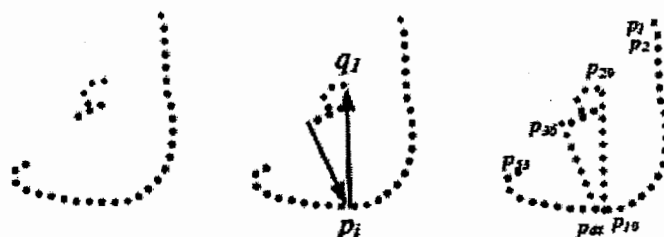


Figure 2.3: Vertical Projection of Diacritical marks [Biadisy et.al, 2006].

Biadisy et.al handled the delayed stroke based on the projection onto the primary strokes shown in figure 2.3. The diacritical marks detection is performed based on size, location of stroke and order of strokes [Biadisy F. et.al, 2006]. Whereas the position is estimated using the vertical projection of diacritical marks on the main stroke. The first point of delay stroke is projected on to the ligature body. Normally the Arabic character appears as part of connected word while they are written using different strokes. This is handled by projecting the starting point of each delayed stroke on to the word part body and integrating it with other later. Mozaffari et.al extracted the dots from the word to reduce the lexicon. The dots position is estimated using based line i.e. up and down [Mozaffari et.al, 2008]. The output of segmented dot is like 1U2D3U1D (1 up, 2 down, 3 up and 1 down). The dots are separated as they appear in the word in sequential order.

Hamad and Zitar presented density shape based approach for the diacritical marks segmentation for Arabic character recognition. As the density shape of diacritical marks smaller than Arabic character. Thus the diacritical marks of density size smaller than constant c are removed by utilizing the x, y coordinate information shown in figure 2.4. The disadvantage of this approach is that the smaller characters are also removed,

whereas in the case of Urdu script, few diacritical marks have the same size and same shape as of primary size [Hamad and Zitar, 2010].

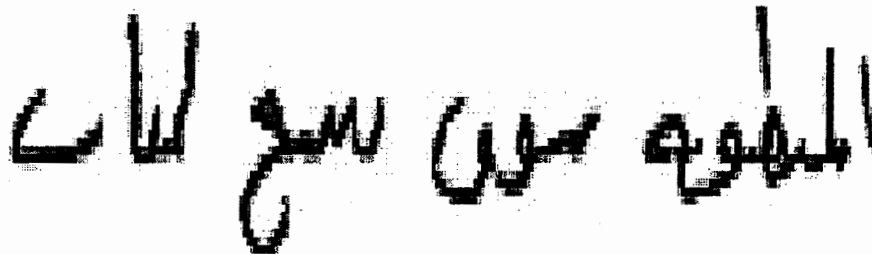


Figure 2.4: Issues diacritical marks separation [Hamad and Zitar, 2010].

2.2 Feature Extraction

The main aim of feature extraction phase is to extract that patterns which is most pertinent for classification. J. Sternby et.al. extracted relative features i.e. relative horizontal position, relative vertical position, with the mean value of segments and corresponding mean horizontal value and each segment is identified by feature angle, arc type, connection angle, length ratio and relative position [Sternby et.al, 2009]. A. Borji et.al extracted simple cells and grow these cells by connecting the associated cells with each other to form the features [Borji et.al, 2008] according to biological visual perception. Malaviay and L. Peters presented a layered approach by extracting features from level 0 (strokes elements) towards unique identified patterns. From the strokes elements, geometrical features are computed by using the fuzzy membership function, and these small geometrical features are combined to form global features by using fuzzy primitives [Malaviay et.al, 1999]. In order to obtain 24 dimensional feature vectors J. Schenk extracted three dimensional feature vectors from both offline and online domain [Schenk et.al, 2008]. The online features are pen pressure, velocity, (x ,y) coordinates,

difference of angles, angle between lines, etc. and for offline features, the strokes is sub sampled into 3x3 along pen trajectory and ascenders and descenders [Schenk et.al, 2008].

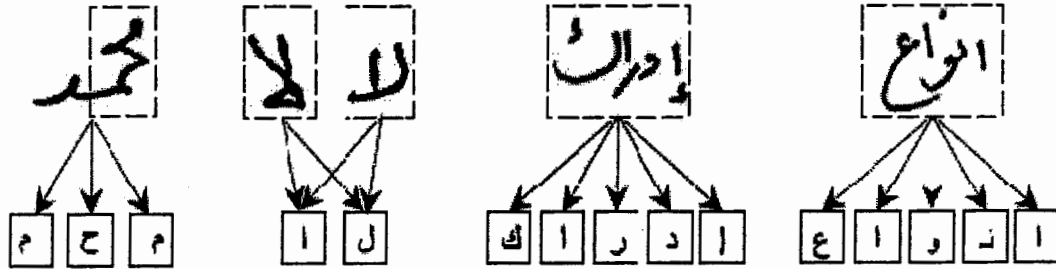


Figure 2.5 Arabic ligature and their constituent character [Hamad and Zitar, 2010]

El-Anwar et.al proposed three types freeman chain codes, long strokes, short strokes and eight pen-up for ensample movement with additional pen down information of succeeding stroke and current stroke [Randa et.al, 2007]. Malik and Khan used three directional features with other features like slope, writing direction, size start and ending coordinates. Hussain et.al extracted 20 unique shape defining features from primary stroke. Sudo T. (2002), proposed two types of pen pressure as a feature for Kanji characters [Sudo, 2002]. One is the pen pressure representing the pen ups and downs in a continuous manner and other is the time derivative of the pressure representing the temporal pattern of pen pressure. Biadisy et.al extracted local-angle, super-segment and loop-presence. Benouareth et.al used uniform and non-uniform segmentation scheme and extracted both statistical and structural features for offline Arabic character recognition [Benouareth et.al, 2008]. Beta-elliptical representation is also used for features extraction from the handwritten strokes [Kherallah M, 2009],[Alimi, 2000],[Binet et.al, 1893],[Bezine et.al,2003],[Kherallah et.al, 2004] and in the second phase Kherallah et.al

transformed the beta-elliptical model trajectory by visual codes based on psychological and cognitive domain and the curvilinear velocity is computed using a second-order derivative. Husam et.al presented efficient neural based segmentation for Arabic handwritten word recognition and the modified direction feature extraction technique combines the local feature vector and global structural information. First the contours are extracted and then directions of line segments comprising the characters are detected and the foreground pixels are replaced with the direction values [Husam et.al, 2010].

Lui et.al extracted gradient based features and for orientation Robert and Sobel operators are used to calculate the gradient component in orthogonal direction. The gradient vector is assigned to discrete direction either by parallelogram decomposition or tangent angle quantization [Lui et.al, 2009]. Sternby et.al extracted relative horizontal and relative vertical position by signifying the mean vertical value of segment and each segment is represented by specified feature angle, arc type, length ratio, connection angle and the relative position shown in figure 2.6 [Sternby et.al 2009].

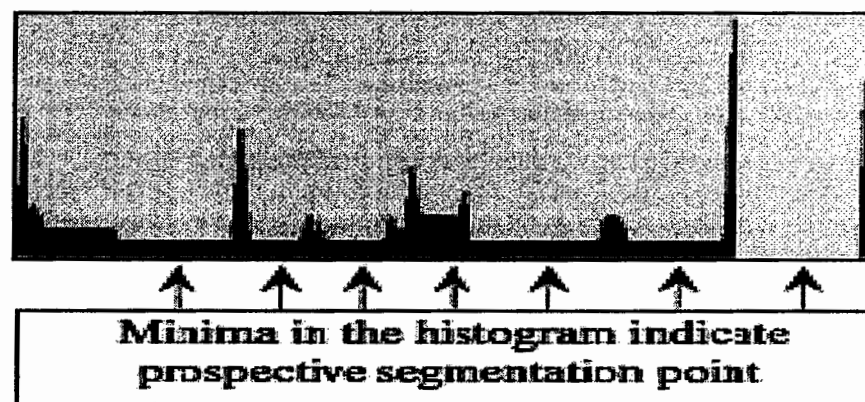


Figure 2.6: Segmentation point estimation based on histogram [Sternby et.al, 2009]

Borgi used the Gabor and DoG filters as in the feature extraction phase to implement biological inspired system [Borji et.al, 2008]. Izadi et.al build shape descriptors and define relative context as in the pair wise distance and angle and rational property spanned by relative context making it possible to include different levels of detail [Izadi et.al, 2008]. Amor and Amara extracted features using Hough transformation which correspond to major line features [Amor and Amara, 2006]. Aly et.al used a technique that normalize normalized the text in bounding box by estimating the partial height of character using which the normalized size and height of the characters are defined [Aly et.al , 2007]. Hanmandlu et.al presented box method for feature extraction for offline unconstrained character recognition and 24 pairs are extracted from handwritten text and all the features are organized in sequential order [Hanmandlu et.al, 2003].

Deepu et.al extract sequence of feature which is used directly for classification and preserved in order in which they appear and PCA is also used with pre-clustering [Deepu et.al, 2004]. Khedher and Al-Talib combined statistical and structural method for feature extraction and 12 features are extracted from both main and secondary parts of the characters. For secondary stroke limited features such as height, width, height to width ratio are extracted [Khedher and Al-Talib, 2004]. Baghshah et.al extracted directional features and relative vertical and horizontal motion features from online handwritten for Persian character recognition. [Baghshah et.al, 2006]. Parui et.al extracted features in sequential order for online handwritten Bangla character and represent the whole stroke using N points whereas for sub stroke sequence, straight line segments are used [Parui et.al, 2008]. Husam et.al presented segmentation based approach combined with neural

network for handwritten Arabic recognition [Husam et.al, 2010]. The words are segment into uniform and non-uniform segmentation. The segmentation points are validated using neural network by fusing the confidence value and several directional features (manifold direction, normalized direction value, etc.) are extracted for successful segmentation.

Features extracted [Hussain et.al], [Malik et.al, Borji et.al] do not fully define the shape of the character in the case of segmented approach; and also fail to recognize the primary stroke when excluding the secondary information. The pen pressure features (Sudo T. 2002) are not useful for feature extraction of Urdu handwritten text due to writer dependent force during writing [Biadisy et.al.] and will not work properly for cursive Urdu script as it is very complex to segment using these features. They also extracted feature matrix based on geometrical processing, while automatic feature extraction based on recognition engine may be used to extract a better feature matrix.

2.3 Classifier

Most handwritten text recognition systems are based on the statistical method or rule based methods. Normally statistical recognizers are more reliable however they also have some disadvantages, as they are very complex and computationally difficult to train and classify a large number of shapes using them because a large vocabulary is needed. The problem in fuzzy rules is the lack of training and it is impossible to form an exhaustive set of rules that can model all possibilities.

Hussain et.al presented segmentation free approach for Nasta'liq font and used backpropagation neural network for classification and recognize 240 basic ligatures and used dictionary based post processing for 18000 words with accuracy of 93%. Malik S.

and Khan S.A recognized basic alphabets and numerals and used tree based dictionary search for word classification and reported 93%, 78% accuracy for numerals and character respectively. Benouareth et.al modeled discrete structure right to left HMM's with explicit state duration and modified viterbe algorithm, while gamma, Gaussian and Poisson distribution are used for state duration modeling [Benouareth et.al, 2008]. K-mean clustering is used for vector quantization of 32 statistical and 6 structural features for Arabic word recognition. Gallies used implicit segmentation based on discrete HMM for cursive word recognition [Gillies et.al, 1992]. Khorsheed presented single left to right HMM with structural features for offline Arabic recognition and each HMM represents one letter from Arabic alphabets [Khorsheed, 2003]. Pechwitz and Maergner presented semi continuous HMM for offline Arabic handwriting with 7 states and each states has three transition i.e. to itself, next and skip state [Pechwitz et.al, 2003]. HMM is the dominant classifier for cursive handwritten recognition especially for online recognition due to additional timing information. Rigoll presented a novel recognition system for online handwritten mathematical expressions based on HMM's by simultaneous segmentation and recognition [Kosmala et.al, 2004]. The segmentation and recognition result is used for the interpretation of the symbols and their spatial relationships. Shu H. described Fujitsu's online handwritten character recognition system and some application software that adopts HMM technology and claimed accuracy of 94.6% for Japanese text [Shu, 1996]. Akira presented new method for on-line handwriting recognition of Kanji characters by employing sub stroke HMMs as minimum units and direction of pen motion is utilized as features [Mitsuru M et.al, 2002]. Mohamed, M. and Gader, P proposed continuous HMMs for handwritten Turkish words using segmentation-free

modeling. They used 12 states HMM for each character and observations symbols are based on transition from white-black and black-white [Mohamed et.al, 1996]. Chen, M.Y. et.al presented explicit continuous density duration HMMs and modified viterbi algorithm. The observation symbols were based on geometrical and topological features. The words are segmented into characters and each sub character is identified with a state which can account for up to four segments per letter [Chen et.al, 1995]. Kherallah et.al presented online Arabic character recognition system based on genetic algorithm and video encoding for different shapes of Arabic script. To estimate the evaluation function, visual indices is developed [Kherallah et.al, 2009].

Husam et.al used neural network for classification of online Arabic recognition based on segmentation and finally the result are fused using fuzzy rules. First the segmentation area of small dimension centered about each heuristic segment point is extracted and classifier is trained on segmentation area either valid or invalid. The segmentation area is extracted from the current segmentation point and previous segmentation point and the second area is segmented about heuristic point [Hussam et.al, 2010].

Al-Habian and Assaleh presented HMM based online Arabic character recognition and decision logic is used to interpret the output of HMMs and converting its output into recognized words shown in figure 2.7. [Habin and Assaleh et.al, 2007]. Deepu et.al used PCA for online handwritten character recognition for Tamil script to find the bases vector for each subspace. [Deepu et.al, 2004]. Ahmad et.al, segmented the printed Urdu character based on the complexity of the word calculated from the topological features i.e. hole, width, height [Ahmed et.al 2007]. To improve the

efficiency, both horizontal and vertical scanning of the word is performed for complex characters. Shahzad et.al presented isolated Urdu character recognition for Urdu Qaeda using Rubine features [Shahzad et.al 2009]. Sattar et.al presented finite state model for printed Urdu characters [Sattar et.al, 2009]. The finite state recognizer remain in a state to accept new input, if the input is white space, it mean incomplete word or accepted state.

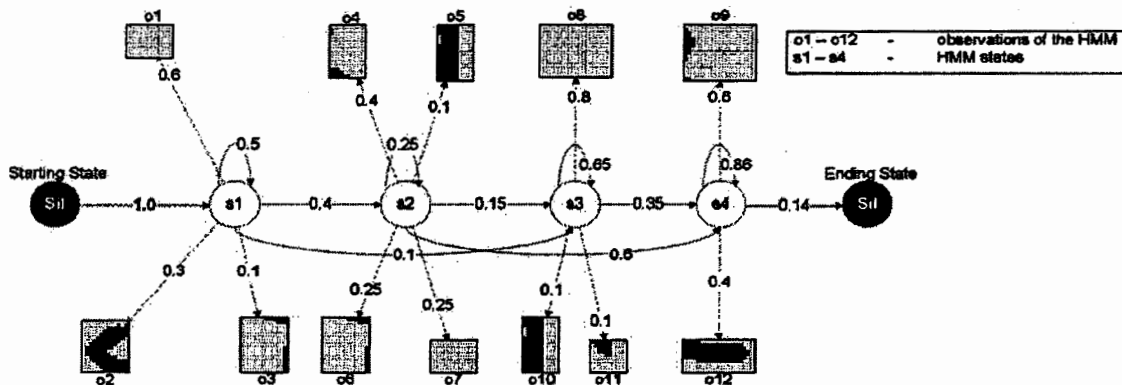


Figure 2.7: HMM modal for digits 4 [Awaidah and Mahmoud, 2009]

Lui et.al used six type of classifier MLP neural network, modified quadratic discriminant function, discriminative learning quadratic discriminant function, polynomial network classifier, SVM classifier and class-specific feature polynomial classifier [Lui et.al, 2009]. MLP is used with one hidden layer with error calculated by back propagation. Sternby et.al presented segmentation based technique for Arabic script based on template matching and each template segment is modeled separately by setting the maximum number of cluster [Sternby et.al, 2009]. Lin et.al presented biological inspired CAPTCHA recognition system for English text i.e. luminance, chrominance and knowledge base filters are applied to model human vision and an efficient biological inspired segmentation technique is presented [Lin et.al 2008]. Wu et.al presented biological inspired hierarchical feature extraction method for localization and they

perform a hierarchical search from coarse to fine in order to minimize the computational complexity. Dong et.al presented biological inspired character recognition for natural image classification and the features are extracted from the image to find saliency map for each image [Dong et.al, 2007]. System simulates the human-like image classification by using the biological inspired knowledge. Huang et.al presented enhanced biologically inspired model which removes uninformative input by imposing constraints and utilizing the weight locally for polling operation with physiological motivations. Moreover it applies the feedback procedure that finds the effective features. Borji et.al presented bio-inspired character recognition system by extracting biologically inspired feature extraction. The scale and translation invariant features are extracted for Farsi characters and standard classifiers such as KNN, SVM are used and for feature extraction, Gabor and DoG filters are used [Borji et.al, 2008]. Al-Khateeb et.al presented a comparative study for handwritten word Arabic feature extraction and different technique are [AlKhateeb et.al, 2009]. Hamada et.al presented biologically inspired model for objection recognition and Gabor filter is used for feature extraction and a multilayered bottom up approach is presented. The model is based on feed backward and feed forward which is similar to human feed forward and backward recognition [Hamada et.al, 2009]. Izadi et.al presented online writer independent character recognition using relational context representation and a SVM classifier is used [Izadi et.al, 2008]. The system is very helpful for a series of recognition task. Amora and Amara presented multifont Arabic character recognition system based on a hybrid method. HMM and ANN are used for classification based on line features extracted using Hough transformation where ANN is responsible for 28 observation to HMM [Amor and Amara, 2006]. Wang et.al presented multiple

classifiers based system for handwritten word recognition to increase the accuracy. The results of multiple classifiers are combined and they use institutive run time weighted opinion for large vocabulary. [Wang et.al, 2002]. The final decision rule is based on discriminant function computation and multiple results are combined using modified ROVER algorithm. Hanmandlu et.al used multilayered BPNN with the input of 48 features extracted from box and the input layer contains 48 neurons, 32 and 16 neurons in two hidden layers and 4 neurons are in the output layer and they used 24 fuzzy sets [Hanmandlu et.al, 2003]. Deepu et.al used PCA for online handwritten character recognition for Tamil script to find the bases vector for each subspace [Deepu et.al, 2004]. Mitoma H. et.al used Eigen faces for handwritten online character recognition and elastic based matching technique are used and category specific deformations are employed to suppress misrecognition. [Mitoma H. et.al, 2004]. Khedher and Al-Talib present fuzzy expert system and fuzzy rules are constructed using AND and OR operations [Khedher and Al-Talib, 2004]. Baghshah et.al used fuzzy LVQ for online Persian character recognition to classify the end token and similarity of two token is calculated using FLVQ [Baghshah et.al, 2006]. Hamadani et.al presented character recognition system by combining both online and offline features and HMM classifier is used and finally classification is compared with single domain systems [Hamadani et.al, 2009]. Husam et.al presented segmentation based approach using neural network for handwritten Arabic recognition [Husam et.al, 2010]. The word is segmented into uniform and non-uniform segmentation. The segmentation points are validated using neural network by fusing the confidence value.

Husam et.al presented several rules to check segmentation area, centered region and right region. Neural network is used to verify the segmentation area which characterized the segmentation area into valid or invalid by analyzing the confidence value and the fusion is performed on the expert value [Husam et.al, 2010]. Lui et.al used three databases for post processing and performance evaluation and proposed that by combining the multiple classifiers the error rate can be reduced [Lui et.al, 2009]. Daifallah presented segmentation algorithm for Arabic online handwritten recognition [Daifallah et.al, 2009]. It is based on arbitrary segmentation followed by segmentation enhancement, consecutive joints connection and finally segmentation point locating. HMM is used for classification of independent subunits. Mezghani performed several preprocessing steps i.e. smoothing, normalization, etc. and presented Bayes classification for online Arabic character recognition using tangent differences histograms and Gibbs modeling of the class-conditional probability density functions [Mezghani et.al, 2008].

Baidy F. et.al used word and word part dictionaries and Arabic dictionary is subdivided into sub dictionary and HMM is used for word modeling [Baidy F. et.al, 2006]. Strenby et.al used static dictionary lookup, the characters are directly constructed from segmentation graph and the system does not support missing characters. The dictionary is complemented with the words in the test set to avoid the error of dictionary [Strenby et.al, 2009].

2.5 Final Remarks

The above study shows that character recognition has been an intensive field of research from early days of computer. Imitating the human abilities in computers is not

an easy task due to high context sensitivity and has extra ordinary viewing devices. The above literature shows that the work done for Arabic script based languages is scarce especially for handwritten character recognition. This may be due to the complexities of this script. All Arabic script based languages are following approximately same character set but written in different styles. From the family of Arabic script based languages, Urdu is more complex and written in Nasta'liq style. In Nasta'liq, the shape of character depends upon the associated character on both side and it may be up to 32 letters. Urdu script contains more than 22000 ligatures and to train a classifier for such a large amount of ligature is a very complex task while most of the characters share same basic shape with Urdu and with its family. Literature shows that a lot of work is required in every phase of the Urdu character recognition system, especially to treat the diacritical marks. As the Arabic script based languages are rich in diacritical marks some researcher has proposed method for diacritical marks separation. These methods can be applied on Naskh writing style while due to the complexity of Nasta'liq and the presence of more diacritical marks for Urdu script; the presented techniques in literatures are not suitable.

On the other hand researchers i.e. computer scientist, neurologist are also working to model the human visual concept in the field of character recognition. The biological theories behind the image processing techniques are an intensive field of research to attain an efficient system like the human visual system. Although there are many complexities involved in modeling the human visual perception but the human reasoning can be utilized to obtained better results close to the natural perception. Due to the large variations in script and writing styles, fuzzy logics can be incorporated for better results.

CHAPTER 3

FUZZY LOGIC AND HUMAN BIOLOGICAL VISUAL PERCEPTION

This chapter gives the background knowledge regarding the fuzzy logic, analysis of human visual perception and its computational modeling. The first part of this chapter presents fuzzy logics whereas the second part describes the terminology and concept of biology of human visual perception. Fuzzy logic is an important classification tool and has been used as a means of interpreting vague, incomplete, noisy and contradictory information into a compromised rule [Pham and Chen, 2002]. The real world is dynamic, enormous and very complex. To deal with it, biological systems have been evolving and self-adjusting over millions of years and adapting to the changes in the environment. Imitation of human capabilities by machine has been field of intense research since the early days of computer and amongst them one of the interesting issues is human vision modeling. In this regard, biology has been an important source of inspiration in image processing and pattern recognition applications. On the other hand the understanding of visual information processing

and perception is one of the challenging tasks of contemporary science. Bio-inspired pattern recognition has been under consideration since mid-eighties not only by computer scientists but also by physiologist and neuroscientists. The deeper review of human vision has helped to advance pattern recognition research, thus making it more robust. Image processing in humans is done in parts at different levels and in different layers. The integration of human vision system in traditional approaches to pattern recognition helps to achieve robustness in term of accuracy and efficiency. The visual perception of humans is very powerful in pattern recognition applications. The selectivity, transformation, speed and context knowledge are the most important features of human visual perception. Moreover, a human is able to detect the familiar and unfamiliar objects even in unfamiliar environment. A computer can perform complex problems more precisely, efficiently and much faster than a human. However, there are certain different types of complex algorithms like pattern recognition and image processing situations where a few months old child can perform much better than the state-of-the-art computer with latest algorithms available today. In order to better understand how human brain performs this extraordinary function the human visual system is briefly described in section 2. While dealing with images, the real time data may be ambiguous and/or incomplete. This may be due to several reasons e.g. blurred edges and/or irregular and missing patterns. Thus, to cater for these data ambiguities and imperfections, it is suggested to model human visual perception in pattern recognition application such as character recognition. Human visual perception can be modeled using fuzzy logic, evolutionary computation and neurocomputing.

Biologically inspired solutions are gaining more interest with wide spread use of artificial intelligence applications [Halem, 2009], [Lin C.W et.al, 2008], [Bosner,

2006], [Dong et.al, 2007], [Wu et.al, 2006], [Bosner et.al 2006], [Forbes, 2005], [Collins et.al, 2004], [Benyus, 1997], [Anastas, et.al, 2000], [Papanek, 1984]. The biology inspired image processing need expertise in two fields, human biological perception and modeling the concept in a computer. Due to the gap between a computer scientist and an expert in biology, both speak different scientific linguast, creating a communication challenge and both use different methods of investigation [Helms, 2009]. To deal with complex problems and involve human-like intelligence, fuzzy logic is a very good approach and due to this property fuzzy logic has been used in many applications of various areas such as pattern recognition, control, information retrieval and quantities analysis, etc.

3.1 Fuzzy Logic

Conventional logic is based on Boolean algebra and this logic is better for modeling the well represented behavioral mathematical models. Whereas, human decision power cannot be represented through Boolean representation thus, the legacy programming languages are not sufficient for modeling the human reasoning because the humans rely on relative instead of discrete values. Fuzzy logic provides the solution to model the human reasoning in multiple categories rather than two categories i.e. very small, small, medium, large, etc. whereas the binary logic is limited to two states: large or small.

In the early 1950s, Herbert Simon, Allen Newell and Cliff Shaw conducted experiments in writing programs to imitate human thought processes. The basic purpose of fuzzy logic was to implement human reasoning power through computer. The comprehension of computer is limited to 0 and 1 whereas the human rely on multiple classes whose boundaries are imprecise or not well defined.

Their experiments resulted in a program called Logic Theorist, “which consisted of rules of already proved axioms. When a new logical expression was given to it, it would search through all possible operations to discover a proof of the new expression, using heuristics. This was a major step in the development of AI. The Logic Theorist was capable of quickly solving thirty-eight out of fifty-two problems with proofs that Whitehead and Russel had devised [Al-Omari, 2004]. At the same time, Shanon came out with a paper on the possibility of computers playing chess [Salah, 2002]. Fuzzy logic is an important classification tool and has been used as a means of interpreting vague, incomplete, noisy and contradictory information into a compromised rule [Pham and Chen, 2002]. It has been applied in many applications like medical imaging, document imaging, robotics applications, etc. [Pham and Chen, 2002], [Mahmoodabadi et.al, 2010], [Demirci, 2010], [Kazemian, 2005], [Patra et.al, 2009], [Ozturk, 2009].

The fuzzy logic is not based on probability i.e. a fuzzy system deals with deterministic plausibility whereas the probability deals with non-deterministic likelihood [Lin, 1996]. For example, the doctor predicts the patient’s diseases by using if-else reasoning, taking into account a number of variables. Probability is the number that measures the certainty of an event whereas the fuzzy logic measures the degree of certainty of particular event. Probability is only meaningful for things that have not happened yet while fuzzy set membership function caters after that [Eberhart, 1996].

The following are the main reasons that make fuzzy logic very helpful in pattern recognition problems.

1. As fuzzy does not require precise results thus it is more robust to noisy and missing data environment.
2. Since the user made rules are governed thus it can be easily changed to improve and new data can easily be incorporated.
3. Due to the rule based approach reasonable inputs and outputs can be processed. Thus it is better to break the problems into smaller chunks of data/categories based on fuzzy sets.
4. It can better modal non-linear systems that are very difficult or impossible to model mathematically.

Fuzzy logic is a multi-valued logic derived from fuzzy set theory to deal with reasoning. Fuzzy logic is approximate rather than exact. The membership function is not only 0 or 1 but have degree of truth of statement. The degree of truth ranges from 0 to 1. The membership function may be discrete or continuous. A continuous memberships function is a mathematical function for example like bell shape, triangular shape, etc. whereas the discrete membership function is a vector.

For reasoning the fuzzy logics use the IF-THEN rules. It is alternative to traditional notion and it is leading towards artificial intelligence. It provides a simpler ways to find conclusion based on noisy and missing input. Fuzzy Logic is the simple IF-THEN rules to solve the complex problems rather than complex modeling i.e. HMM, ANN.

3.1.1 Linguistic Variables

Linguistic variable represents the crisp information in precision and is essential for fuzzy logic. Unlike conventional variables in mathematics, the fuzzy logic also uses non-numerical linguistics variables to model the natural languages e.g.

“Saeed is tall”, “today is hot”. The linguistic variables such as height may be “tall”, “short”, “small”. The linguistic variables can be modified via linguistic hedges. For example, suppose x is the linguistic variable with label weight; and the fuzzy terms of these linguistics variables are light, normal, heavy and very heavy. For example each term in fuzzy variable is ranged from 30Kg. to 100 Kg. as shown in figure 3.1.

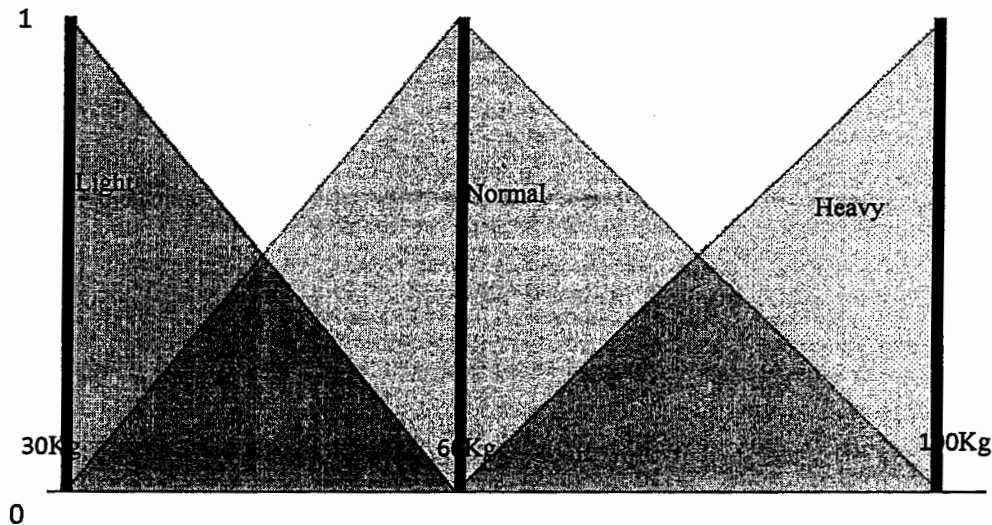


Figure 3.1: Fuzzy Membership function

Due to its tremendous properties, fuzzy logic is used in pattern recognition, image processing, quantitative analysis, control inference and information retrieval. The fuzzy logic is appropriate in the following problems [Lin, 1996].

1. The problems which are concerned with continuous phenomena and cannot be easily broken into discrete steps.
2. The problems where mathematical model is very complex or does not exist
3. The high noise problems.
4. The problems which involve human interaction.
5. To model the human behavior.
6. The problem where experts are available and can model the rule to implement the natural behavior.

The main disadvantage of fuzzy logic is the lack of learning capability to tunes the fuzzy variables and membership function. Basically, fuzzy variables and membership functions are defined by experts according to the application requirement. Thus, it is very difficult, and some-times impossible to define and model all the rules and membership functions, required for application due to complexity, uncertainty and ambiguity. The lack of learning capability of fuzzy logic make it little weaker but the human reasoning involved in fuzzy logic makes it power full tool to solve complex problems.

3.2 Human Visual Perception

Computers can solve most problems quite precisely, efficiently and much faster than humans however there are still many problems such as pattern recognition where few years old child can perform much better than the latest algorithms available today. Character recognition is one of the important pattern recognition applications. The handwritten script may be ambiguous and incomplete due to several reasons i.e. blurred edges, irregular and missing patterns. Even Humans may find it difficult to read handwritten script. Use of context knowledge, adaptability and learning capability leads humans to robust recognition in complex scenarios. The interpretation of complex patterns and transformation of these complex patterns into behaviorally understandable signals is important in our daily life. The human brain has achieved this robustness through thousands of years of evolution. Even the brain of very small animals is so complex that it is very difficult to imitate it in the artificial system [Mumolo, 2006]. To understand how human brain can achieve these extraordinary abilities and to model the biological visual perception has been a

research issue from many decades. Biological study of human vision has provided important source of inspiration in the field of pattern recognition.

The visual perception system of humans is very powerful in pattern recognition problems. The selectivity, transformation invariance, speed and context knowledge are the most important features of human visual perception. Moreover, human is able to detect the familiar and unfamiliar objects even in variable environment. It has exquisite selectivity and context knowledge that helps to identify even similar shapes such as twins' face identification. Moreover, humans can distinguish more that tens of thousands of different shapes [Biederma, 1987] while it is very difficult (almost impossible in current age) to build a generalize system that works like them. The computer system developed so far can identify only minor changes i.e. handwritten characters where machine have any context knowledge, face recognition problems where system can identify slight expression, pose or illumination variation. Human perception is very fast and can recognize objects in 100-200 ms [Thorpe, 1996], [Potter, 1969]. The additional context knowledge of human makes the recognition very easy i.e. the missing character can easily be interpreted from the sentence knowledge. Moreover the transformation invariant feature makes the human visual system more robust. The object remains the same even after changing its position, scale and rotation. Moreover, the human neurons are tuned to become selectively efficient for the shapes that are more frequent [Kuo et.al. 2003], [Baker et.al, 2007].

3.2.1 Autonomy of Human Visual System

The human visual system has been under study for approximately four decades to find the response of neurons in different part of the cortex. To model the

human visual perception into computer vision applications, the detailed knowledge of human visual system is very useful even though it is very difficult to attain the level of human brain. When the light from an object falls on to the retina it is converted into electrical signal. The retina is driven from the CNS (central nervous system) to process the incoming light into signal and conveyed to the brain for further processing by retinal ganglion cells [Rieke, 2006],[Meister,1995],[Koch, 1982]. The information is sent to the V1 (primary visual cortex) after processing through the thalamus [Reinagel, 1999], [Sherman. 2001], [Alonso, 1996], [Lesica, 2004]. The primary visual cortex also called striate or V1 cortex receive feedback information from higher cortical areas that consists of six layers. Dorsal and ventral are the two information processing units after V1 [Mishkin, 1982], [Haxby , 1991]. The dorsal is responsible of spatial localization of object in the environment and guiding the actions towards the object whereas the ventral is involved in recognition of the objects [Goodale, 1992] and both are dependent on each other. The thalamical component called lateral geniculate nucleus (LGN) send the information to the cortical area V1. Thus LGN receives the output of retina. V1 sends the signal to ventral through visual areas V2 and V4 and also projects back to thalamus. The information is further processed from V1 to V2, V3, V4, IT (infero-temporal cortex). V1 is responsible for static and moving object processing and is excellent in pattern recognition. A frog's behavior for catching with and without barrier diagram is shown in figure 3.2.

Generally the primary visual cortex neurons are more selective to simple features than neurons in the higher cortex area. V1 and V2 have many common properties. V2 responds to spatial properties i.e. contours, orientation and whether the stimulus is the part of background or object [Borji et.al, 2008].

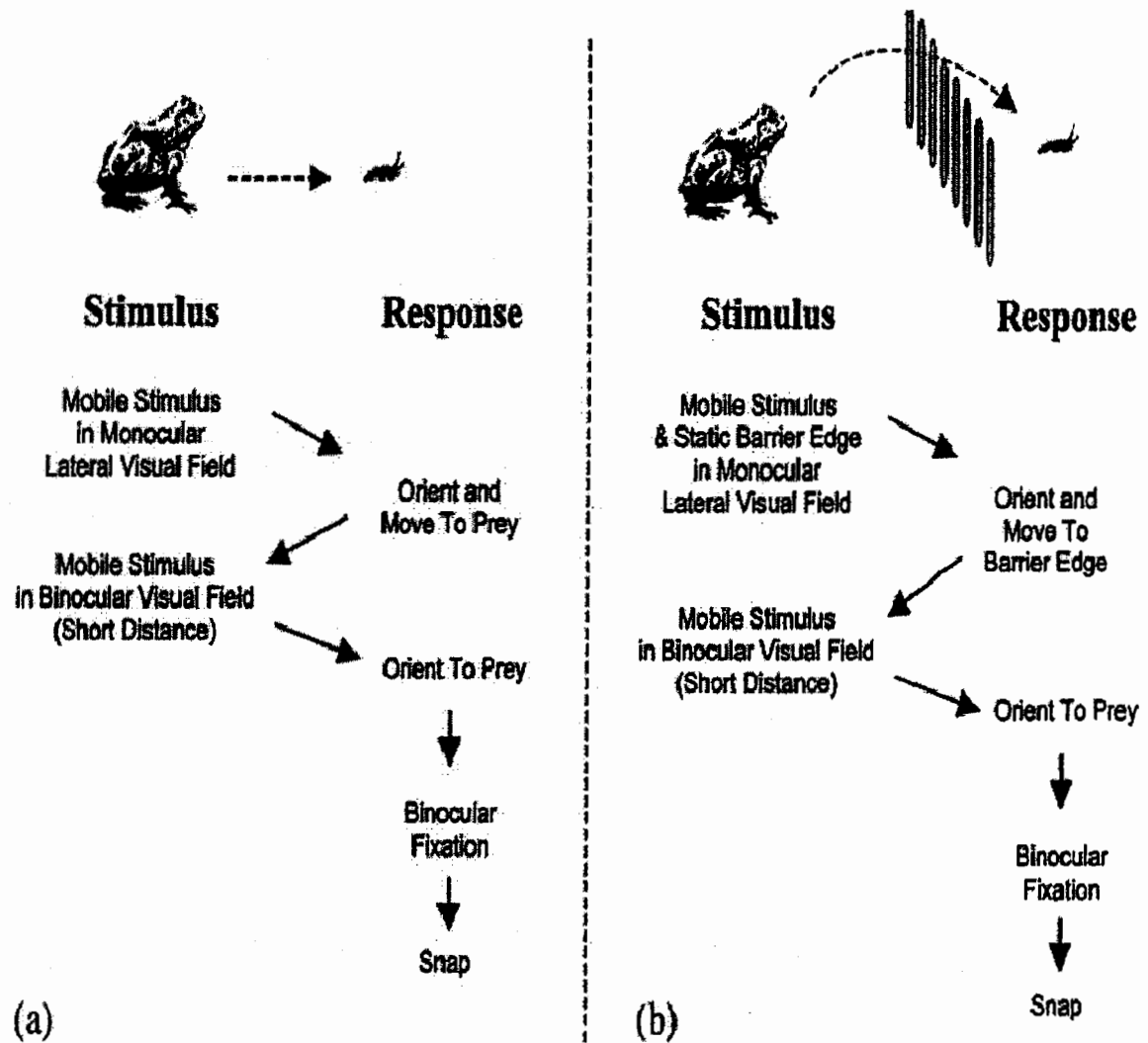


Figure: 3.2: Stimulus diagram for frog. (a). prey acquisition behavior (b) detouring around the barrier in reaching stimuli [Weitzenfeld, 2008]

The cortex area V3 is responsible for motion detection, visual stimuli combination and also deals with color sensitivity whereas V4 neurons are responsible for spatial frequency, orientation and color. The visual area MT (V5) is sensitive to motion perception and the integration of local motion signals into global signals. Inferotemporal cortex is the last stage of ventral stream and responds to complex signals like handwritten marks, faces, etc.

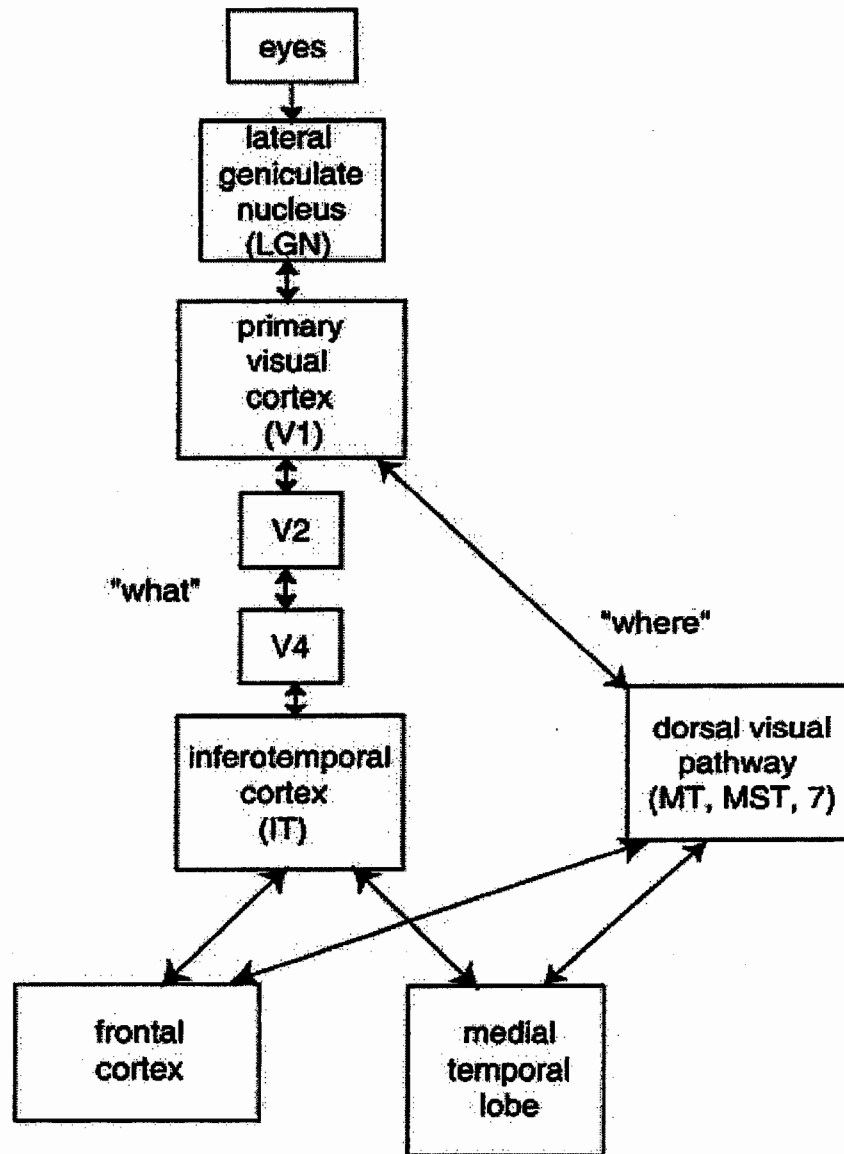


Figure 3.3: The autonomy of human visual object recognition [Borji, 2008]

Dorsal stream also called *where stream* starts from V1 and passes through V2 and then to dorsomedial area and MT. It is associated with special properties, object location, control and focusing of eyes. Whereas ventral stream also known as *what pathway* starts from V1 and goes through V4 via V2 and finally to infero-temporal. Mainly, it is responsible for object identification by extracting the features and performing step by step categorization. The complex objects are classified in IT and is sensitive to complex patterns [Borji, 2008].

3.3 Bio-Inspired Computer Vision

The term “bio-inspired” is also referred to as “biomimetic” or “bionic” however the former is preferred [Zhou, 2000]. Due to the high growth of computer vision and pattern recognition applications, classification has become one of the important and critical issues. Lots of research has been done and a number of algorithms have been proposed. The biological theories behind the computer vision algorithms are an intensive field of research to attain an efficient vision system like that of the humans. There are many complexities involved in the modeling of human visual perception. The traditional image processing systems are based on 512×512 whereas the human visual system has 126 photoreceptors. It is well known that better resolution of images helps to obtain better recognition result.

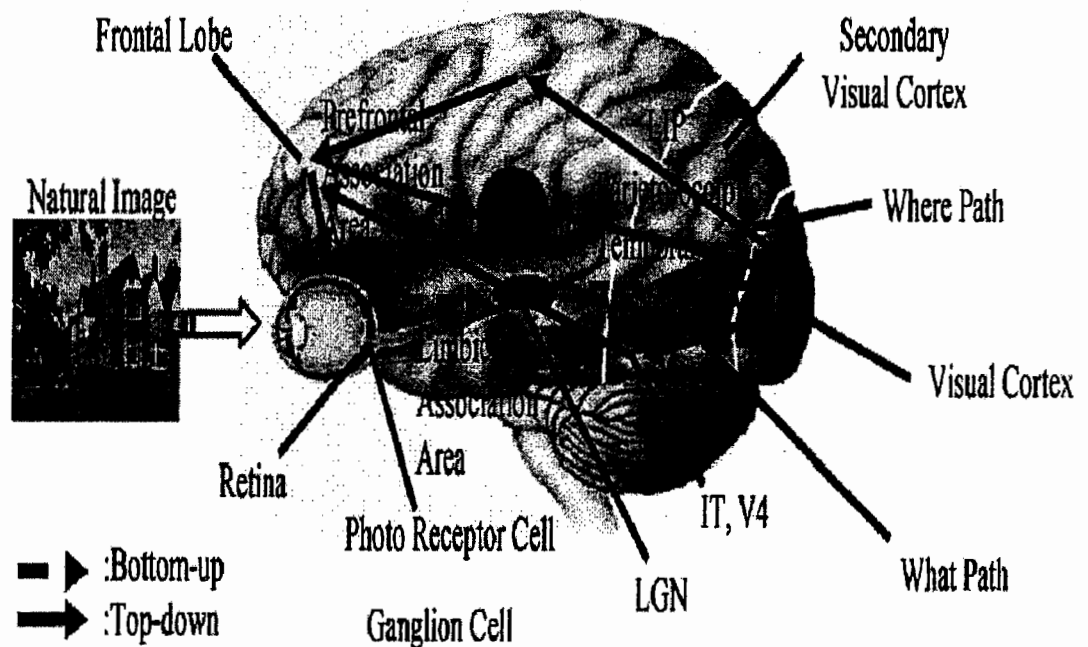


Figure 3.4: Biological perception mechanism of human mind [Dong, 2007]

Developing an efficient handwritten character recognition system is still a challenging task in case of many languages i.e. Arabic, Farsi, Urdu, Chinese, etc. Traditional character recognition systems face problem in handwritten script due to lots of variations and complexity. A lot of research effort has been done but still there is a room for improvement in the accuracy and efficiency and the research has focused to biological visual perception. Fuzzy logic and neuro-computation have shown promising results to cater to flexibility and ambiguity.

Not much research is available in character recognition inspired by biological vision. Omari and Jarrah presented probabilistic neural network for handwritten recognition of Indian numerals [Al-Omari, 2004]. Salah et.al presented selective attention based method for handwritten English character recognition [Salah, 2002]. The brain activity in major neural circuit in ventral and dorsal stream did not show much difference for English [Chan et.al, 2009] and similarly for Urdu and Arabic script. The brain activities related to difference of structure has not been investigated yet for Arabic script based languages. Lin et.al presented biologically inspired CAPTCHA recognition system for English text, i.e. luminance, chrominance and knowledge based filters are applied to model human vision concept and efficient biological inspired segmentation technique have been presented [Lin et.al 2008]. Wu et.al presented biologically inspired hierarchal feature extraction method for localization. They performed the hierarchical search in order to minimize the computational complexity. Dong et.al presented biologically inspired character recognition for natural image classification and features are extracted from the image to find saliency map for each image as shown in figure 3.5 and 3.6 [Dong et.al, 2007].

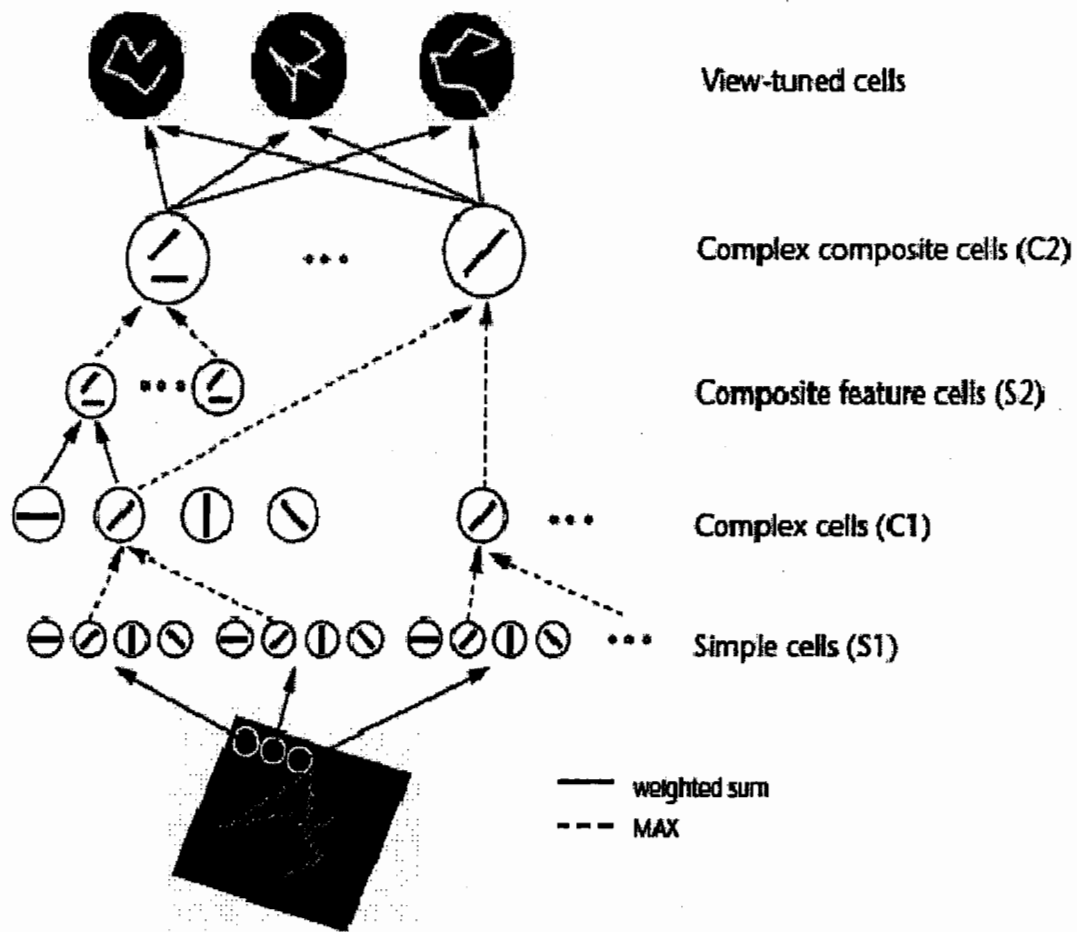


Figure 3.5: Ventral stream based object recognition [Riesenhuber and Poggio, 1999]

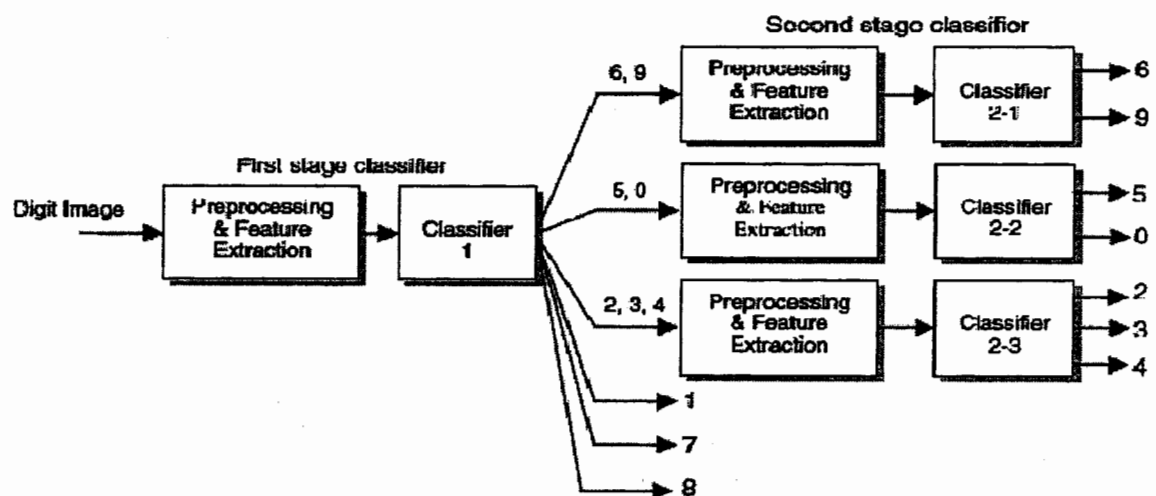


Figure 3.6: Bio-inspired Numeral recognition [Borji et.al, 2008]

A system simulates the human-like image classification by using the biologically inspired knowledge. Huang et.al presented an enhanced biologically inspired model which removes the un-informative input by imposing constraints and utilizing the weights locally for polling operation with physiological motivations. Moreover, it applies the feedback procedure that finds the effective features. Borji et.al presented a bio-inspired character recognition system by extracting biologically inspired features. The scale and translation invariant features are extracted from Farsi characters and standard classifiers such as KNN, SVM are used. For feature extraction, Gabor and DoG filters were used [Borji et.al, 2008]. Hamada et.al presented a biologically inspired model for objection recognition and Gabor filters were used for feature extraction and multilayered bottom up approach is presented. Finally, the model is based on feed backward and feed forward similar to human feed forward and backward [Hamada et.al, 2009].

3.4 Bio-Inspired Fuzzy Logic System

Modeling human visual perception in real world image processing algorithms with the help of some principle present in our visual system is an interesting field of research. The human biological vision can be modeled using fuzzy logic, evolutionary computation and neuro-computing. Fuzzy logic is the important tool due to reasoning and has high range of values instead of definite. There are two main reasons for success of fuzzy logic inspired biological vision.

1. The human visual system works layer by layer as discussed above i.e. from retina to V1 to V2, etc. This layer by layer processing can easily be modeled using fuzzy logic.

2. The fuzzy logic can better handle the noisy and irregular information by using the range of value instead of definite values.

Thus modeling fuzzy logic using human biological vision (the combination of fuzzy logic and human biological perception system) is more power full to handle the complex patterns than other conventional methods. Figure 3.5 briefly illustrates the layered approach by human ventral stream. The robustness is due to the reasoning mechanism. [Riesenhuber and Poggio, 1999]

3.5. Final Remarks

In this chapter an introduction to fuzzy logic and human biological vision system is presented. The fuzzy logic remains a powerful tool to solve the unconstrained complex pattern recognition problems. The handwritten character recognition in an unconstrained environment is very a complex pattern recognition problem due to large variations involved. The real time data may be ambiguous; incomplete due to several reasons i.e. blurred edges, irregular and missing patterns. This fact suggests that modeling of the human visual perception in real world pattern recognition applications such as character recognition by using principles present in our visual system can present improved results. The human visual modeling using fuzzy logics helps to handle complex problems in a better way. The happy relation of reasoning based fuzzy logics and biological vision modeling can better model the complex problems in layered approach.

PROPOSED METHODOLOGY

CHAPTER 4

GHOST CHARACTER RECOGNITION THEORY

Arabic script is used by more than 1/4th population of the world in the form of different languages like Arabic, Persian, Urdu, Sindhi, Pashto, Afghani, etc. but each language has its own vocabulary, writing rules and set of alphabets. The set of Urdu alphabets is a superset of the alphabets sets used by all other Arabic script based languages. Character recognition of Arabic script based languages is one of the most difficult tasks due to complexities involved in this script. This chapter presents a novel technique i.e. the Ghost Character Recognition Theory that helps develop a Multilanguage character recognition system for Arabic script based languages. The main benefit of proposed approach is that it works well for all Arabic script based languages by treating the ghost shapes and diacritical marks separately and using dictionary for each language for ligature and word formation. On the other hand, treating all Arabic script based languages has several issues as compared to a system for a specific language (i.e. Urdu or Persian) and specific writing style (i.e. Nasta'liq or Naskh). The layered word formation can be followed by the valid ligatures recognition and these valid ligatures are combined to form a valid word. Moreover,

Ghost Character Recognition Theory also helps to mimic the biological character recognition system.

4.1. Urdu Script Based Languages

There are at least 26% Muslims in the world having directly or indirectly interaction with Arabic script based languages, it may be due to the fact that Islam originated from the Arabia and Arabic is the language of the Quran. Generally this script is followed in many countries of Arabian Peninsula, Iraq, Iran, Pakistan, Afghanistan, India, Uzbekistan, Tajikistan, Kazakhstan, Turkey, etc. and is followed by many languages like Arabic, Persian, Urdu, Punjabi, Sindhi, Pashto, Blochi, etc. Arabic script based languages especially Urdu and Arabic are used in almost every part of the world.

Arabic script based languages are written in a cursive style from right to left in both machines printed and handwritten forms. These are context sensitive languages and are written in the form of ligatures that may comprise a single or many different characters to form a word. Most of the characters have different shapes depending on their position in the ligature e.g. a letter may appear differently depending on its position as an isolated, middle, center, or ending character as shown in Figure 4.1. Arabic script based languages also use punctuation marks to separate sentences and leave white space between ligatures and words for separation. Furthermore, characters may overlap with each other and are very rich in diacritical marks; i.e. Urdu contains 22 diacritical marks and these additional diacritical marks associated with ligature represent short vowels or other sounds. Some diacritical marks are compulsory whereas some diacritical marks are optional and only added to help in pronunciation. Optional diacritical marks are not often used by the native speaker i.e. Arabic and

Urdu speaker who do not use the optional diacritical marks which are only added for the nonnative speaker.

ب	ب	ب	ب
ع	ع	ع	ع

Figure 4.1: Different Shapes of (ب and ع) with respect to position from left to right
isolated, start, mid, end

ا	ب	ت	ث	ج	ح	خ
alif	baa	taa	thaa	jiim	haa	kha
د	ذ	ر	ز	س	ش	ص
daal	thaal	raa	zaay	siin	shiin	saad
ض	ط	ظ	ع	غ	ف	ق
daad	taa	thaa	syn	ghayn	faa	qaaf
ك	ل	م	ن	ه	و	ي
kaaf	laam	miim	nuun	ha	waaw	yaa

ا	ب	پ	ت	ٹ	ج	چ	ح	خ	د	ذ
alif	ba	pe	ta	se	jeen	che	ha	xa	da	zai
ر	ز	ژ	س	ش	ص	ض	ط	ظ	ع	غ
re	ze	ze	sin	shin	sad	sad	ta	za	ayn	ghayn
ف	ق	ک	گ	ل	م	ن	و	ه	ي	
fe	qa	kar	gar	lam	mim	nun	vaav	he	ya	

Figure4.2. a) Arabic Alphabets

b) Persian Alphabets

Arabic is spoken in many countries i.e. Saudi Arab, UAE, Oman, Jordan, Kuwait, Iraq, Egypt, Syria, Yemen, etc. Arabic is also the language of the Quran, divine commandments revealed to the last prophet Muhammad [PBUH] and later compiled in the form of a book. Thus this script is mostly used by Muslims either directly (Arabic) or indirectly (in the form of other language like Urdu, Persian etc.). Arabic is ranked at 5th most spoken language in the world and is written in Naskh style. It consists of 28 alphabets shown in Figure 4.2.a. Historically, it was written without diacritical marks, but later on diacritical marks were added for non-natives readers. Arabic has great influence on many other languages especially those used in

Muslim countries and is a major source of vocabulary for many languages i.e. Spanish, Persian, Urdu, Hindi, Punjabi, Sindhi, Pashto, Malay, Turkish, Gujarati, Kurdish, Bengali, Kashmiri, etc.

Persian also known as Farsi and is the official language of Iran, Tajikistan and Afghanistan; it is written in Arabic script (Nasta'liq style) and has 32 basic alphabets as shown in Figure 4.2.b. It also has a large influence on Urdu, Punjabi and Sindhi and other south Asian language [Lazard, 1975].

Urdu is the 2nd most spoken language of the world and is written in two main scripts; Arabic Script, and Devanagari script. When written in Arabic script, it is called Urdu which is mostly used in Pakistan and when Devanagari script is followed then it is called Hindi which is used in India. The language scholars categorize Urdu as standard version of Hindi. Actually Urdu has different versions that depend upon regions instead of the writing script [Durani, 2008]. Urdu is the national language of Pakistan and official language of many Indian states. It is written in Arabic script (Nasta'liq style) and consists of 58 basic letters as shown in Figure 4.3.a. Other languages that share the Urdu script are Sindhi, Pashto, Punjabi and Balochi. Punjabi is the local language of Pakistan and India. It is written in Gurmukhi and Shahmuki in Indian and Pakistani Punjab respectively. Shahmukhi is based on Arabic script and mostly written in Nasta'liq style as shown in figure 3.b. Punjabi consists of 47 alphabets and is ranked 11th most spoken language in the world.

Sindhi is the local language of India and Pakistan written in both Arabic and Devanagari script. It is official language of the province of Sindhi, and also some states of India. In Pakistan; Sindhi is written in Arabic script and contains 52 alphabets as shown in Figure 4.a. and is ranked as 23rd most spoken language. Pashto

is written in Arabic script (Naskh) and is spoken in parts of Afghanistan and is a local language of Pakistan. It is influenced by Farsi and Avastan however most of the words belong to itself. It consists of 39 alphabets as shown in figure 4.b. and it is ranked at 33.

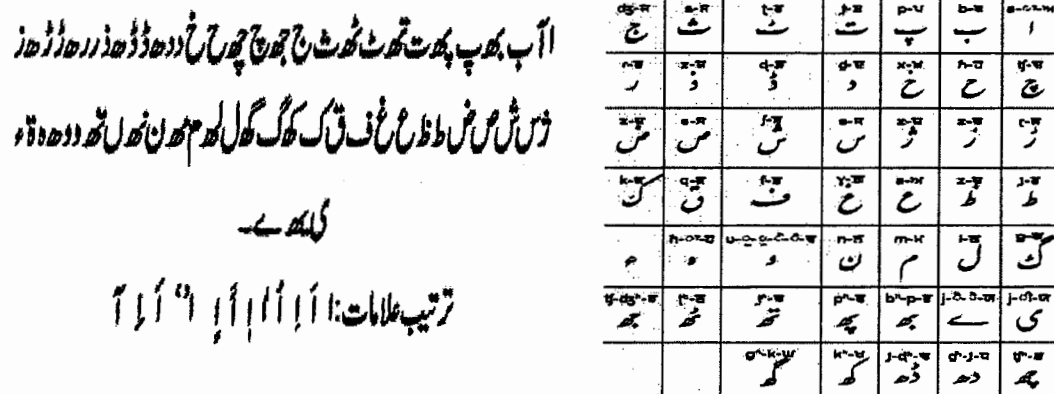


Figure 4.3. a) Urdu Alphabets [Durani 2008] (b) Punjabi Alphabets (Shahmukhi)

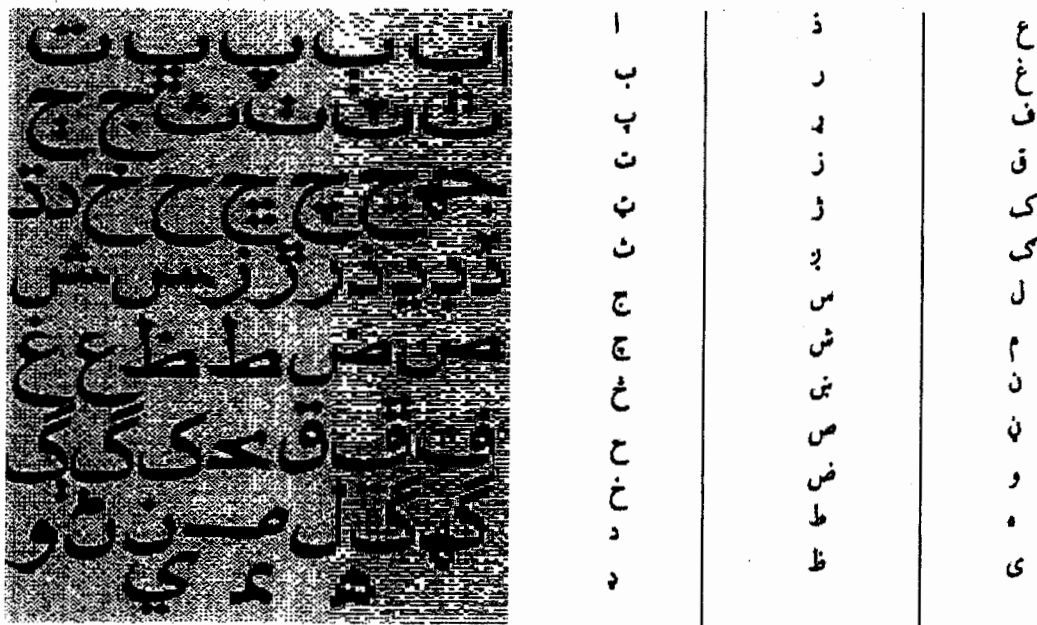


Figure 4.4. a) Sindhi Alphabets (b) Pashto Alphabets

From the above discussion, it is concluded that Urdu is the superset of all Arabic script based languages as it contains all basic shapes that are used by other

Arabic script based languages. Local languages of Pakistan like Punjabi, Sindhi, and Pashto have different letters as compared to Urdu but with the same basic shapes and different diacritical marks. Moreover, it is also concluded that all Arabic script based languages have same basic structure and Urdu is complex in the family of Arabic script based languages.

4.2 Character Recognition of Arabic Script Based Languages

Character recognition is the branch of pattern recognition to imitate the computer in reading the graphical marks written by human or printed by machine so that that the machine can perform like human skills in reading. It has been an ongoing research problem for more than four decades. Basically character recognition is classified into two classes with respect to input namely online (handwritten), offline handwritten/printed recognition. In offline; input is in the form of image while in online case coordinates as well as timing information. The offline printed character recognition is little easy task as compared to handwritten either online or offline due to large variation in writing styles. The recognition for Arabic script based languages is much more complicated than any other language like English, Chinese etc due to complexities of this script. These complexities are context sensitive shape, Cursiveness, Overlapping, large number of diacritical marks, segmentation of words itself and mapping of diacritical marks. Moreover the recognition for handwritten Nasta'liq is much more complicated as compared to Naskh writing style due to its complexity. Naskh may contain only 4 basic shapes of character whereas Nasta'liq contains 32 shapes depending upon the associated character shown in figure 4.9.

Limited research efforts have been carried out in Arabic script based languages character recognition especially for handwritten recognition even though

no Multilanguage character recognition system available while there exist high similarity between Arabic script based languages. Both segmentation base [Safabakhsh 2006, Haraty 2002, 2004, Sari 2002, Fahmy 2000, Miled 2001, Abuhaiba 1998] and holistic [Haji 2005, Meslati and Farha 2004, Adeed 2004, 2002, Khorsheed 2003, Pechwitz 2003, Dehgan 2001 and Al-Badr and Haralick 1998,1995] approaches are discussed for Arabic script based languages (both printed and handwritten) by using diacritical marks as features points with other features. There is no such (separate the diacritical marks from ghost character and map these diacritical marks with respect to position after recognition separately) effort proposed in the literature that leads to multilingual character recognizer.

4.3 Ghost Character Theory

"There are some problems in Urdu ASCII code plate, when I analyzed that some symbols and all the language of Pakistan is possible from one code plate and one font. Then I proposed the idea of Ghost Character. [Durani 2008]".

Nasta'liq and Naskh are two basic and different writing scripts that have their own fonts. Urdu is not a subset of Arabic [Durani, 2008] moreover Urdu alphabets are the super set of alphabets of all Arabic script based languages. Nasta'liq is more complicated than Naskh, due to a higher number of shapes for each character in Nasta'liq as compared to Naskh where each character may have 4 shapes i.e. "Bay" has 32 shapes in Nasta'liq while only 4 in Naskh.[Durani, 2008].

All Arabic script based language can be written with only 44 ghost characters. Ghost character consists of 22 basic shapes called Kashti as shown in Figure 4.6 and 22 diacritical marks [Durani, 2009]. Basically, the idea is 700 years old when diacritical

marks were first applied on Quran to make it easy to read for the non-native readers by Hajaj Bin Yousif. Before his period, there were no dots and diacritical marks used in Arabic. Arabs were using only 19 characters, and they were reading this dot-less character by their cultural habits and had no difficulty in reading. The philosophy behind dots was that the first character with similar shapes had one dot, 2nd character had 2 dots and 3rd had 3 dots. Persians also followed the Arabic script after Islam in Persia and some dots on character were added that were not in Arabic. Similarly in Urdu 4 nuqtas were added on ghost character which were later on converted to line and then to Urdu letter "Tota" as shown in Figure 4.5.a. Some of the basic shapes were added in Urdu and Persian as shown in Figure 4.5.b [Durani 2008].



Figure 4.5: a. Convergence of four dots to "Tota" (b). Additional shapes in Urdu and Persian.

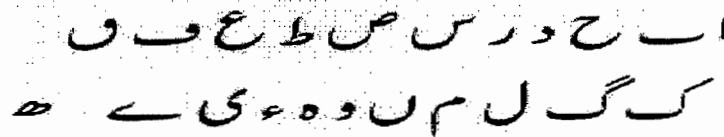


Figure 4.6. Ghost characters used in Arabic script based language [Durani 2008].

The 22 ghost characters used in Arabic script based languages are shown in figure 4.6. All the Arabic script based languages like Persian, Urdu, Punjabi, Sindhi, Persian, Balti, etc. can be written with these 22 ghost characters, 22 dots and diacritical marks. Thus, these 44 basic shapes and diacritical marks can represent all Arabic script based languages.

4.4 Ghost Character Recognition Theory

Character recognition of Arabic script based languages is very difficult task due to complication involved in this script and it has a large variety of shapes. Only Urdu script has more than 22000 ligatures [Durani, 2009]. No research efforts have been done in the field of Multilanguage character recognition even though there are minor differences between scripts followed by these languages. Most of the work done is language specific while Multilanguage system can easily be achieved by making some more effort in the preprocessing and post processing phases. To overcome language specific character recognition with Multilanguage character recognition for Arabic script, ghost character recognition theory is presented.

As mentioned above, all the Arabic script based languages can be written with the 22 ghost character and 22 diacritical marks but each base ligature has its own phonemes and meanings in each language with the same or different number of diacritical marks. Thus the basic shapes (glyph) are same for all Arabic script based languages with only difference in font i.e. Naskh, Nasta'liq and diacritical marks. Nasta'liq is mainly followed by Urdu, Persian, Sindhi and Punjabi and is more complicated than Naskh i.e. "Bey" has 32 shapes as shown in Figure 4.8. Ghost character theory has great influence on Arabic script based languages character recognition in order to develop Multilanguage character recognition.

Based on the ghost character theory the ghost character recognition theory is divided into five basic steps are

- I. First step is to segment the additional marks i.e. dots, diacritical marks from the word. Now the word consists of only ghost characters (khali kashti) and diacritical marks.

- II. Recognize the separated basic shape through a classifier.
- III. Recognize the diacritical marks associated with recognized ligature.
- IV. Map the diacritical marks on to the recognized ghost character using the language specific ligature formation dictionary.
- V. Word formation from valid ligature using language specific dictionary.

The above process is shown in Figure 4.7 for 2nd ghost character of Figure 4.6 used in all Arabic script based languages.

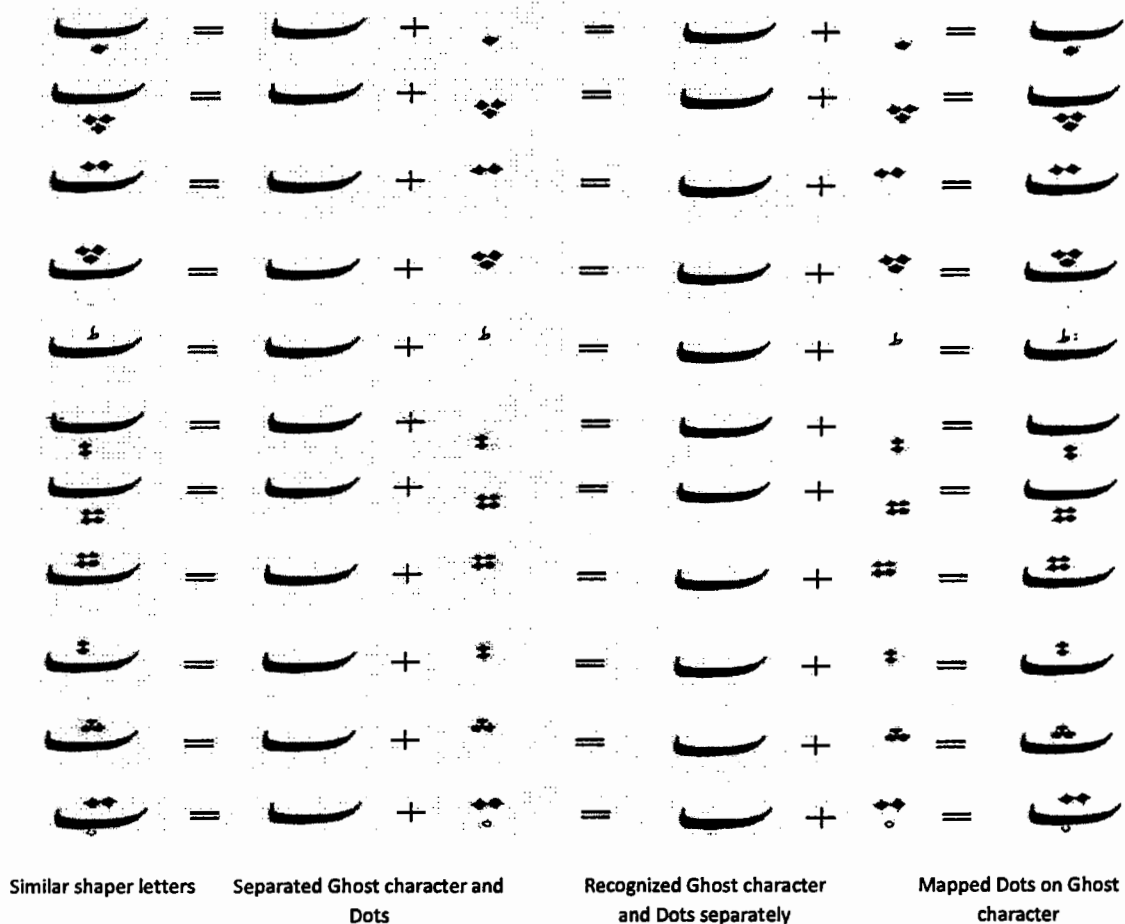


Figure 4.7: Recognition of 2nd ghost character letter with associated dot

Figure 4.7 shows the character recognition process for 2nd Arabic Character based on Ghost character theory. Total number of character for all Arabic script based

languages are 11 having same basic shape and different diacritical marks i.e. Urdu contains 5, Arabic contains 3. Thus it reduces the data set from 11 to 1 character. The association of diacritical marks may be based on the language rules or language dictionary.

As it is a very difficult task to classify Arabic script based languages due to complexities involved in the script, especially for handwritten text. The training for every language puts a big overhead on recognition engine to classify different writing styles like Nasta'liq, Naskh by a single classifier. This increases the complexity and reduces the recognition rate. This issue is resolved by implementing the ghost character theory and extracting the style independent structural features like loop, cusp, end points, line shapes, etc. In other words this is done by developing two separate systems for mostly used writing styles Naskh and Nasta'liq. Nasta'liq style is more complex than other styles followed by Arabic script based languages shown in Figure 4.8 and Figure 4.9. The character appears in Nasta'liq style may also appear in Naskh style with a little variation. Thus, the system developed for Nasta'liq by using structural features may also work for other writing styles.

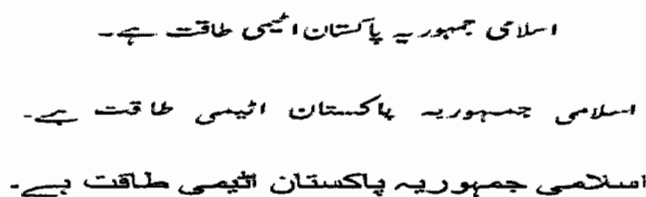


Figure 4.8: Urdu Samples in three different styles. Urdu Nasta'liq, Urdu Nasq, Naskh

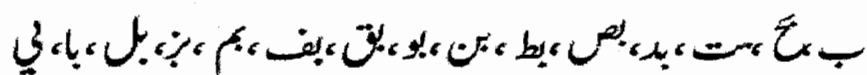


Figure 4.9: different shapes of "ب" in Nasta'liq Font with respect to neighbor character

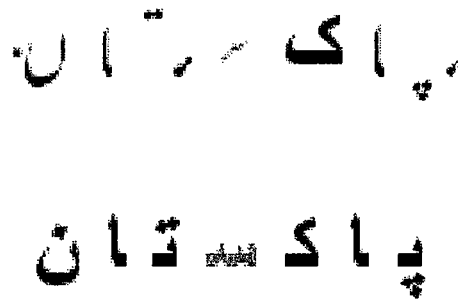


Figure 4.10: Feature comparison of Nasta'liq and Naskh

Basically the structural features i.e. loops, cusp, endpoint, etc. are intuitive aspects of writing and computed from the skeleton of the ligature. Furthermore the extraction and mapping of diacritical method is also based on the structural features especially for Arabic script based languages which are healthy in diacritical marks. Due to this reason structural features are mostly used for Arabic script based languages in literature. By deeply analyzing the both Nasta'liq and Naskh, it is concluded that structural features for Urdu script written in Nasta'liq font may also work for other script written in either Nasta'liq or Naskh style. This is due to the complexities in the Nasta'liq script. The shapes in Nasta'liq are more complex and vary up to 32 with respect to its associated character and position whereas in Naskh shapes are only four deepening upon the position of the character. Thus those structural features along with directional features are extracted that can help in both writing style i.e. loop, cusp, endpoint, branch point, etc.

4.5 Effect of Ghost Character Theory

Urdu script based language character recognition system for both Nasta'liq and Naskh writing style. Normally Naskh and Nasta'liq are mostly followed by Arabic script based languages. Nasta'liq is mostly followed for Urdu, Punjabi, Sindhi, etc. whereas Naskh is mostly followed for Arabic, Persian, etc. Thus, this work is selected because

Algorithm: Diacritical Marks Fuzzy C-Mean Clustering

Input: Online handwritten strokes (base and diacritical marks)

Output: Diacritical marks with associated character i.e. position

Algorithm:

Step I: Select the number of cluster based on estimated characters.

Assign every cluster to every expected character.

Step II: Calculate the center of fuzzy cluster using position.

Step III: Repeat step IV and V for each diacritical mark

Step IV: Phase-1: Perform fuzzy estimation of diacritical marks clusters with each other.

Phase-II: Perform fuzzy estimation of diacritical cluster with characters.

Step V: Draw fuzzy projection on to the character to find the associated character.

5.1.5. Slant Correction

The character inclination that is normally found in cursive writing is called slant. Slant correction is one of the important steps to reduce the variation in the handwritten text and crucial part of preprocessing. The main purpose of slant correction is to reduce the variation in handwritten script specifically to improve the diacritical marks segmentation process, etc. which in turn can yield higher recognition accuracy. Since for Urdu script especially Nasta'liq the feature extraction and segmentation is highly dependent on the direction of writing thus it is important that vertical movement of pen is as perpendicular as it is possible, instead of diagonally slanted. Otherwise for diagonally slanted strokes written features will be poorly extracted or segmentation cannot be

performed correctly. For example for implicit segmentation it is very difficult to adjust the criteria for segmenting the stroke into subunits i.e. the input stroke is divided into slice/subunits for HMM based recognition so that accurate splitting is obtained; the strokes must be vertically perpendicular.

Generally handwritten strokes are not uniform, thus the slant may change with words and it may even be different within ligatures or strokes. Therefore the global uniform slant estimation is not guaranteed to provide good slant correction for handwritten text. Due to the complexity of Urdu handwritten script we need local slant correction instead of global slant correction. Whereas it's a difficult to estimate the slant on stroke bases due to limited information available for estimation. This problem mostly occurs where short strokes in height and width are involved i.e. ر، د، و. These kind of short strokes contain limited information that helps in slant estimation. Thus the need is to normalize the stroke by estimating the slant locally based on the current stroke and some clue from the previous stroke to normalize the stroke on to the vertical axis. We present an algorithm that computes slant locally based on vertical projection with the help of neighbor strokes information.

For slant estimation only vertical and horizontal movements of pen are analyzed while slant correction is performed on only vertical axis. As the Urdu ligatures are written from top to bottom and left to right. The writers may write the ligatures little right or left slanted shown in Figure 5.12. The slant estimation is performed on only primary strokes because some of the secondary strokes are slanted naturally whereas alignment correction is performed on both primary and secondary strokes as shown in figure 6. For slant estimation, gravity fall of topmost strokes elements is used with additionally

most popular styles for these languages i.e. Naskh is used for Arabic while Nasta'liq is used for Urdu, Punjabi and Persian. The overall ligatures hence decrease.

$\text{Ligature}_{\text{Multilanguage}} = \text{No of total ligatures by Arabic script based languages}$

$\text{Ligature}_{\text{Arabic}} = \text{No of total ligatures of Arabic}$

$\text{Ligature}_{\text{Urdu}} = \text{No of total ligatures of Urdu}$

$\text{Ligature}_{\text{Persian}} = \text{No of total ligatures of Persian}$

$\text{Ligature}_{\text{Punjabi}} = \text{No of total ligatures of Punjabi}$

$\text{Ligature}_{\text{other Arabic script based languages}} = \text{No of total ligatures of other Arabic script based languages like Pashto, Sindhi, etc.}$

$\text{Ligature}_{\text{Multilanguage}} \lll \text{Ligature}_{\text{Arabic}} + \text{Ligature}_{\text{Urdu}}$

$+ \text{Ligature}_{\text{Persian}} + \text{Ligature}_{\text{Punjabi}}$ (4.1)

$+ \text{Ligature}_{\text{other Arabic script based languages}}$

Moreover by using the character theory, it can also be written as

$\text{Ghost Ligature}_{\text{Urdu}} \ll \text{Ligature}_{\text{Urdu}}$ (4.2)

4.7 Limitations

The ghost character recognition theory has also few disadvantages.

- I. As there are multiple languages in one classifier, the number of ligatures is increased.
- II. It is a difficult and complex task to develop classifier for multiple fonts for Arabic script based languages.

- III. The recognition rate is little low due to multiple fonts and large number of diacritical marks as compared to language specific character recognition system.

4.8. Final Remarks

Every fourth person in the world is Muslim and Arabic script is used directly or indirectly by Muslims. Urdu is the language that contains 58 alphabets; the basic shapes in Urdu also exist in other languages. Thus Urdu is the superset of all other Arabic script based languages. This chapter presents a novel technique; ghost character recognition theory that helps to develop Multilanguage character recognition system for all Arabic script based languages. The main advantage of the proposed technique is that the recognition system works for all Arabic script based languages by classifying ghost character and mapping the associated diacritical marks and dots latter with respect to selected language. By deeply analyzing the shapes, appearance and structure of both Nasta'liq and Naskh script with Urdu linguistics it is concluded that structural features for Urdu script written in Nasta'liq font may also work for other script written in either Nasta'liq or Naskh style. This chapter discussion shows that, an efficient multilingual character recognition for Arabic script based languages system can be developed by including some preprocessing and post processing steps. Moreover, Ghost Character Recognition Theory also helps to build HVS inspired character recognition system.

CHAPTER 5

PREPROCESSING AND FEATURE EXTRACTION

This chapter discusses the issues involved in preprocessing and feature extraction. Preprocessing forms level-1 and level-2 whereas the feature extraction comprises level-3 and level-4 of the bio-inspired character recognition system. The complete bio-inspired character recognition is presented in chapter 6. The level-1 is responsible for the segmentation of diacritical marks and basic level preprocessing; level-2 is responsible for complex preprocessing steps. Level-3 is responsible for dividing the strokes into sub patterns and extraction of unique sub patterns whereas the level 4 is responsible for complex feature extraction and fusion of features. Pre-processing of the raw input strokes is crucial part for better feature extraction and for success of character recognition system. First phase (level-1 and level-2) describes several preprocessing steps for online character recognition by considering the input strokes from both online and offline perspectives to reduce the variation by normalizing the input stroke. The proposed technique is also a necessary step towards character recognition, person identification,

personality determination where input data is processed from all perspectives. The second phase (level-3 and level-4) presents the feature extraction method used for classification. Feature extraction in pattern recognition problems involves the extraction of unique and salient patterns from the raw data in order to enhance the discriminatory power and reduce the data for classification. Success of classifier depends upon the feature extraction. Fuzzy based reasoning is presented for both preprocessing and feature extraction.

5.1 Preprocessing

Preprocessing is one of the most important phases of character recognition and directly influences the recognition result. It is required in order to compensate for variations in the stroke elements. Several preprocessing steps are performed and the preprocessing phase has been dealt with both offline and online data.

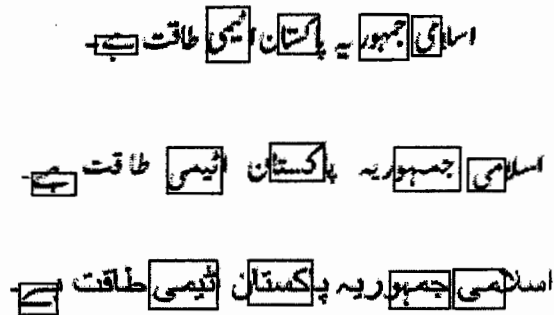


Figure 5.1: Urdu Samples in three different styles. Urdu Nasta'liq, Urdu Nasq, Naskh

The preprocessing in both domains includes compensating for variability in size, smoothing from irregular patterns, skew detection and base line estimation, etc. The preprocessing prepares the text for segmentation and feature extraction. Unlike English,

Japanese, etc. the cursive Arabic script has not received much attention in the last few years [Bouchareb et.al, 2006] and within the family of Arabic script; studies dealing with Urdu Nasta'liq script are scarce.

The main aim of the preprocessing step is to produce a clean version of raw input data for either online strokes or images of handwritten text, so that they can be utilized efficiently for feature extraction and classification. Generally, preprocessing phase is a set of filters that produce simplified patterns from the noisy patterns for classification. The raw input strokes for online character recognition consist of stroke elements i.e. x, y coordinates, timing information, force. The writing force and speed is important for personality detection, person identification, etc. but it is not important for character recognition due to a variety of writing speeds in writing of the same strokes by the same person as well as different people. Thus the accuracy of recognizing text either by humans or machine is highly dependent on the quality of input as even humans cannot read easily the noisy, irregular or unfamiliar input data.

To minimize the unnecessary elements for classification, several preprocessing steps are presented for online Urdu character recognition. These unnecessary elements occur due to a vast variety of writing styles and inconsistencies that exist in the natural way of writing. It is not an easy task to reduce noise and variations to simplify the input pattern for Urdu character recognition due to complexities in the script.

Generally preprocessing phase consists of some operation such as baseline detection, de-hooking, smoothing, secondary strokes separation, slant correction, stroke mapping. The resultant outcome of pre-processing phase is a normalized, skewed

smoothed and clean sequence of stroke element x', y' and can be used directly for feature extraction and classification.

5.1.1. De-Hooking

Hooks are most common artifacts that occur at the beginning and ending of strokes during handwriting. Due to the sensitivity of pen/pc tablet, hooks are generated by inexperienced user or due to fast writing speed during the pen-down and pen-up movements as shown in Figure 5.2. In other words hooks are the noisy pattern that may occur at the start and end of handwritten strokes and are common in handwriting. The presence of hooks may lead to incorrect feature extraction especially hooks create problems in recognizing the character start with jeem 'ج' and ayen 'ع' due to the natural presence of hooks in these characters. Thus before removing the hooks it is necessary to take the decision whether it is a hook or not.

If variation in the chain code (last 6 chain code in length) at the beginning or end is less than the specified threshold ($T = 4$), then the part considered a hook is removed by either discarding it or replacing the respective co-ordinates with the neighboring ones. There are some issue in removing the hooks i.e. it may be possible that small up of the stroke 'ج' is removed as a hook which is the most important part in the detection of stroke 'ج' as shown in Figure 5.3. So, to avoid de-hooking in jeem, de-hooking at beginning is not performed on those ligatures which are written from left to right for some length like ج. The isolated jeem is written from left to right and the ligature جر is also written left to right at the beginning while the remaining ligature is written from right to left.

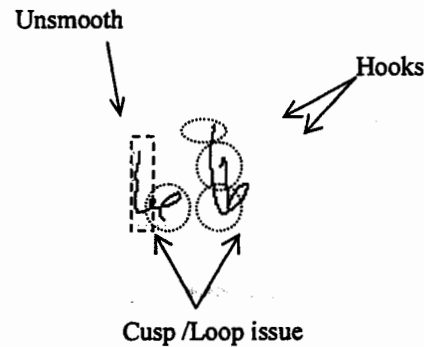


Figure 5.2: Before Preprocessing issue (a) (b)

5.1.2 Smoothing

The main aim of smoothing is to decrease the variation occurred due to fast writing, or natural hand shivering during writing and it is common in literature to normalize the input text. Smoothing is the most significant operation of preprocessing and it attempts to normalize the variation that does not affect the basic structure of the stroke. In the case of offline character recognition, averaging filter is applied on the image whereas for online character recognition the smoothing filter is only applied to stroke element array, thus the smoothing operation does not affect the width of the stroke elements x , y . Smoothing is performed using weighted average over strokes elements as shown in Figure 5.4 using filters shown below where x_i and y_i are the corresponding x , y coordinates.

We have used 5x5 windows to select the weight. The weight is 3 at vertical and horizontal and weight eight is 2 at exactly diagonal whereas to move the less diagonal line towards horizontal or vertical straight line the weight is 3.

$$P'(x_i, y_i) = \frac{\sum_{j=\max(i-k,1)}^{j=\min(i+k,n)} P_{i+j}(x, y) \times w_j}{\sum_{j=-k}^k w_j} \quad (5.1)$$

5.1.3 Interpolation

Sometimes, intermediate points are skipped by pen due to the high writing speed and low processing power of pen. Moreover, the points are not equidistant between them but they are equidistant in time. Thus the number of points varies depending upon the writing speed. This missing and non-equidistant data may create some problems in feature extraction, especially in the case of loop and cusp detection and it's a crucial step of preprocessing phase. Figure 5.3 shows the missing points. Mezghani et.al presented an approach for re-sampling the points at equal distance [Mezghani et.al, 2008]. The intermediate points are estimated after fixed distance whereas in loop extraction some problem may occur when the point of intersection is not the same. Thus proposed approach interpolates the missing points. To compute the missing data between two points interpolation is performed using Bresenham's line drawing algorithms shown in figure 5.4. Moreover, filtering step is also necessary to eliminate the duplicated points by forcing a minimum distance between consecutive points. The duplicates points are removed using low pass filter on raw input data.



Figure 5.3: Missing points due to fast speed of writing.



Figure 5.4: Smoothed and Interpolated stroke

5.1.4. Secondary Stroke Segmentation

Delayed strokes (secondary strokes) are necessary to differentiate the similar strokes in inter languages i.e. Arabic, Urdu and Persian, etc. and also in intra language i.e.

The same ghost stroke may have different interpretations based on the diacritical marks. Arabic script based languages contain zero or more secondary strokes associated with primary strokes. A study on 50 native readers of Arabic, Urdu, Punjabi, Sindhi and Pashto written in two mostly followed fonts Naskh and Nasta'liq showed that the first step in recognition of Arabic script based languages is the recognition of basic shapes without diacritical marks. Thus secondary strokes handling is very important for the classification of Arabic script based languages. Generally Arabic script based languages contains zero or more secondary strokes corresponding to one primary stroke and smaller in size as compared to primary strokes. The separation of secondary strokes is not an easy task, because the dots may not appear exactly above or below the character in the stroke and they may have different order furthermore some delayed strokes have same or similar shapes as of primary strokes i.e. " " may appear as primary strokes and secondary stroke as well. The positioning of secondary stroke is also very critical especially when related character is of very small size or related stroke may contain more than one secondary

stroke shown in Figure 5.5. It is very difficult to decide, whether the secondary stroke belongs to one character in the ligature or belongs to multiple characters in the ligatures.

If Premises P_0 then Conclusion C_0

If Premises P_1 then Conclusion C_1

(5.2)

Where Premises $P_1 \in C_0$

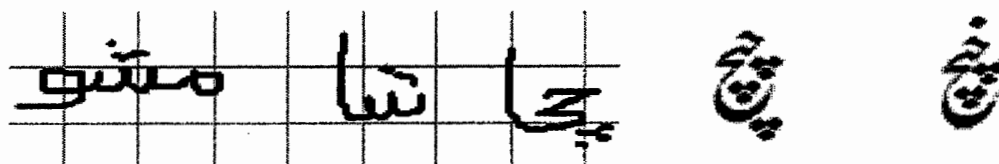


Figure 5.5 (a): Nasta'liq character with three dots

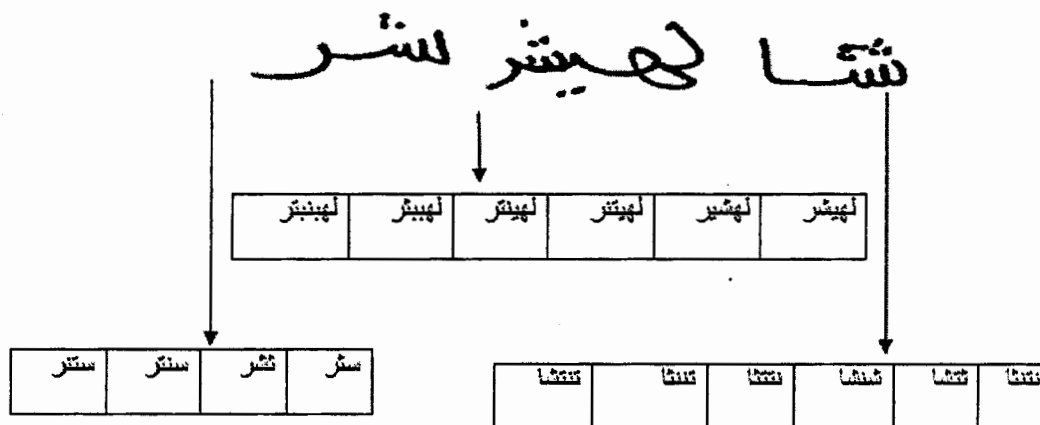


Figure 5.6: Diacritical marks issue with three dots in Naskh style

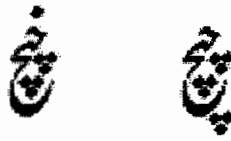


Figure 5.7: Diacritical complexity

The separation of secondary strokes is not easy task, because the dots may not appear exactly above or below the character in the stroke and they may with different order furthermore some delayed strokes have same or similar shapes as of primary strokes i.e. " " may appear as primary strokes and secondary stroke as well. There are two main issues in diacritical marks segmentation.

- I. Segmentation
- II. Localization

5.1.4.1. Segmentation

As in Urdu script; secondary strokes are written above or below of the word and in some cases they may appear little before, little after (closely touching the right or left side respectively), or within the word-part with respect to the horizontal axis as shown in figure 5.7. The strokes smaller than the size θ (where θ is $1/10$ of base stroke) are considered as candidate for secondary strokes. If the size is greater than θ then location of the occurring stroke is compared with the primary stroke, it will be considered a candidate for secondary strokes if it lies little up, ending must be from right to left and thirdly it's ending point lies within range θ_1 (where θ_1 $1/6$ of stroke width). By combining the fuzzy terms through the logical operators (OR, AND) on fuzzy variables i.e. location,

timing, stroke size, etc, a number of fuzzy rules are constructed for delayed stroke segmentation. The fuzzy membership function is explained in table 5.1 and table 5.2.

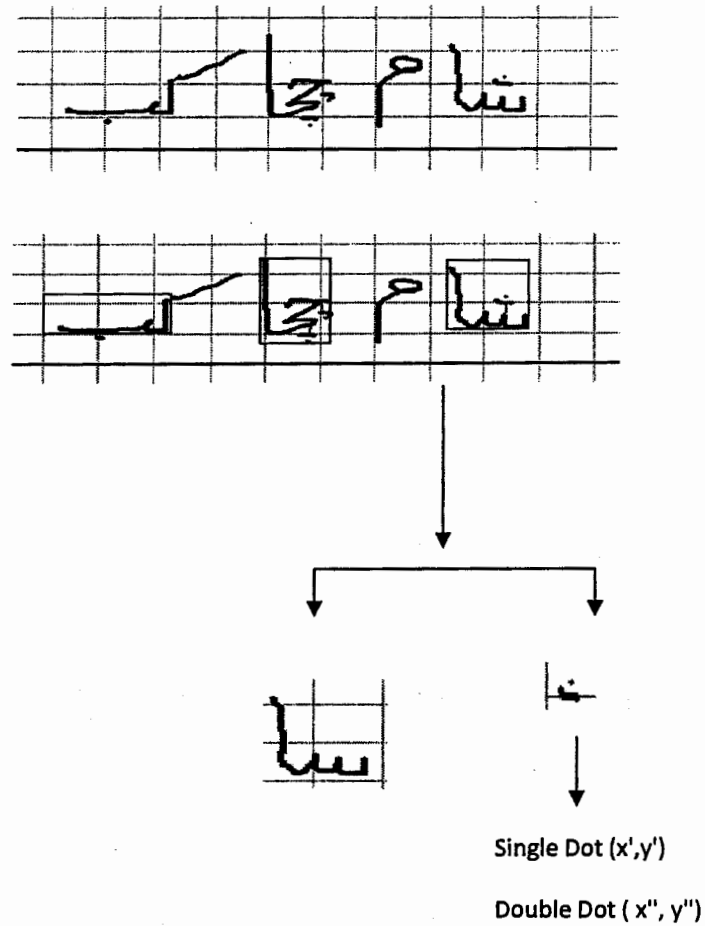


Figure 5.8: Secondary stroke separation

Table 5.1: Distance Vs Position Fuzzy Rules

Distance/Position	Up	Down	Inside	Left	Right
Small	MS-1	MS-1	MS-1	MS-1	MS-1
Medium	MS-1	MS-1	MS-1	MS-2	MS-2
Large	MS-2	MS-2	MS-2	NS	MS-3
Very Large	MS-3	NS	NS	NS	NS

Inherited result

Table 5.2: Table 1 Result Vs Size Fuzzy Rules

Size/Score Result	MS-1	MS-2	MS-3
Small	SS	SS	NS
Medium	SS	SS	NS
Large	SS	NS	NS
Very Large	SS	NS	NS

*MS May be Secondary Stroke

**NS not Secondary Stroke

Table 5.3: Diacritical Marks Separation Algorithm

Algorithms: Delayed Stroke Segmentation
<p>For All P, P_{ij}</p> <p>P_{ij} is the delayed Stroke, P_i is the concerned primary stroke.</p> <p>$RMER(P_i)$ is the point on right side of minimum enclosing rectangle</p> <p>$LMER(P_i)$ is the point on left side of minimum enclosing rectangle</p> <p>IF $P_{i+j} < \theta$</p> <p style="padding-left: 40px;">P_{i+j} may be the candidate for secondary stroke $C(S_j P_i)$</p> <p>Else</p> <p>IF $P_{ij} (x_f, y_f) - RMER(P_i) < \theta_1$ OR $P_{ij} (x_f, y_f) - LMER(P_i) < \theta_1$</p> <p style="padding-left: 40px;">P_{ij} may be the candidate for secondary stroke $C(S_j P_i)$</p> <p>Otherwise</p> <p>$P_i = P_{i+1}$</p>

5.1.4.2 Localization

The positioning of secondary stroke is also very critical especially when related character is of very small size or related stroke may contain more than one secondary stroke shown in figure 5.9. It is very difficult to decide, either the secondary stroke belong to one character in the ligature or belong to multiple character in the ligatures. Fuzzy logic is used to handle delayed-stroke through the vertically projection on the stroke to find its corresponding character surrounded by considering vertical projection along with timing information and corresponding stroke shape. The propose technique is

divided into two parts fuzzy clustering and estimation of associated character. In the first phase the local cluster are estimated on diacritical family with respective to relative position and associated characters. The basic aim of clustering is to resolve the conflict of dots such as the one shown in Figure 5.9. Unlike the direct projection or direction estimation of dots family, the fuzzy c-mean clustering finds the degree of belonging to each family (cluster). Thus by using the fuzzy c-mean clustering first optimized the family of the each diacritical mark and then finds the projection on to the character. Relative position of diacritical marks is very important in order to find the associated cluster. By considering the biological knowledge of human visual perception, the diacritical position and associated character are estimated. The relative size and position of diacritics and character is most important to map the diacritics onto the associated character. The fuzzy c-mean algorithm has been successfully applied to many clustering problems. It is technique where each point belongs to cluster to some degree. Fuzzy c-mean is a simple unsupervised learning technique that can be used for classification or data grouping when the clusters size is known.

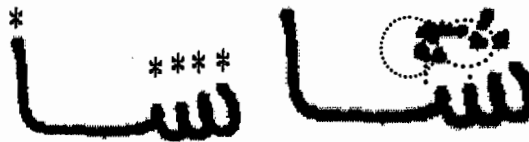


Figure 5.9: Diacritical marks issue with three dots in Naskh

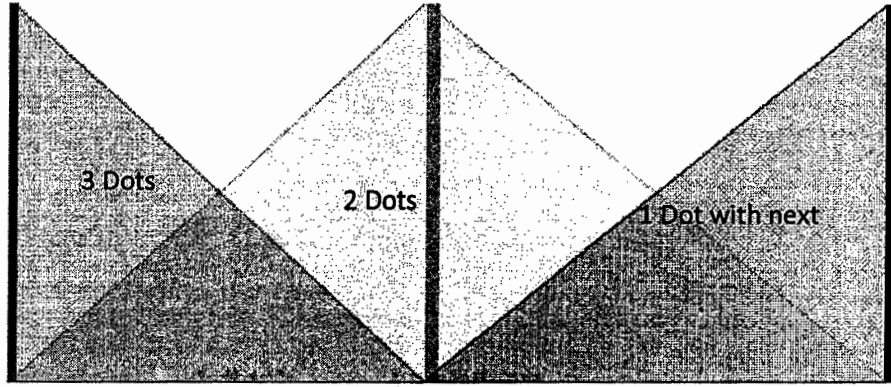


Figure 5.10: Fuzzy membership function for diacritical marks

For supervised learning distance metric, we used some training samples obtained from 10 users which contain single dot, double dot and triple dot, taa, etc.

$$X = \{x_1, x_2, \dots, x_n\} \text{ and } C = \{c_1, c_2, \dots, c_m\} \quad (5.4)$$

Where n is the number of diacritical marks associated with m base characters. x_i may belongs to the current character, previous or next character or dot of same, previous or next character.

$$S = \{(x_i, x_j) | x_j \text{ and } x_j \in x_i\} \quad (5.5)$$

$$S = \{(c_i, x_j) | x_j \text{ and } x_j \notin x_i\} \quad (5.6)$$

$$C = \{(c_i, s_j) | s_j \text{ and } s_j \in c_i\} \quad (5.7)$$

$$C = \{(c_i, s_j) | s_j \text{ and } s_j \notin c_i\} \quad (5.8)$$

The decision is based upon the local clustering of diacritical marks and clustering with the character. In the first phase, the above equations performed diacritical marks clustering within the diacritical family. The second phase run on the clustered diacritical

marks and character in the ligature to find the associated character has shown in equation 5.5 and 5.7. We also discarded many pair in first phase based on the projection on to the character. The distance can be computed as

$$d(x_i, c_j) = \|x_i - c_j\| = \sqrt{f_{(i,j)}^T A f_{(i,j)}} \quad (5.9)$$

Where $f_{i,j}$ is the feature vector

One of the main issues is handling the diacritical marks associated with each character. We present a technique to handle diacritical marks with respect to relative position and using fuzzy c-means clustering. Fuzzy logic is robust enough for scatter and irregular data. We elevated the performance of the character recognition system and experiments show that the proposed technique performs well to deal with diacritical marks moreover the result can also be improved by using dictionary during the mapping of diacritical on to the base stroke to avoid the incorrect clustering of dots.

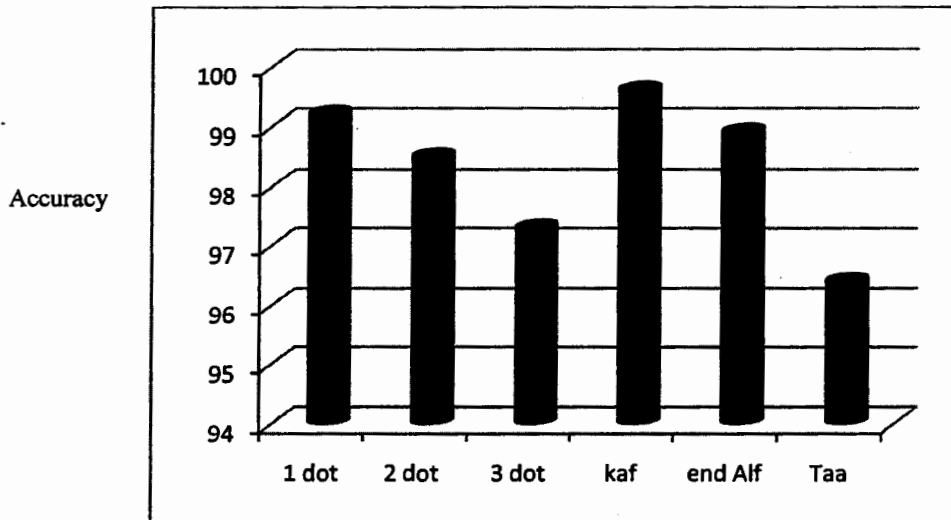


Figure 5.11: Recognition rate of diacritical marks

Algorithm: Diacritical Marks Fuzzy C-Mean Clustering

Input: Online handwritten strokes (base and diacritical marks)

Output: Diacritical marks with associated character i.e. position

Algorithm:

Step I: Select the number of cluster based on estimated characters.

Assign every cluster to every expected character.

Step II: Calculate the center of fuzzy cluster using position.

Step III: Repeat step IV and V for each diacritical mark

Step IV: Phase-1: Perform fuzzy estimation of diacritical marks clusters with each other.

Phase-II: Perform fuzzy estimation of diacritical cluster with characters.

Step V: Draw fuzzy projection on to the character to find the associated character.

5.1.5. Slant Correction

The character inclination that is normally found in cursive writing is called slant. Slant correction is one of the important steps to reduce the variation in the handwritten text and crucial part of preprocessing. The main purpose of slant correction is to reduce the variation in handwritten script specifically to improve the diacritical marks segmentation process, etc. which in turn can yield higher recognition accuracy. Since for Urdu script especially Nasta'liq the feature extraction and segmentation is highly dependent on the direction of writing thus it is important that vertical movement of pen is as perpendicular as it is possible, instead of diagonally slanted. Otherwise for diagonally slanted strokes written features will be poorly extracted or segmentation cannot be

performed correctly. For example for implicit segmentation it is very difficult to adjust the criteria for segmenting the stroke into subunits i.e. the input stroke is divided into slice/subunits for HMM based recognition so that accurate splitting is obtained; the strokes must be vertically perpendicular.

Generally handwritten strokes are not uniform, thus the slant may change with words and it may even be different within ligatures or strokes. Therefore the global uniform slant estimation is not guaranteed to provide good slant correction for handwritten text. Due to the complexity of Urdu handwritten script we need local slant correction instead of global slant correction. Whereas it's a difficult to estimate the slant on stroke bases due to limited information available for estimation. This problem mostly occurs where short strokes in height and width are involved i.e. *و، د، ذ*. These kind of short strokes contain limited information that helps in slant estimation. Thus the need is to normalize the stroke by estimating the slant locally based on the current stroke and some clue from the previous stroke to normalize the stroke on to the vertical axis. We present an algorithm that computes slant locally based on vertical projection with the help of neighbor strokes information.

For slant estimation only vertical and horizontal movements of pen are analyzed while slant correction is performed on only vertical axis. As the Urdu ligatures are written from top to bottom and left to right. The writers may write the ligatures little right or left slanted shown in Figure 5.12. The slant estimation is performed on only primary strokes because some of the secondary strokes are slanted naturally whereas alignment correction is performed on both primary and secondary strokes as shown in figure 6. For slant estimation, gravity fall of topmost strokes elements is used with additionally

diagonally vertical up or vertical down movement. The slant is corrected by vertical estimation of line from lowest measure towards new upwards estimated points as shown in figure 5.12 and figure 5.13.

If the topmost point angle with starting point angle is greater than 70 then no slant correction is performed. Otherwise slant correction is performed by using the gravity fall from estimated position to the lowest position and the next segment is attached with the vertically dropped position as shown in figure 5.13.

Vertical angle detection.

$$\theta = \tan^{-1} \left(\sqrt{\frac{(x - a)^2}{(x - a)^2 + (y - b)^2}} \right) \quad (5.10)$$

If $\theta \leq \alpha$

$$(x', y') = (x, b) \quad (5.11)$$

$$(x'_i, y'_i) = (x_i, b) \quad \text{For all intermediate point } i$$

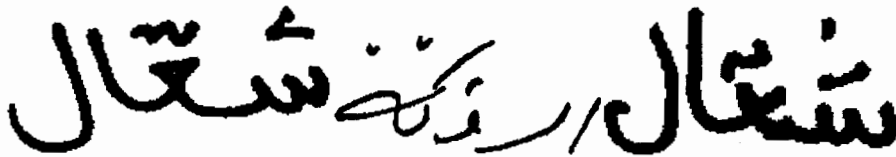


Figure 5.12: Left slanted, right slanted and normal words.

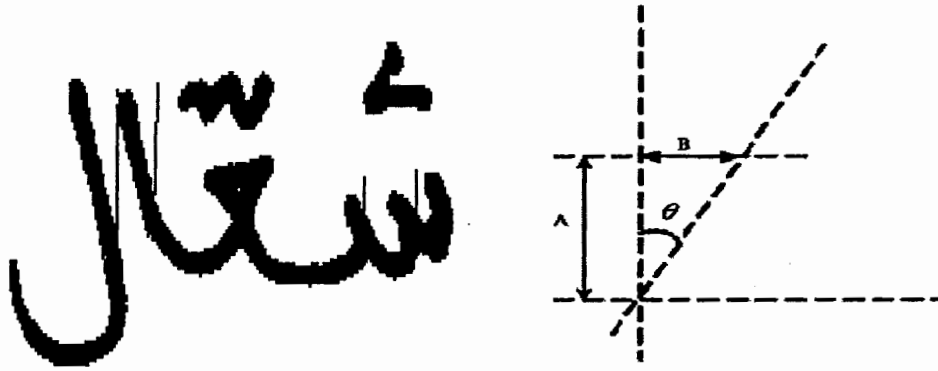


Figure 5.13: Locally angle computation for Slant normalization

As the secondary strokes are small and small in number thus there is no need to correct the slant of secondary strokes. But the issue is association of secondary strokes with it concerned character. The position of the secondary strokes also need to be improved based on the slant estimation of primary strokes correction angle.

5.1.6. Stroke Mapping

The input strokes are transformed into image to perform offline preprocessing steps. Only online finding are not enough for the Urdu online character recognition due to its characteristics. It is more suitable to involve offline information along with online information to increase the recognition rate.

It is difficult to write some ligatures i.e. لا, بصر, etc. shown in figure 5.14. Without lifting the pen as shown in figure 5.15, whereas online character recognition does not permit lifting of pen during the writing of ligature for Urdu script based languages. Thus to overcome this issue, an algorithm is proposed based on two rules:

Rule 1: Two consecutive ligatures are treated as one ligature if they overlap each other and starting point of second ligature is close to the ending of the first stroke and one of the strokes is vertical ending. This value threshold β is considered twice in vertical than in horizontal.

Rule 2: Two consecutive strokes are combined to form one stroke if the previous strokes ending point is very close to the start of second stroke, these two ligatures are combined and considered as a single ligature as shown in figure 5.15.



Figure 5.14: Base storks that are difficult to write in single stroke



Figure 5.15: (a) Resolved problem shown in figure 5.14. (b) Error in combining

$$\text{If } dist \mid (x, y), (m, n) \mid < \beta \quad (5.12)$$

If both strokes end at vertically then both strokes are combined

$$dist \mid (x, y), (m, n) \mid < \phi \quad (5.13)$$

Where α is depends upon the size of the strokes.

The issue in figure 5.15.b can be resolved by using the word level processing.

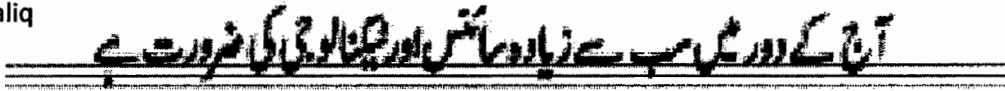
5.1.7. Baseline Estimation

The most significant operation in preprocessing is the estimation of baseline. Baseline is the virtual line on which characters are aligned to form the ligatures or words and it is the necessary requirement for both readers and writers. Character in different languages follow different baseline for example Hindi follows baseline in which the characters are aligned at the top while Roman characters are aligned at the bottom. Naskh writing style of Arabic follows bottom alignment whereas Nasta'liq is aligned at the center. Generally, baseline is kept in mind during both writing and reading. Baseline detection is not only used for automatic character recognition but is also helpful for human reading. If the characters are written without following the baseline, it is difficult to read the text even for human and error rate increase up to 10% even with the context sensitive interpretation. In automatic character recognition where no context knowledge is involved, baseline detection is the necessary part for better classification especially for Arabic script based languages. Moreover the localization of diacritical marks also requires the baseline estimation. Several baseline detection methods based on horizontal projection have been proposed in literatures as discussed in chapter 2 but are valid for complete text lines.

Urdu script based languages are written in many styles but mostly Nasta'liq and Naskh are followed. Urdu, Punjabi, etc. is written in Nasta'liq whereas Arabic, Pashto, etc. are written in Naskh style. We have proposed a novel technique for baseline estimation based on extraction of some primitives extracted from the ghost characters.

Nasta'liq and Naskh styles are mostly followed by Urdu script based languages. Figure 5.16 shows the modeling of Urdu on baseline and two descender lines. Different characters may appear at different descender lines. Due to the complexity of Nasta'liq over Naskh, one character in Nasta'liq may appear at different descender lines depending upon the associated characters whereas in Naskh style each character follows the baseline and does not depend upon its associated characters as shown in figure 5.17. Thus baseline estimation for Nasta'liq written text is more complex than Naskh style. Without pre knowledge of word structure it is difficult to estimate the baseline. The blue line in figure 5.16 shows the baseline for Nasta'liq and Naskh style. Figure 5.17.a and figure 5.17.b show that characters in Naskh follow the same baseline whereas in Nasta'liq they may follow different baselines depending on their context.

Nastaliq



آج کے دور میں سب سے زیادہ سائنس اور ٹیکنالوجی کی ضرورت ہے

Naskh



آج کے دور میں سب سے زیادہ سائنس اور ٹیکنالوجی کی ضرورت ہے

Base-line

First-descender-line

Second-descender-line

Figure 5.16: Baseline and Descender lines for Nasta'liq and Naskh font for Urdu [4]



Figure 5.17.a Baseline (Red) Naskh Style and blue line shows issues in baseline

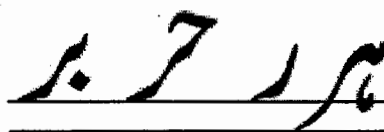


Figure 5.17.b Baseline (Red) for Nasta'liq

A novel baseline estimation method is presented that uses only the ghost character. The character has been made rotation invariant. The proposed approach is divided into two phases.

Phase I: Primary baseline estimation

Phase II: Locally baseline refinement

Phase III: Alignment Correction

5.1.7.1 Phase I. Primary baseline estimation

For primary baseline estimation we have used a projection based method. The primary baseline estimation gives a rough baseline for refinement by locally baseline estimation. The horizontal projection based method counts the number of pixels on the horizontal line. The line with the maximum number of pixels is considered as the initial

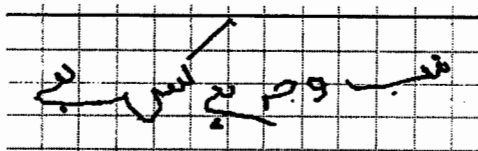
baseline. As the secondary strokes are removed during the phase-I, thus projection baseline is estimated based on ghost characters and gives good results by eliminating the influence of diacritical marks on baseline estimation as followed in earlier research[].

5.1.7.2 Phase II. Locally Baseline Refinement

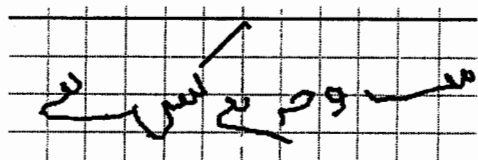
Although, projection baseline is robust and is easy to estimate but this method requires a long straight line of text. In case of handwritten text especially for online handwritten text the length of line/words may be very short and it is difficult to find a single baseline due to large variation in handwritten text. Locally baseline estimation is based on extraction of additional features like shape ending feature hedge, ray, long left to right. These features are used to refinement the baseline locally with the help of primary baseline. The features and locally baseline estimation is fully dependent on the style of script i.e. Nasta'liq has different sets of features with different rules as compared to Naskh. Figure 5.18 describes the proposed baseline extraction method.

Features lying on the baseline are extracted i.e. ray, bey, etc. as shown in figure 5.19 and then baseline is computed for these features instead of the whole stroke. As in Nasta'liq, the last character may not lie on the baseline, thus for Nasta'liq, the baseline for the last character is extracted and has less influence on the final baseline. Whereas for Naskh font, the local baseline has more influence on the final baseline because the major portion of each character lies on the baseline. The projection angle computed previously by the global projection method is further corrected based on the local baseline.

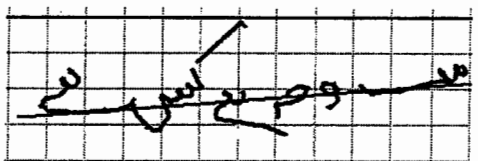
(a)



(b)



(c)



(d)

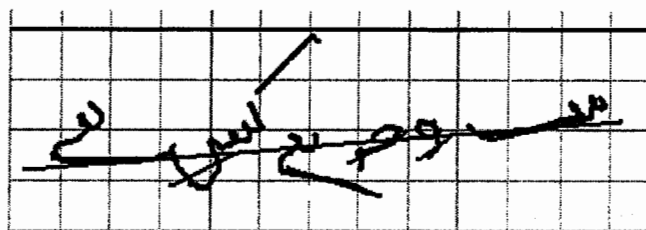


Figure 5.18: (a). Raw input strokes. (b): Ghost shapes after separation of secondary strokes. (c). Primary baseline estimation based on projection. (d): Locally baseline estimation based on features.

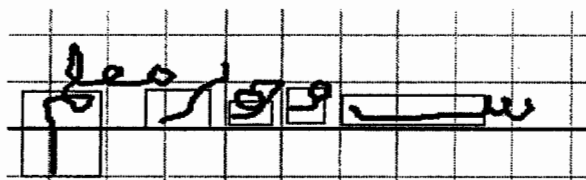


Figure 5.19: Features for locally baseline estimation

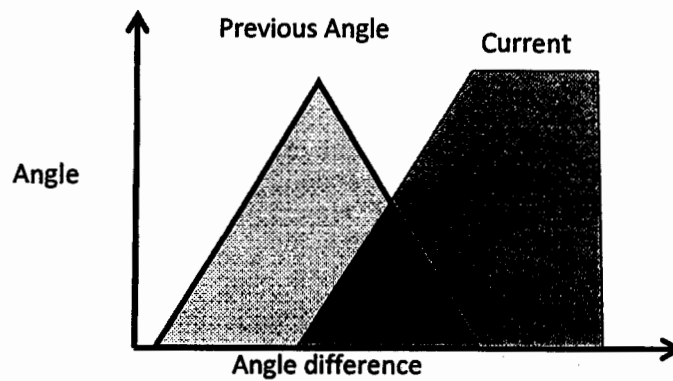


Figure 5.20 Fuzzy Membership Function for Angle fusion

For primary baseline estimation α is computed based on horizontal projection as shown in Figure 5.18.c on ghost character. The primary baseline estimation is used to reduce the error occurred by using the feature based approach. Then the second baseline is computed based on the features and previous baseline. The role of primary baseline is to compute the exact angle (β) for local baseline by using the following relation.

For each ligature

IF $|\alpha - \beta_i| < \theta$ then estimated angle $= (\alpha + \beta_i)/\phi$

ELSE estimated angle $= \alpha$

Where α , β are globally and locally computed angles of each ligature and ϕ is computed using the fuzzy membership function $\alpha' = (\phi_i \alpha + \phi_j \beta_i)$ as shown in figure 5.20 based on the previous and current angles.

If $|\alpha - \beta_i|$ is very small then $\alpha' = (\alpha + \beta_i)/2$

If $|\alpha - \beta_i|$ is small then $\phi_i = 0.2 \sim 0.3$ $\phi_j = 0.7 \sim 0.8$

If $|\alpha - \beta_i|$ is medium then $\phi_i = 0.3 \sim 0.4$ $\phi_j = 0.6 \sim 0.7$

If $|\alpha - \beta_i|$ is large then $\phi_i = 0.7 \sim 0.9$ $\phi_j = 0.3 \sim 0.1$

5.1.7.3 Alignment Correction

Alignment correction is performed on both ghost strokes and secondary strokes. The angle of secondary strokes is the same as the angle of the associated ghost stroke and performed using the following relation. (Association of diacritical marks based on writing order)

For ghost strokes and diacritical marks.

$$\begin{aligned} x' &= x \cos \alpha'_i - y \sin \beta_i \\ y' &= y \cos \beta_i - x \sin \alpha'_i \end{aligned} \quad (5.14)$$

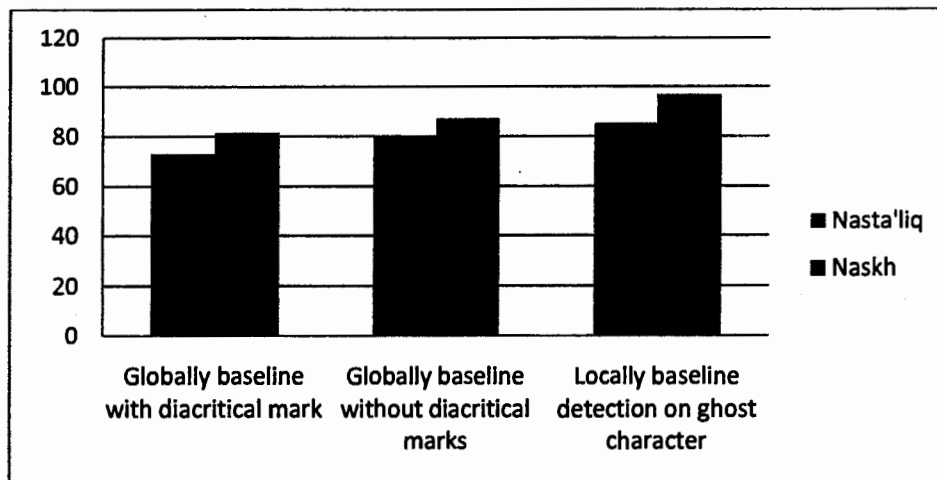


Figure 5.21: Recognition result of baseline estimation

5.2 Feature Extraction

A successful character recognition methodology depends on the choice of features used by the pattern classifier. Feature selection in pattern recognition problems involves the derivation of salient features from the preprocessed data in order to reduce the amount of data used by the classifier for classification and simultaneously provide the enhanced discriminatory power. In pattern recognition problems, feature selection involves the

evaluation of most discriminant features from the extracted. The direct recognition of handwritten stroke is almost impossible due to high variability of handwritten strokes. The feature selection phase also called dimensionality reduction. The aim of feature extraction from the input strokes is reducing the input pattern to avoid from complexities while keeping the high accuracy and the extraction of those distinct patterns that uniquely define the strokes and most important for classification while the task of human expert is to select those features that allow effective and efficient recognition. The feature extraction phase must be robust enough so that extracted features are small in number, produce less error during extraction and uniquely describe the shape of strokes. Generally, handwritten strokes are oriented lines, loops, and curves. These orientations play an important role in the classification of strokes. The directional and structural features i.e. loop, cusp, etc. plays important role in the classification and previous work shows that improved results have been obtained by combining structural and directional features. Structural features are the shape defining features and these are based on the instinctive aspects of writing and include loops, cusp, endpoints, starts points, etc. Sometimes, further preprocessing is required on this feature matrix to remove the unnecessary and noisy features by using language rules. Fuzzy rules have been used to extract the unique and meaningful shape defining features. Our main focus is on the structural features. Time variant structural feature are extracted from both online and offline stroke elements to identify on-line handwritten ligatures. The additional time information is utilized to extract features from the stroke elements in the order of occurrence. This arrangement helps us in combining the diacritical marks through estimation. Finally, post processing is applied on extracted features matrix to remove the

unnecessary or noisy features by modeling the biological concept and language rules i.e. communication between upper level (Level-4) and lower level (Level-3) to refine pattern at upper level. Figure 5.22 briefly describes the directional feature extraction process.

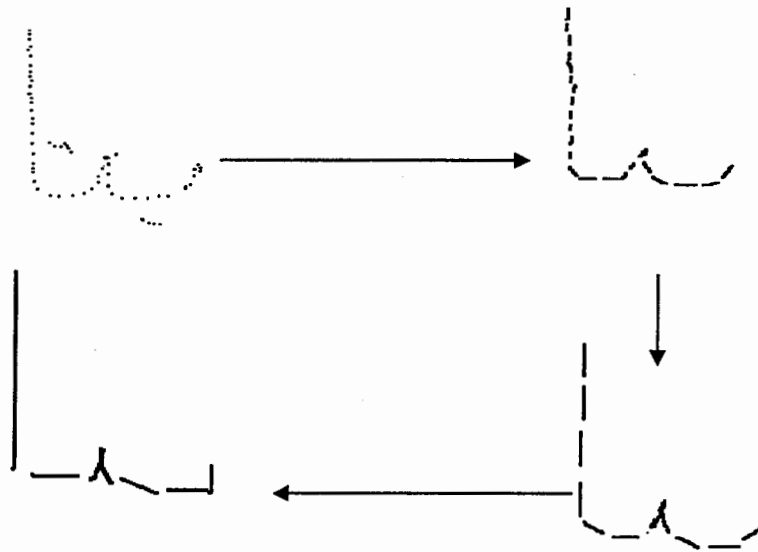


Figure 5.22: Biologically inspired feature extraction (divides the above figure in 4 parts and describes the whole picture in paragraph)

5.3.1 Directional Features

Handwritten strokes are oriented lines, curves, or poly-lines which play an important role in differentiating between various characters. For a long time, orientation or direction has been taken into account in handwritten character recognition. In early stages, character recognition using directional features was called directional pattern recognition [Fujisawa and Lui, 2003]. Chain code is the important way to extract the directional features as shown in figure 5.24. We have used chain codes and length for directional feature extraction based on fuzzy logic and context knowledge of the ligature itself,

instead of length of current feature as shown in figure 5.25. The relative fuzzy decision of directional features is performed in two steps. The first stage extracts small patterns and the second stage combines these small patterns to form large patterns. The decision of large patterns is performed in level-4 based on the associated patterns detected in the ligature using the fuzzy language rules.

As humans have high context knowledge and they can better differentiate the different sizes based on their context knowledge, similarly we have tried to model the context knowledge of ligature itself. Although it is very limited yet it helps to differentiate and extract more discriminating features.

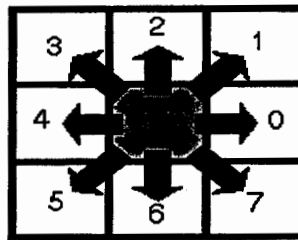


Figure 5.24: Chain code directional vectors

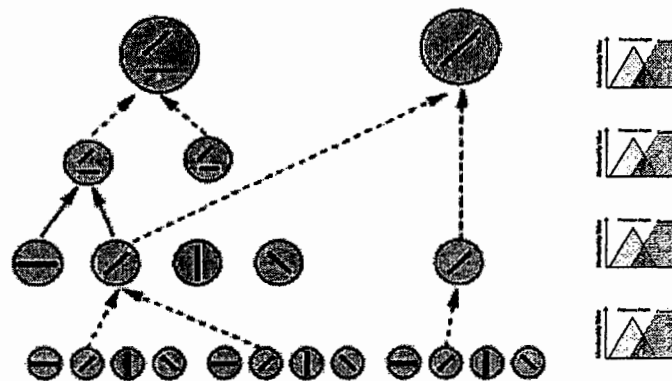




Figure 5.25 (a) Layered based low level to complex level feature extraction (b): Level-3 and Level-4: Small pattern to larger pattern extraction

- **Eight_Small_Movement:** This feature depends upon the small movement at start of the ligature either left, right, top, down or diagonal direction e.g. .
- **Start_Long_Down:** This feature is selected when the ligature was a straight long vertical downward in the beginning. e.g. .
- **Start_Long_Up:** This feature is selected when the ligature was a straight vertical upward in the beginning. As there is no word which starts from upward but this feature is used to differentiate numerals like and having same shapes.
- **Ending_Vertical_Long_Up:** This feature is selected when the ligature was a straight long vertical upward in the end. For e.g. .
- **Ending_Vertical_Long_down:** This feature is selected when the ligature was a straight long vertical downward in the end. E.g. .
- **Long_Horizontle_Left:** This feature is selected if during writing the ligature; the pen movement is very long and from right to left horizontally e.g. in .
- **Long_Digonal_Left:** This feature is selected if during writing the ligature the pen movement is long and from left to right diagonally like in .
- **Long_Left_Right:** If during writing the ligature, the pen movement is long and from left to right horizontally then the horizontal left to right is selected e.g. in .

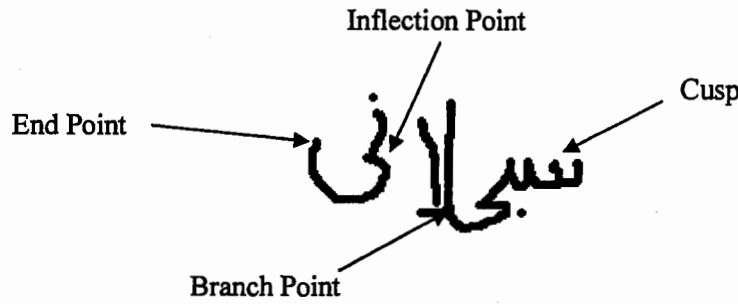


Figure 5.26: Some Structural feature (inflection point, cusp, end point, branch point)

5.2.2 Structural Features

- **Hedge_Right:** Some shapes have semi-circle. For such characters, feature called the semicircle are extracted when left to right semicircle present from right to left like .
- **Hedge_Left:** Similarly the ligatures containing the semicircle from left to right i.e.
- **Curve_Right:** The direction of writing of curves varies from right to left and from left to right. Therefore, curve right to left is selected for characters those having shapes like .
- **Curve_Left:** This feature is selected if the curve direction is from left to right then like

$$\alpha = \Theta_{(i+1,S)} - \Theta_{(i-1,S)} \quad (5.19)$$

$$\Theta_{i,S} = \text{Arc tan} \left[\frac{(y_i - y_{(i-k)})}{(x_i - x_{(i-k)})} \right] \quad (5.20)$$

- **Cusp_Up:** Cusps are the sharp turning point in a stroke. This feature is selected for the ligature which contains the up side cusps such as those present in س, .

$$\alpha = \Theta_{(i+1,T)} - \Theta_{(i-1,T)} \quad (5.21)$$

$$\beta = \Theta_{(i+1,T)} - \Theta_{(i,T)} \quad (5.22)$$

$$\beta - \alpha < 8 \quad (5.23)$$

- Cusp_Down: This feature is selected for the ligature which contains the downward cusps such as those present in *سہ، پھ، بھ*. Similarly using the above value of α, β .

$$\alpha - \beta < 8 \quad (5.24)$$

- Intersection: Whenever an intersection is encountered in a primary stroke this feature is selected i.e. *فل، ط*.
- Rey: This feature is selected if there is ray shape at the end of the stroke. If any ligature is a combination of ray or dal then this feature is also selected i.e. *ر، بد، د*.
- Loop_D1: In order to differentiate the loop the clockwise written loop is selected i.e. *ف، ق*.
- Loop_D2: In order to differentiate fay, qaf and meem, this feature is selected for Meem. The writing direction of the loop in meem is anti-clockwise. i.e. *ہم، جم، تم*.
- Loop_D3: To differentiate the loop in Swad from the other loop, this feature was identified and selected for Swad. As the swad loop is egg shaped so it is identified to separate the Swad from other loops like *ص، بھ*. Some time it is very difficult to differentiate between the loop-swad, loop-down, and circular hey and loop-up. If loop selection problem occurs then final decision is based on the fusion of directional and structural features describes in section 5.2.3.
- Loop_D4: In order to differentiate the loop in fee, Qaf, Meem and hey, this

feature is identified when isolated loop occur i.e. Hey “ہ”.

- C_Ray: This feature is selected if some characters combine with ray like رچر .
- C_Ye: This feature is very difficult to extract because there are very small variation involved. It is selected if some characters combine with choti ye like جیبی .
- Ayen_Jeem: This feature is used to resolve the problem in differentiating between jeem and ayen i.e. ع، عب. Its selection depends upon the previous features and extracted having using the following equation.

$$\beta = \Theta_{(i+1,s)} - \Theta_{(i,s)} \quad (5.26)$$

$$(5.27)$$

$$150 > \alpha < 180 \quad \& \quad 200 > \beta < 250$$

- Tawn: This feature is selected on the presence of tuan in any ligature like ط، سط . Tuan is selected when loop exist after long up and down.
- End_Hey: When character combines with hey this feature is selected. i.e. چہ بہ .

$$\alpha = \Theta_{(i+1,s)} - \Theta_{(i-1,s)}$$

Loop is an important feature to differentiate some similar numerals, loop is divided into two categories small (only for zero written in Urdu Numerals) and large for other numerals. To detect loops, the pen movement is recorded. If points from the recorded points intersect occurring point, previous recoded and current point forms a circle, thus a closed loop is detected. If the internal width and length is greater than α ($\alpha \leq 7$ pixel), then this closed loop will be considered as a large loop otherwise this closed loop will be considered as a small loop used for zero. Another issue is the open loops which always exist at the beginning of the stroke. To find open loop, from the beginning of the stroke,

the distance between recorded points are calculated with the current point, if it is less than β then it is considered at closed loop, then internal width and length of this closed loop is calculated to find greater loop or small loop. As the open loops exist at the beginning of the stroke, thus recorded length will be $\gamma < 6$ to avoid loop diction in some case like figure 5.23.

For closed loop,

$$\alpha_1 = \text{Max}(x) - \text{Min}(x) \quad (5.16)$$

$$\alpha_2 = \text{Max}(y) - \text{Min}(y)$$

For closed loop

$$\alpha = \left\{ \begin{array}{l} \leq \alpha_1, \alpha_2 \\ \text{other small loop} \end{array} \right\} \quad (5.17)$$

For open loop,

$$\alpha = \left\{ \begin{array}{l} \leq 6 \\ \text{other other feature} \end{array} \right\} \quad (5.18)$$

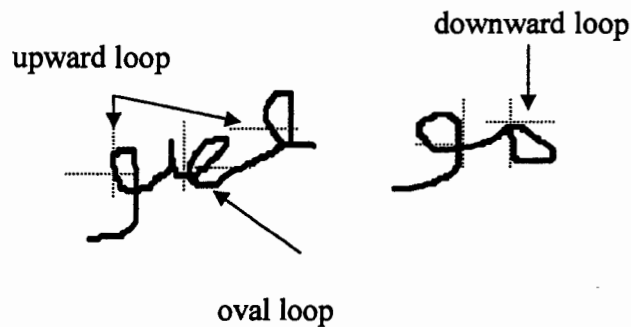


Figure 5.23: Loop confliction and detection

5.2.3. Feature Fusion:

Fusion of structural and directional features is performed based on relative position information of each feature.

$$F_S = \{f_1, f_2, \dots, f_n\} \quad (5.28)$$

$$F_D = \{f'_1, f'_2, \dots, f'_m\} \quad (5.29)$$

Where

$$f_i = (Z_i, x_a, y_b) \quad (5.30)$$

$$f'_j = (Z_j, x_a, y_b) \quad (5.31)$$

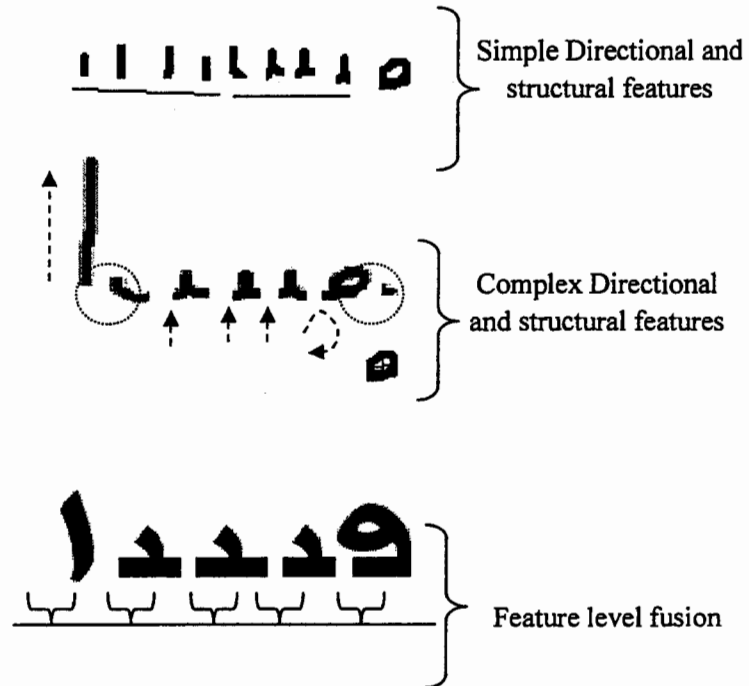


Figure 5.27: Level-4 Fusion of structural features and directional features



Figure 5.28: Fusion of structural and direction features

Z is the feature and x, y is the position of the features. The fused feature matrix is given as

$$F_F = \{ f_1^F, f_2^F, \dots, f_k^F \} \quad (5.32)$$

$$F_F = \{ f_1^F, f_2^F, \dots, f_k^F \} \quad (5.33)$$

5.3. Final Remarks

The first part presents the issues of preprocessing steps for handwritten online character recognition. As Arabic script based character recognition is difficult task due to complex language structure to attain good accuracy in pre-processing of the raw input strokes, which is crucial part. This phase describes the preprocessing steps for online character recognition by considering the input strokes from both online and offline features and their fusion to increase the efficiency of the recognition system. The proposed technique is also a necessary step towards character recognition, person identification, personality determination where input data is processed from all perspectives. A number of

preprocessing steps are presented in this chapter i.e. secondary stroke separation, baseline estimation, etc. The secondary strokes localization algorithm is based on fuzzy logic to handle secondary strokes through the vertically projection of the stroke. In Urdu script; secondary strokes are written above or below the word and in some cases they may appear little before or within the word-part with respect to the horizontal axis. Baseline estimation is one of the most difficult steps and it has direct influence on accuracy. We have presented a novel technique for baseline estimation for cursive handwritten Urdu script written in Nasta'liq and Naskh styles. Firstly the secondary strokes are segmented from the raw input strokes. Then primary baseline is extracted using the horizontal projection on ghost shapes. Finally the locally baseline of each ligature is estimated based on the features and primary baseline estimation. The presented approach gives good results due to mixture of local baseline estimation over global baseline estimation and reduction of diacritical marks. The proposed method provides accuracy of 85.3% and 96.7% for Nasta'liq and Naskh font respectively for baseline estimation.

The second phase of this chapter present biologically inspired feature extraction. As successful character recognition methodology depends upon the particular choice of features used by the classifier. Thus it is the most critical part in any recognition problem. Due to the complexity and variation, direct recognition of handwritten stroke is almost impossible. The literature shows that directional and structural features i.e. loop, cusp, etc. played important role in the classification. Fuzzy rules have been used to extract the unique and meaningful directional and structural features and shape defining patterns i.e. loops, cusp, endpoints, starts points, etc. Further post processing is also applied on

extracted features to remove the unnecessary and noisy features by modeling the language rules and fusion of directional and structural features.

CHAPTER 6

BIO-INSPIRED ARABIC CHARACTER RECOGNITION

Bio-inspired pattern recognition has been under consideration by computer scientists and neuroscientists since the mid-eighties so as to model the human vision system into applications of pattern recognition. The understanding of visual information processing and perception principle is one of the challenging tasks for contemporary science. A deeper review of biological human vision helps to advance the pattern recognition research and make it more robust. The natural world is enormous, dynamic, incredibly diverse, and highly complex. Image perception in human indicates that processing is integrated at different levels. The integration of principles of biological vision with image processing helps to achieve robust system in terms of speed, accuracy and efficiency. The visual perception of humans is very powerful in pattern recognition problems. The selectivity, transformation invariance, speed and context knowledge are the most important features of human visual perception. Moreover humans are able to detect the familiar and unfamiliar objects even in variable environments. The computer can perform complex problems more precisely, efficiently and much faster than humans

in the domain of data processing however it still lags behind in many pattern recognition problems.

Despite the inherent challenges of surviving in real world, biological organisms evolve, self-organize, self-repair, navigate and have been evolving over millions of years and adapting to an ever-changing environment. The real time data may be ambiguous and incomplete due to several reasons i.e. blurred edges, irregular and missing patterns. The high context knowledge, adaptability and learning capability of humans provides robust intelligence to complex problems. The interpretation of complex patterns and transformation of these complex patterns into behaviorally understandable signals is important in our daily lives. The human brain has achieved this robustness through thousands of years of evolution. Even the brain of lower animals is too complex to imitate. To understand how human brain can achieve these extraordinary abilities, human visual system has been described here briefly in chapter 3. This fact suggests modeling of human visual perception in the real world image processing algorithms. These systems can be modeled using fuzzy logic, evolutionary computation and neuro-computing.

Biologically inspired solutions are gaining more interest as wide spread use of artificial intelligence applications is materializing [Halem, 2009], [Bosner, 2006], [Bosner et.al 2006], [Forbes, 2005], [Collins et.al, 2004], [Benyus, 1997], [Anastas, et.al, 2000], [Papanek, 1984], [Marinakis et.al, 200]), [Borji et.al, 2008], [Jeng et.al, 2009], [Gupta, 1998], [Chang et.al, 2008]. The biologically inspired image processing needs expertise in two fields, human biological perception i.e. neurology and modeling the concept in computer i.e. computer science. Computer scientists and neurologists strive to build systems that can efficiently recognize the graphical marks. We performed the study

for character recognition based on biologically inspired character recognition as shown in figure 6.1 and model the human visual system into our system.

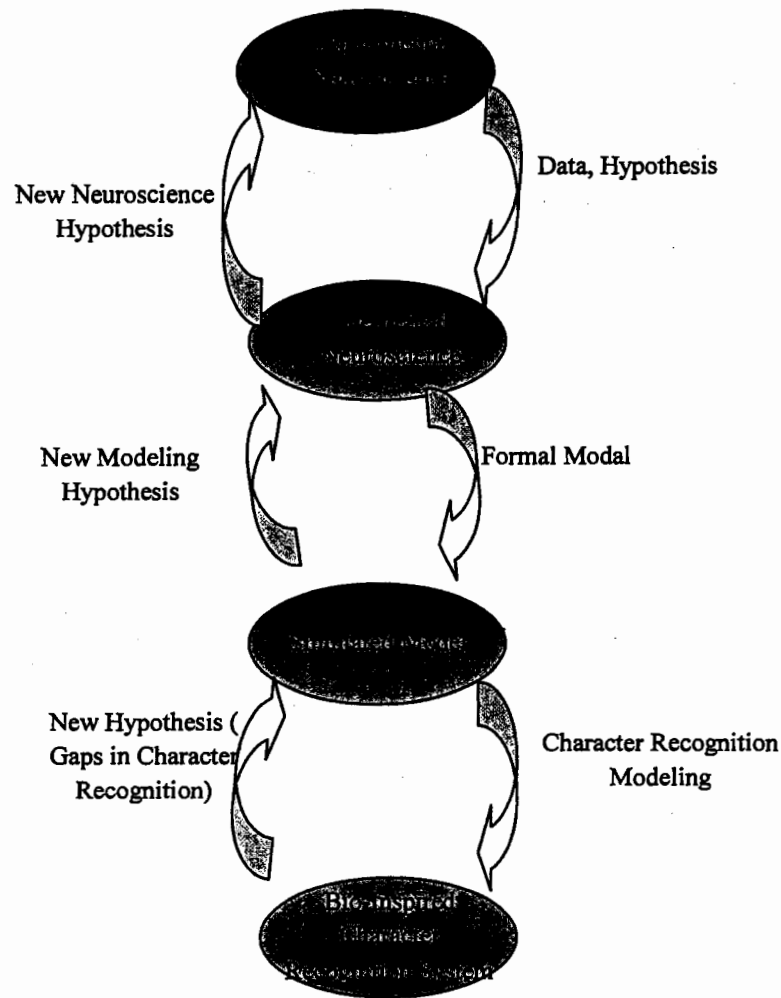


Figure 6.1: Framework for the study of human vision for character recognition

We performed the study on human visual concept and model into computer world. Several experiments have been performed and model is simulated. We model the previous study performed for object recognition.

Several uncertainties exist in cognitive systems due to numerous reasons i.e. static, random, physical, chaotic and fuzzy. These uncertainties have different effect depending on the type of problem whereas it is very difficult to classify the problem into different level of uncertainty. Vast variation and inconsistencies make handwritten Arabic script based character recognition more complex than any other language. Thus a careful observation is required to divide the complex handwriting recognition issue into uncertainty levels. We use the human ventral stream concept to divide the character recognition problem into different uncertainty levels.

This chapter presents a multilayered fuzzy rules based expert system for handwritten Multilanguage character recognition i.e. Arabic script based languages inspired by human biological perception.

6.1 Multilayered Human Visual System

Several complex image processing applications have been studied under human biological perception of object especially in motion detection, object recognition, etc. As a model of V1 neurons (as described in chapter 3) Gabor filter has been used for many applications of image processing i.e. edge detection, writer identification, texture classification, etc. Brain does not store the picture of object; it only possesses the spatial information of the objects. This spatial information of object is encoded in neurons i.e. various orientations, lines, edges and endpoints, etc. Recognition of the object is based on these spatial relationships stored in the mind instead of image copy. The visual sense is activated in response to the relative position and orientation of these characteristics.

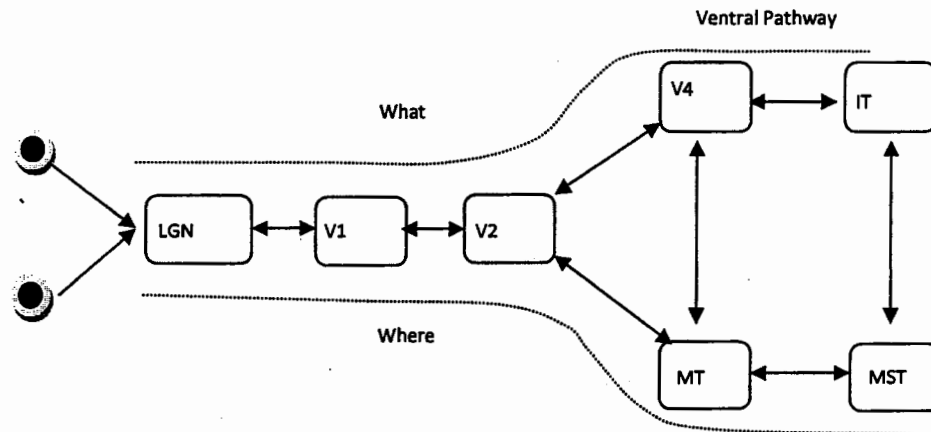


Figure 6.2: Human visual system

The study on monkey has shown that object recognition is feed forward for fast recognition whereas the visual perception is done in hierarchical manner and is sensitive to edges, lines, etc. [Weitzenfeld, 2008]. Sharp contours are more discriminating than gradients and straight lines are more obvious as compared to irregular patterns. Early visual parts such as retina, LGN, V1, V2 take part in extraction of simple features like edges, orientation, color, etc. The neuron structure becomes more complex and its receptive field becomes larger in the upper side of ventral stream. As the human visual system is more sensitive to curved line, other geometrical points along with the context knowledge, therefore the human recognition system is more powerful for complex shapes i.e. handwritten character recognition, motion detection, etc.

The human recognition system consists of seven or eight major layers and four layers are used for computation. Lower layer neurons are more precise for simple features as compared to upper layer features. Basic division of human visual system is shown in figure 6.2. Visual information covered by the two eyes is transmitted on to the lateral

geniculate nucleus (LGN). The striate cortex V1 is the primary visual cortex responsible for static and moving objects and receive the information from LGN.

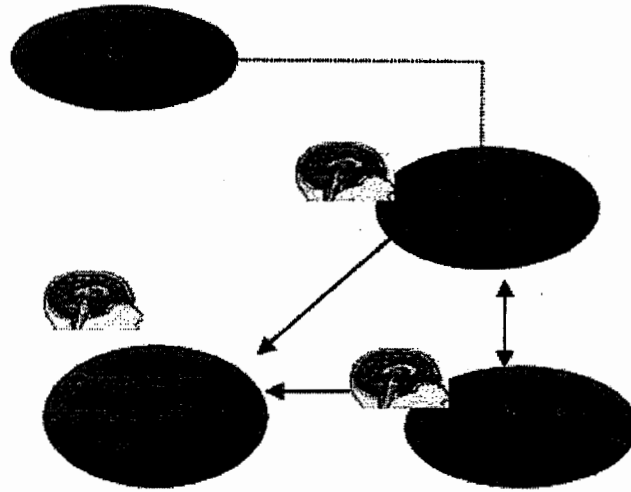


Figure. 6.3 Human behavior diagram for character recognition

Lower layer neurons (V1) are more responsive for simple features as compared to upper layer features. V2 is responsible for orientation and spatial properties such as contour and boundary computation, etc. The contour extraction and feature grouping is performed at first level by the combination of V1 and V2. Therefore the first stage is responsible for object motion detection, contour extraction, background and object detection. The second stage is responsible for complex features and third stage involves combining these large number of features extracted during stage 1 and stage 2 to form a pattern. In the fourth stage “object generalization” in a view specific manner is performed. The high context knowledge database plays an important role in the recognition of complex patterns, without which; human visual perception gives much worse results. Although it is very difficult to build such a big memory as of human that

can be used as clue for recognizing the complex objects but we can divide the complex task into layers inspired by human biological system and each layer can be focused by different uncertainty levels using fuzzy logics.

6.2 Multilayered Character Recognition

Patterns of object are stored in human mind instead of copy of objects and recognition is based on layer by layer computation of features. Unlike the traditional machine the human recognize the text in different phases. The early phase separates the needed character information from the background i.e. segment the character from noisy background. The second stage is recognition which obtains minute details from the earlier levels. The information in higher vertex revises the information from the earlier phase during recognition. Naturally, every language has its own properties i.e. writing rules, features, and unique complexities and therefore the classification process is dependent on the language rules.

Traditional character recognition systems face problems in handwritten Arabic script based character recognition due to variations and complexity in the script. The biological concept to recognize the Arabic script in multilayered manner is inspired by human biological system. Fuzzy logic with additional linguistics rules of Arabic script is used to handle uncertainties involved in handwritten character recognition. The language rules play an important role in solving the uncertainties using fuzzy logic. By using the fuzzy logic, help of linguistic rules and biological division of recognition process, complex handwritten text recognition such as Arabic and Urdu can be performed. The

processing of levels L_1, L_2, \dots, L_6 are based on the fuzzy triangular member function, context knowledge of previous word and knowledge of word itself for next sub-part.

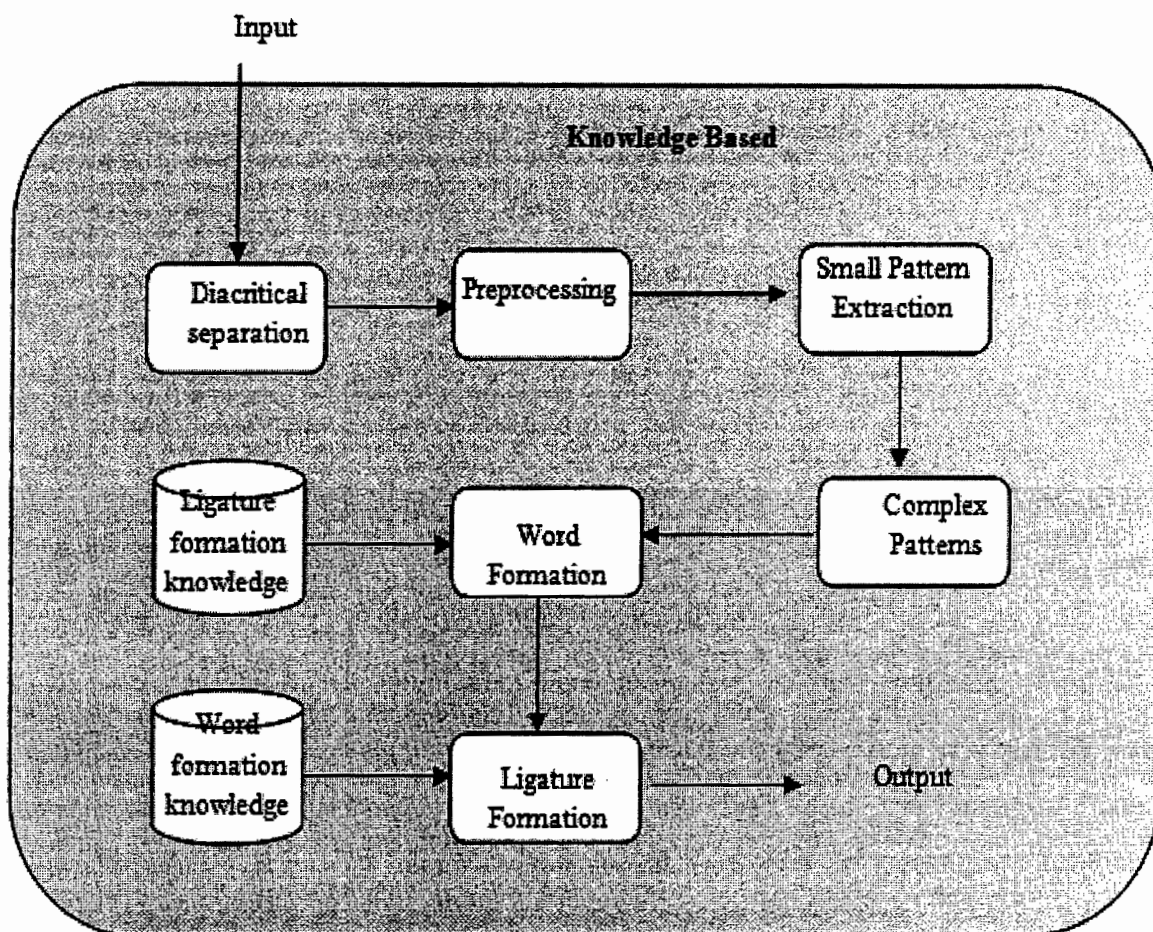


Figure. 6.4 Bio-Inspired Machine based character recognition (describe)

At level-0, the stroke information (x,y) from the input device is acquired. Level-1 is responsible for segmentation of diacritical marks whereas level-2 accounts for geometrical processing. Level 3 is responsible for dividing the strokes into sub patterns and extraction of sub patterns. Chain code is used at level 3 to find the small sub patterns. This level combines the small sub patterns (chain code) to form large patterns (lines and

curves). The complex pattern extraction from simple pattern and their combination to form identified sub-unit (features) is performed at level 4. Level 5 is responsible ligature formation based on the recognized features, and association of diacritical marks. For further refinement; context clue at level 6 can be used for word formation. The rules formed at one level form the input at other levels, where these rules are combined to form new rules with the additional information of linguistic rules and summarization of rules is layer by layer as shown below.

Level-1: *If* Premise_{1i}*then* Conclusion_{1i} and *If* LinguisticPremise_{1a}*then* LinguisticConclusion_{1a}

Where Premise_{1i} \in Level 0 and LigusticPremise_{1a} \in Language

Level-2: *If* Premise_{2j}*then* Conclusion_{2j} *If* LinguisticPremise_{2b}*then* Linguistic Conclusion_{2b}

Where Premise_{1k} \in Conclusion_{1k} and LigusticPremise_{2b} \in Language && LinguisticConclusion_{1a}

Level-6: *If* Premise_{6z}*then* Conclusion_{6z} *If* LinguisticPremise_{2g}*then* Linguistic Conclusion_{2g}

Where Premise_{6z} \in Conclusion_{5y} and LigusticPremise_{6g} \in Language && LinguisticConclusion_{5f}

Where Premise_{1i} is the set of rules for diacritical marks separation and Conclusion_{1i} the separated diacritical marks with associated positions and Linguistic Premise, Linguistic

conclusion are the primacies and conclusion based on linguistics rules. The above rules are the basis of multilayered biological inspired character recognition system.

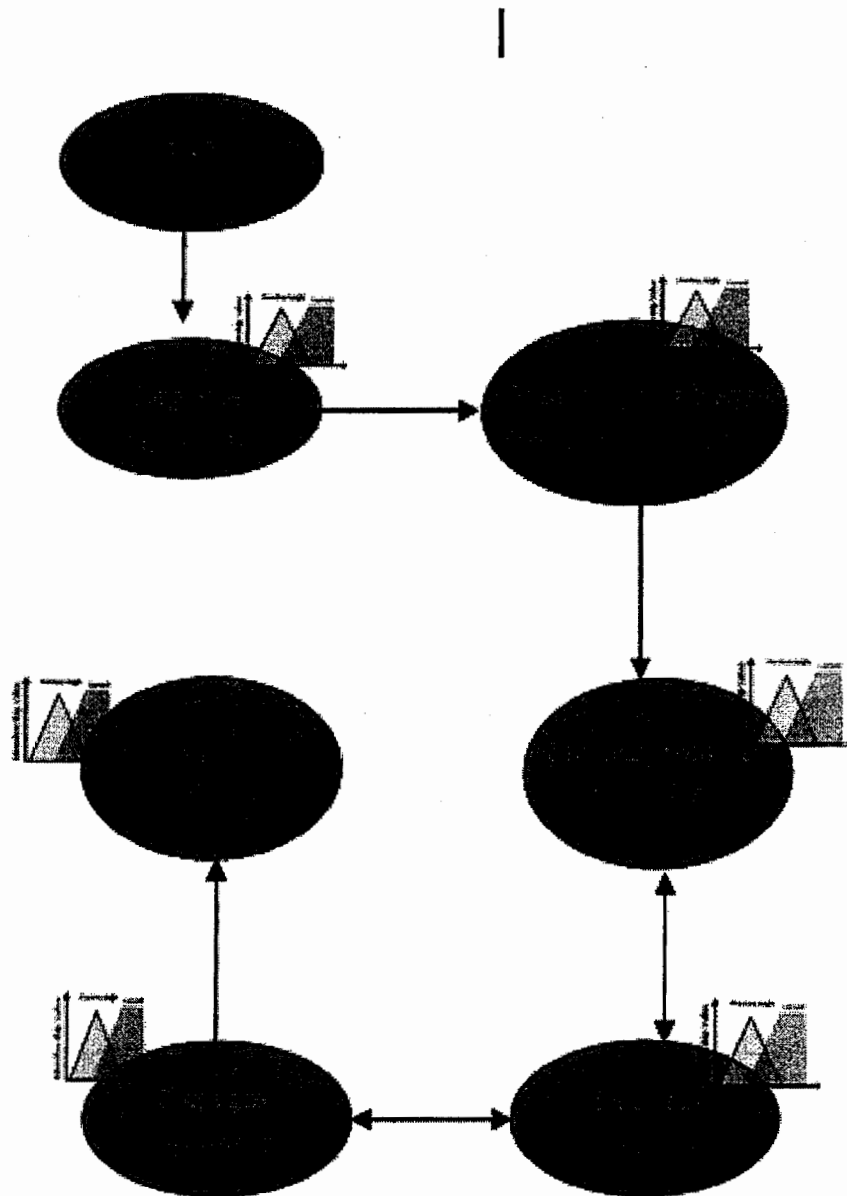


Figure 6.5. Machine based character recognition behavior diagram

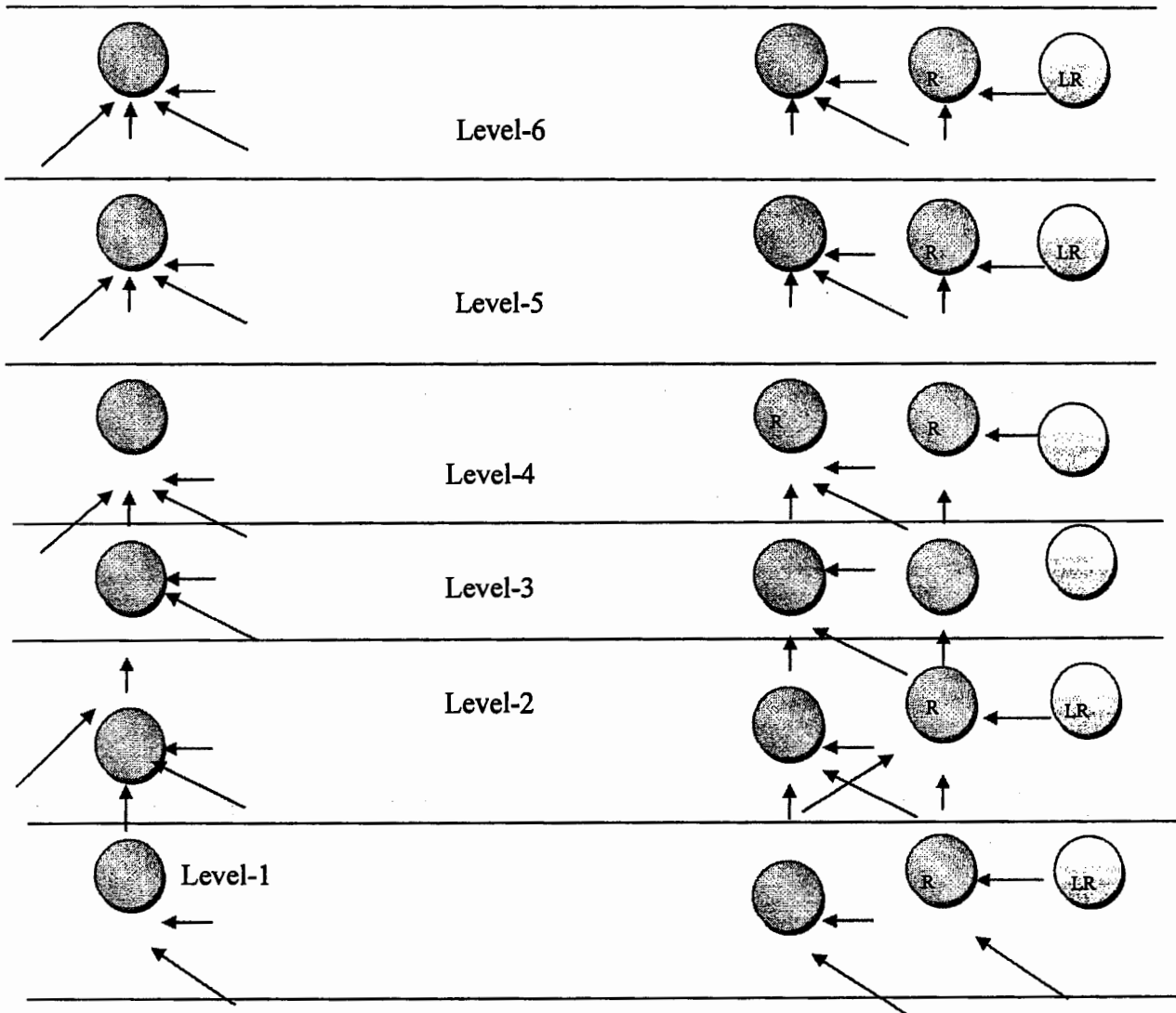


Figure 6.6: Multilevel bio-inspired rule model

Level-1: *If* $Premise_{1i}$ *then* $Conclusion_{1i}$ and *If* $LinguisticPremise_{1a}$ *then* $LinguisticConclusion_{1a}$

Where $Premise_{1i} \in \text{Level } 0$ and $LinguisticPremise_{1a} \in \text{Language}$

Level-2: *If* $Premise_{2j}$ *then* $Conclusion_{2j}$ *If* $LinguisticPremise_{2b}$ *then* $Linguistic Conclusion_{2b}$

Where $Premise_{1k} \in Conclusion_{1k}$ and $LigusticPremise_{2b} \in Language$ &&
 $LinguisticConclusion_{1a}$

:

:

Level-6: *If* $Premise_{6z}$ *then* $Conclusion_{6z}$ *If* $LigusticPremise_{2g}$ *then* $LinguisticConclusion_{2g}$

Where $Premise_{6z} \in Conclusion_{5y}$ and $LigusticPremise_{6g} \in Language$ &&
 $LinguisticConclusion_{5f}$

Where $Premise_{1i}$ is the set of rules for diacritical marks separation and $Conclusion_{1i}$ the separated diacritical marks with associated positions and Linguistic Premise, Linguistic conclusion are the primacies and conclusion based on linguistics rules. The above rules are the basis of multilayered biological inspired character recognition system.

6.2.1 Level-0: Stroke Acquisition

The input strokes are obtained through the digital pen/Wacom 4 electronic digitizing tablet contains (x, y) coordinates and timing information, writing force, speed. Force and speed are not helpful for character recognition due to variety in writing speed and writing force and are mostly used for person identification, personality determination, etc. Thus, time and stroke elements can be used for handwritten character recognition.

6.2.2 Level-1: Low Level Processing:

Level 1 is responsible of diacritical marks segmentation and primary stroke smoothing. Both operations have been described in chapter 5. Arabic script based languages are rich in diacritical marks. The same ghost stroke has different interpretation

based on the diacritical marks. Arabic script based languages contains zero or more secondary strokes associated with each primary stroke. We observed the human perception for Arabic script recognition. A study on 50, native readers of Arabic, Urdu, Punjabi, Sindhi and Pashto on the mostly followed Naskh and Nasta'liq fonts was performed. The study showed that the first step in recognition of Arabic script based languages is to separate the diacritical marks from the basic shapes. We have proposed fuzzy based separation of diacritical marks using delayed stroke segmentation algorithm as discussed in 5.1. This layer is responsible for separation of diacritical marks and primary stroke noise reduction based on the time and position information.

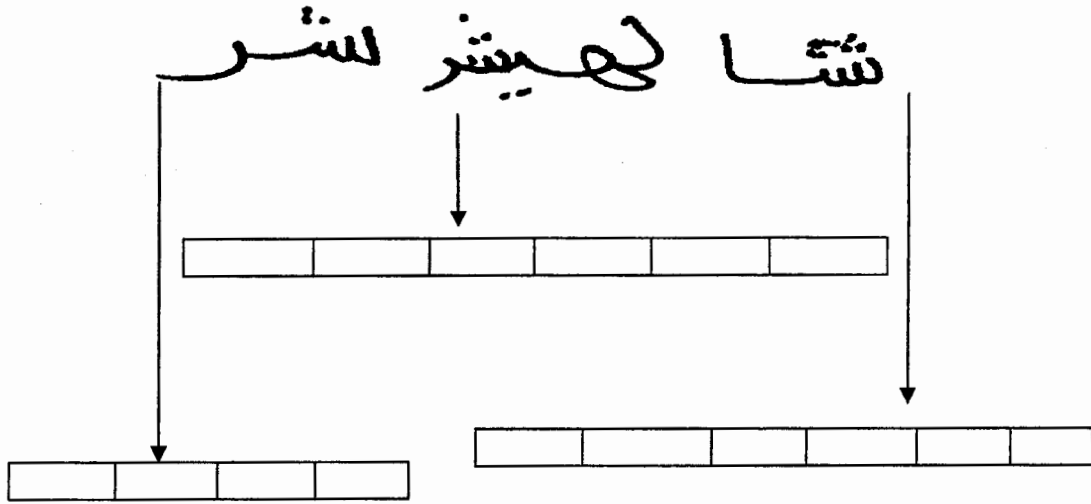


Figure 6.6: Diacritical marks issue with three dots in Naskh style

ا ح د ر س ص ط ع و
ک گ ل م ن و ه ی ع ه

Figure 6.8: Basic ghost character for Arabic script based languages written in Nasta'liq

Figure 6.6 shows the result of different ligature based on different diacritical marks association based on the position of diacritical marks. Figure 6.7 shows the ghost character for Urdu script written in Nasta'liq script.

6.2.3 Level-2: Geometrical Processing

The level 2 is responsible for geometrics computation, which is performed similar to human visual processing where geometrical processing is performed in visual cortex (lower area). It is very difficult to recognize the irregular, non-skewed, non-baseline character. Even for the humans, who have brilliant image capturing device and context knowledge, reading speed and accuracy of non-normalized text is lesser as compared to normalized text and that written on baseline. Although the human visual perception is based on translation and rotation invariant features, yet baseline has some influence on their perception power. At level-2, we have performed preprocessing steps like baseline estimation, stroke mapping and slant correction. These operations are based on the fuzzy logic and context knowledge. Although the context knowledge used in the proposed system is limited as compared to the humans, but it can resolve most of the uncertainties involved in stroke recognition i.e. the baseline is computed by the locally computed angle with additional knowledge of previous word angle. The knowledge of previous word may or may not take part in the correction of baseline as discussed in detail in chapter 5. The decision is based on the fuzzy rules i.e. if the baseline angle computed for current word is closer to the angle of previous and word length is large, then the previous has influence on current angle estimation. The fuzzy rules used for baseline estimation have been discussed in detail in chapter 5.

6.2.4 Level 3: Simple Feature Extraction

A successful character recognition methodology depends on the particular choice of features used by the pattern classifier. Feature selection in pattern recognition problems involves the derivation of salient features from the raw data input in order to reduce the amount of data used by the classifier and simultaneously provides the enhanced discriminatory power. The aim of feature extraction from the input strokes is to reduce the input patterns to avoid complexities while keeping the accuracy and the extraction of those distinct patterns that uniquely define the strokes. Generally handwritten strokes are oriented line, loops, and curves. These orientations play an important role in the classification of strokes. The directional and structural features i.e. loop, cusp, etc. plays important role in the classification and previous work shows that result are improved by combining structural and directional features. Level 3 and level 4 is responsible of feature extraction. Level 3 extract the small sub patterns and input to the level 4 to find complex patterns from these small patterns.

Level 3 is responsible for dividing the strokes into sub patterns and extraction of sub patterns. Raw data consist of stroke elements thus, the first level pattern extraction is chain code computation from the stroke elements. Now we have two types of sub patterns, chain code and stroke elements. Firstly, the stroke elements are unitized to extract more complex patterns i.e. loops, hedge. At level-3, the simple structural patterns e.g. intersection point, close passing point, cusp, loops, curves, etc. are extracted shown in figure 6.10. Secondly, extracted chain code is further used to find bigger and complex patterns from small patterns like directional features small vertical lines, horizontal lines

shown in figure 25 and figure 26. Moreover, the position of each pattern is kept in record which will be used for fusion of structural and directional features.

Table: 6.1 Biological Inspired Character recognition

Algorithm: Multilayered Biological Inspired Urdu Character Recognition System

Step-1: RUN Level-1 Segmentation of diacritical mark from the words and primary stroke noise reduction

INPUT: Raw input Strokes, Low Level Linguistics Rules

OUTPUT: Primary Strokes, Associated Diacritical Marks, Position Coordinates of each Diacritical Mark

Step-2: RUN Level-2 Geometrical Computation for Baseline estimation, Slant correction, Stroke Mapping

INPUT: Primary Strokes, Associated Diacritical Marks, Position, Low Level Linguistics Rules

OUTPUT: Normalized Word.

Step-3: RUN Level-3 Simple Feature Computation

INPUT: Normalized Words, Low Level Linguistics Rules

OUTPUT: Simple Features

Step-4: RUN Level-4 Complex Feature Extraction, Formation of Valid Sub Part

INPUT: Simple Features, High level Linguistic Rule

OUTPUT: Complex features, Recognized sub parts

Step-5: RUN Level-5: Formation of Valid Ligature [Character Level Recognition]

INPUT: Recognized Sub Parts, Diacritical Marks, Position Information

OUTPUT: Recognized Ligature

Step-6: RUN Level-6: Context Level Correction [Word Level Recognition]

INPUT: Recognized Word, Linguistic Rules

OUTPUT: Recognized Word at Context Level

Level-0

Level-1 Diacritical Marks Separation, Noise Reduction

Level-2: Baseline, Slant etc. estimation

Level-3: Simple Feature Extraction


Loop, Cusp, Cusp, Cusp, Vertical Up, Vertical Up, Vertical Up

Level-4: complex Feature Formation

Loop Upward, Cusp Upward, Cusp Upward, Cusp Upward, Long Vertical Up

Level-5: Character Level Recognition

Fuzzy Diacritical Strokes Mapping on Character, Rule based Character from Complex Features

Level-6: Context Level Computation

Figure: 6.9: Simple example based on multilayer based approach

6.2.5 Level-4: Complex Features Extraction

The level-4 is responsible for deduction of complex pattern extractions from simple patterns and then the complex patterns are organized to deduce the independent character. These semi-independent characters are candidate for independent character based on the surrounded character. The candidacy of character is confirmed based on the diacritical marks. The level-4 can only represent the possible character as shown in figure 6.9 and figure 6.10.

The simple pattern obtained at level 3 are combined to form complex patterns i.e. long vertical up, long vertical down, small vertical up, cusp, loop, hedge, etc. Moreover complex pattern directions are also extracted i.e. hedge cusp directions, direction, loop direction and loop shapes. Finally the fusion of directional and structural features is performed to get more robust features. The feature level fusion is based on the position information recorded at each feature point shown in figure 6.10. L4 gives the recognized subpart i.e. characters in the strokes. These recognized dot-less characters are candidate character. The decision is based on the mapping of diacritical marks at level 5.

6.2.6 Level-5: Character Formation

These recognized characters are combined using the linguistics rules and rule based approach at level 5. The linguistics rules, associated dot, recognized sub unit with associated units and diacritical positional information based on projection is used for mapping the secondary strokes. Figure 6.11 and Figure 6.13 describe the mapping of diacritical marks on to the primary strokes as similar to the human visual system shown in figure 6.12. The position information of diacritical marks and candidate character

confirm the candidate character as a recognized character based on diacritics. The knowledge of surrounding character is very important for character mapping i.e. if the three independent characters (cusp) are combined to form one character 'seen' if all have no diacritical marks or considered as 'sheen' if three dots appear and there will be no surrounding dots. The dots may appear as combination of two and one. In this case the associated character are three, base on the associated character, the decision is 'sheen'. It may happen that the associated characters are two and dots are in order of two and one. Based on this the result is 'tey and 'noon'.

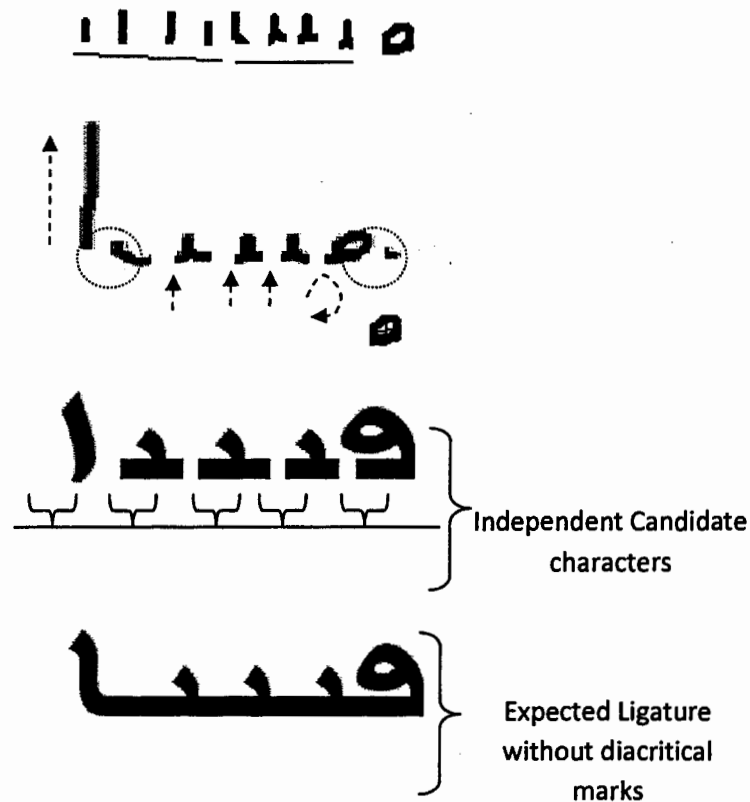


Figure 6.10: Level-4 Computation of complex features and candidate independent character

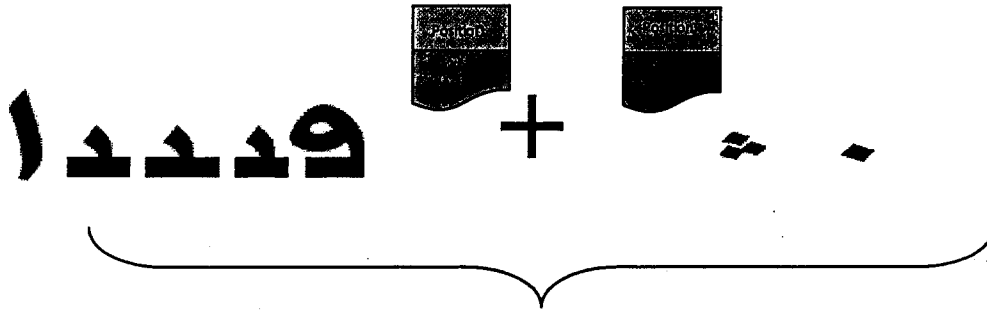


Figure 6.11: Combination of diacritical marks with respect to languages

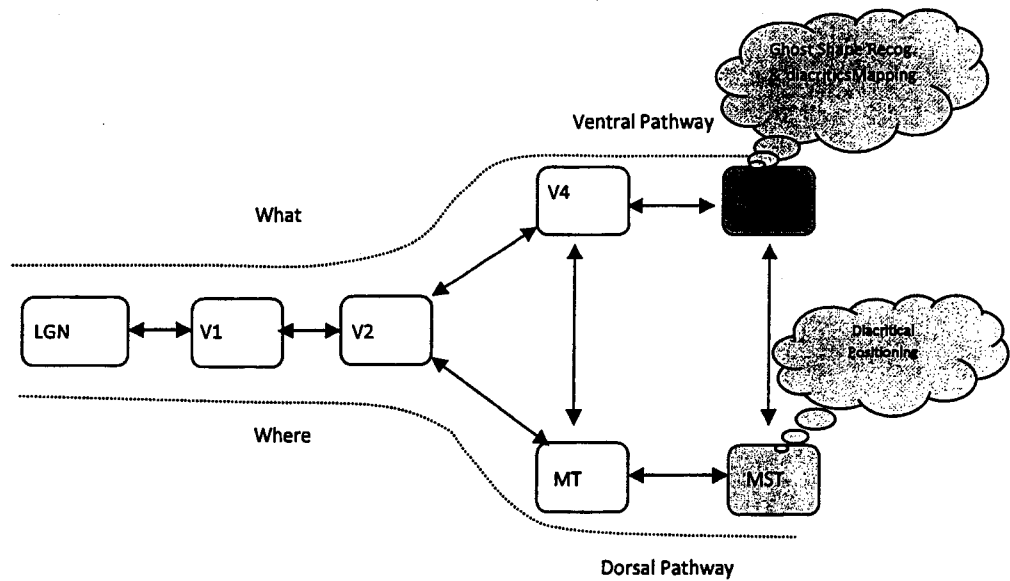


Figure 6.12: Bio-Inspired diacritical mapping

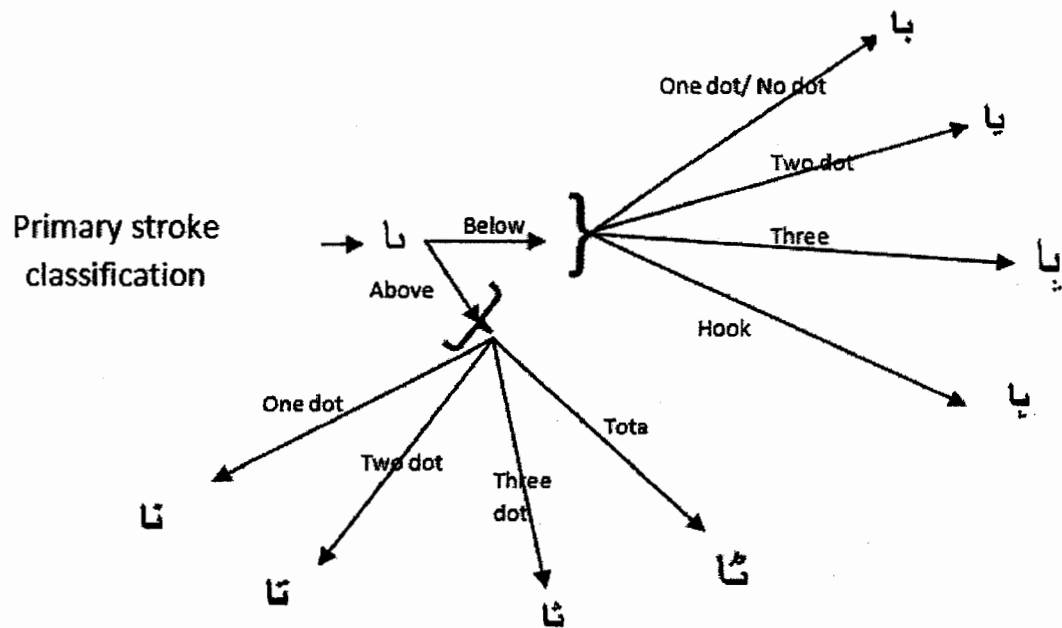


Figure 6.13: Combination of diacritical marks with respect to languages

6.2.7 Level-6: Word Formation

HMM can be used for word level recognition shown in figure 6.14.

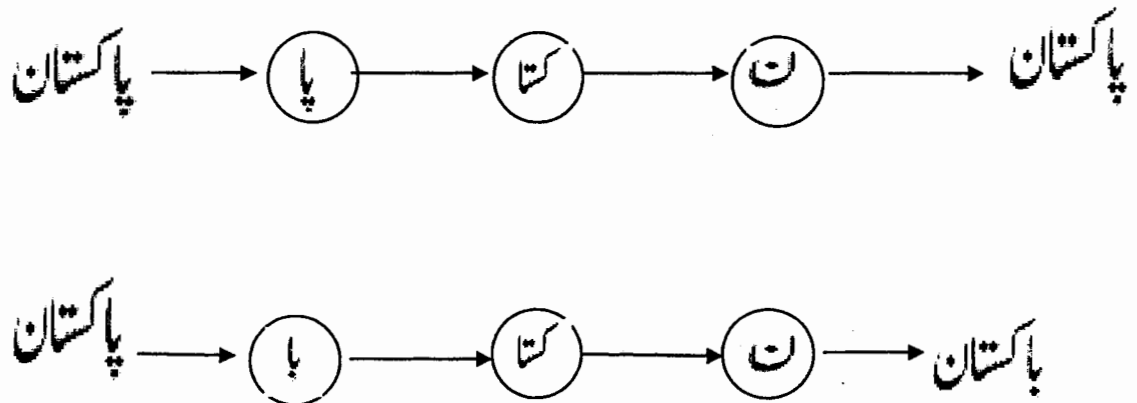


Figure 6.14: (a) Ligature formation by combination of diacritical marks and ghost character (b) Word level recognition HMM for same word in Urdu and Arabic.

6.3 Discussions

Fuzzy logic has proved itself to be a powerful tool for classification of irregular and complex patterns. Naturally human recognition system is translation, rotation and position invariant. Thus it is very difficult to imitate the human perception capabilities in machines. Translation, rotation and position have some effect on the recognition system in terms of reading speed. With the help of fuzzy modeling, linguistic rules and biological division of handwritten ligatures, complex handwritten recognition such as Arabic can be performed. We divide the recognition system into multilayered architecture inspired by human visual perception system. One example is discussed using the layered architecture and layer by layer process is shown in figure 6.7.

Level-0 acquiring the strokes information (x,y) from the input device. Level-1 is responsible for diacritical strokes segmentation and noise reduction. The four diacritical marks are separated from the basic shapes and its positional information is calculated.

Translation, rotation and position also effect at some level on human recognition, while it is very difficult to extract the translation and rotation invariant feature. We perform translational and rotational preprocessing to make the text translation and rotation invariant. Level 2 is responsible of baseline, slant, and mapping, etc. The simple geometrical features i.e. Loop, , Vertical Up, Vertical down, vertical left, Down left horizontal attached with horizontal left are extracted from the basic shape. With the help of linguistic rules, these simple features are combined to form the complex features i.e. loop may be downward, upward, oval, etc. Level 5 is responsible for forming the character from these local complex features and mapping for diacritical marks. Every diacritical mark is mapped onto the associated character. This mapping is based on fuzzy

triangular member function and linguistics rules are used to help the fuzzy member function in diacritical mapping.

Table 6.2: Recognition Rate

System Name	Classifier	Recog. Rate	Data Set
Multilanguage Arabic System (Proposed)	Biological Inspired Multilayered Fuzzy Logic	83.2	Unlimited Full Urdu, Arabic etc. Only Nasta'liq Style
Hybrid	HMM + Fuzzy	87.1	1500 Ligature only
OLUCR	BPNN	93	240 Ligature only
Online Urdu handwritten Recognition	Tree Based Dictionary Search	96	49 Ligature only

6.4 Final Remarks

This chapter presents multilayered fuzzy rules based expert system for handwritten Multilanguage character recognition i.e. Arabic script based languages inspired by human biological system. The biological visual perception is studied on 50 people, native readers of Arabic, Urdu, Punjabi, Sindhi and Pashto written in two different scripts Naskh and Nasta'liq to visualize the level of diacritical marks. The concept of ghost character recognition is built on this study and this concept helps to build biology inspired Multilanguage character recognition system. The character recognition system is divided

into layers according to human visual perception. The linguistics rules are added at each level to add the help of language behavior. The presented biological inspired provide 86.2% accuracy for Arabic, Urdu, Persian and Punjabi written in only Nasta'liq style. The result can be increase by adding the context clue at level 6.

CHAPTER 7

HOLISTIC APPROACH FOR URDU CHARACTER RECOGNITION USING MODIFIED HMM

Automatic recognition of cursive handwritten script remains a challenging problem even with the promising improvement in classifier and computational power. Segmentation based approach for recognition of handwritten Urdu script has considerable computational overhead and has lower accuracy as compared to Roman and Chinese script. Presence of complimentary characters in Urdu language makes it complicated as they have to be segmented into secondary strokes which are associated with a primary stroke. This first phase introduces a ligature based approach using Hidden Markov Model that provides solution for recognition of Urdu script. HMM database is divided into 54 subclasses based on the starting and ending shapes of the ligature. Twenty six time variant features have been selected for the base strokes. The sub division in classes reduces the time complexity and increases the efficiency. The second part of this chapter presents a segmentation free approach for recognition of Online Urdu handwritten script

using hybrid classifier, HMM and Fuzzy logic. Trained data set consisting of HMM's for each stroke is classified into 62 sub pattern based on the primary stroke shape at beginning and ending using fuzzy rule. Fuzzy linguistic variables based on language structure are used to model features and provide suitable results for large variation in handwritten strokes. The fuzzy classification into sub patterns increases the efficiency and decreases the computational complexity due to reduction in data set size.

7.1 Modified HMM Based Urdu Classification

After preprocessing phase, intended to simplify and reduce the noise from the input strokes, several structural and directional features are extracted and arranged in a sequence of score time. These features are further processed to remove unnecessary features from the extracted matrix based on the language structure with the help of language rule as described in chapter 5. Finally feature matrix is quantized to 32 discrete observation symbols and fed to HMM classifier to recognize the primary strokes. We have used the segmentation free approach i.e. ligature based approach in which the input stroke is not broken into subunits to avoid segmentation errors. Successful use of HMM in speech recognition led the researchers to use it for complex patterns especially for handwritten text. The segmentation free system extracts features for each ligature as described in chapter 5 that are passed on to the Hidden Markov Model for ligature classification. Using the strokes (x, y) co-ordinates and the chain codes, we extracted twenty six directional and structural features for primary strokes.

HMM based Arabic script handwritten recognition system illustrated in figure 7.1. We build separate HMM for each ghost ligatures each having 15 states. The HMM

database is divided into 54 classes based upon the structural shape of the primary stroke instead of size of the stroke. The additional timing information plays central role in recognition of primary stroke. First the ghost characters (basic shapes) are recognized (ligature based approach) and then secondary strokes associates with recognized primary strokes are recognized and then associated with primary strokes. We have used a grammar based approach to map the secondary strokes on to the associated recognized primary strokes to form valid ligature.

The task of classifier is to map the extracted features on to the trained data to find the best matched class. Hidden Markov Models (HMMs) are finite state machines and powerful statistical models for modeling the sequential or time-series data, and have been

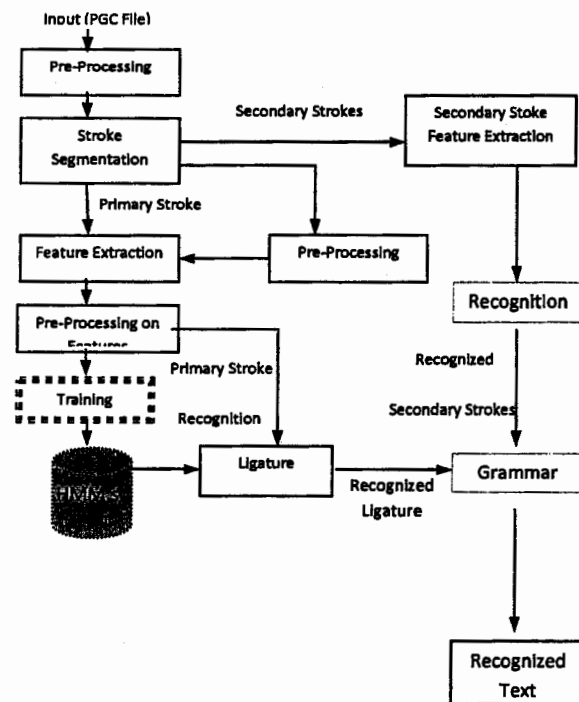


Figure 7.1. HMM based Recognition system architecture

Successfully used in many tasks such as speech recognition, OCR, information extraction and robotics where the visible observations and state transitions generating the observations form a doubly stochastic process with a sequential nature. HMMs are defined with a finite number of states where the features are assumed to be generated by a probability distribution depending on the current state. HMM can be differentiated by model unit, type of HMM (discrete or continuous) state meaning, topology, state duration modeling and its dimensionality. HMM's is dominant in the field of handwritten character recognition thus due to its success in handwritten character we proposed a solution for Urdu written in both Nasta'liq and Naskh style based on right to left, discrete with implicit state duration on time variant observation symbols arranged in sequential order. Handwriting strokes obtained through Tablet PC/ digital pen are segmented into primary and secondary strokes based on location, timing and size of stroke. HMM's used for the classification of primary strokes.

To increase the efficiency and decrease the computational complexity the primary strokes are divided into 54 classes based on the shape of the stroke instead of size of the stroke due to large variation in size of handwritten strokes. This classification is based on the directional and structural features at beginning and ending as shown in table 1. The sharp classification into 54 classes has disadvantage: if incorrect directional features are obtained, then definitely the wrong class will be selected and end result will be incorrect classification and advantage is that the computational complexities is decreased and the recognition result is more accurate due to less number of HMM's selected classes. Before forwarding the feature matrix to HMM classifier, lookup table is used to select the appropriate HMM classifier. The following is the HMM

N : Number of States

a_{ij} : Transition probabilities

b_{jk} : Observation probabilities

π_i : Expected frequency in state s_i at time $k=1$

$$\gamma_k(i) = \frac{P(q_k = s_i, o_1 o_2 \dots o_k)}{P(o_1 o_2 \dots o_k)} = \frac{\alpha_k(i) \beta_k(i)}{\sum_i \alpha_k(i) \beta_k(i)} \quad (6.1)$$

$$a_{ij} = \frac{\text{Expected number of transitions state } s_j \text{ to state } s_i}{\text{Expected number of transitions out of state } s_j} = \frac{\sum_k \zeta_k(i, j)}{\sum_k \zeta_k(i, j)} \quad (6.2)$$

$$b_{ij} = \frac{\text{Expected number of times observation } V_m \text{ occurs in state } s_j}{\text{Expected number of times in state } s_j} = \frac{\sum_k \zeta_k(i, j)}{\sum_{k, o_k = V_m} \zeta_k(i, j)} \quad (6.3)$$

$$\pi_i = (\text{Expected frequency in state } s_i \text{ at time } k = 1) = \gamma_1(i). \quad (6.4)$$

Observation sequence $O = o_1 o_2 \dots o_k$

Where $O_k \in FM \{\text{Loops, directional features, cusp, etc.}\}$

As we used the discrete HMM's therefore we have to map continuous features to discrete values and this mapping is done by vector quantization process. The code book of size 32 is used as shown in table 7.2. For example, the continuous data of loop is discrete based on width, height, position, size, and status of loop activation bit. If the loop is below the intersection line, it is considered as downward loop. If the loop is upward then size, width and height of internal side of loops are the important parameter to identify whether it is oval or round loop.

Table 7.1: Class division of strokes

Class	Features for division	Class	Features for division
1	Starting Loop and ending with hedge i.e. .	7	Starting with long down and ending with long up.
2	Starting with Jeem and ending with ray .	8	Ending with long left to right direction
3	Starting with Jeem and ending with hedge i.e. .	9	Starting with jeem and ending with long down
4	Starting with small down and ending with long up i.e. .	10	Starting and ending with long up. ,
5	Ending with ray	11	Starting with jeem and ending with left to right hedge
6	Starting with loop and ending with right to left direction	12	Starting with ayen and ending with long right to left direction .

The parameter estimation is concerned with training the HMMs through optimizing the model parameters. There is no known approach to solve the model parameter that maximizes probability of observation. However an iterative method such as the Baum–Welch, also known as expectation maximization (EM), is used to locally maximize the likelihood $P(O | M)$ of the chosen model $M=(A, B, \pi)$. The Baum-Welch training estimate the parameter $\lambda=(A,B, \pi)$ for each primary stroke that best describe the data.



Figure 7.2: Simple right to left HMM

Table 7.2: Code Book for HMM

Class	Description	Class	Description
1	Start vertical down	14	Interaction
2	Start vertical up	15	Ray\Dal
3	End vertical down	16	Loop -Up
4	End vertical up	17	Loop -Down
5	Diagonal right to left	18	Loop -Swad
6	Diagonal left to right	19	Hey
7	Horizontal right to left	20	Bay-Ray
8	Horizontal left to right	21	Bey-Ye
9	Hedge right to left	22	Ayen Shape
10	Hedge left to right	23	Circular Hey
11	Curve left to right	24	Tuan
12	Curve right to left	25	Shape of Madda Exist
13	Cusp	26-32	Eight chain code direction at start

The HMMs are estimated by giving some observation sequences $O = o_1 o_2 \dots o_k$ and general structure of HMM (numbers of hidden and visible states), determine HMM parameters $M = (A, B, \pi)$ that best fit training data, that maximizes $P(O | M)$. The observation sequence of each ligature is the sequentially occurrence of features. The important issue is to decide the number of states of each ligature. In our research the no of states are fixed for each ligature and we have used 15 states with skip state on simple right to left HMM.

The training data for HMM has been obtained by taking input from Urdu literate users for online Urdu handwritten script on the selected words. For training purpose, input data is acquired from 15 people so that more variations are obtained. HMM is built for each ghost character. The ghost character data set is divided into 54 classes depending on the starting and ending shape of the ghost ligature. Thus there are 54 HMMs classes, and each ligature HMM is placed in its relevant class depends upon the starting and ending shape of a character for example, "مت", "فپ" and "صب" are placed in the same class because they start with loop and ends with ب.

Recognition system is trained in two phases. In the first phase, database consist of 1500 ghost ligatures is obtained by taking the samples from 15 trained users. During data collection for training, baseline is considered by providing the interface and the user wrote on the provided baseline. In the second phase, secondary strokes are recognized and mapped onto associated primary strokes. For testing purpose, samples were taken from 15 skilled users. The testing is divided into four parts on classical HMM and HMM with sub divided database for both

Table 7.3: Recognition result in % on two data set A and data set B

HMM Structure	Input Type	Sample Data A	Sample Data B
Classical HMM	Primary Stroke	81.2	68.7
Classical HMM	Ligature with mapped dots	75.4	63.2%
HMM with Sub divided Database	Primary Stroke	85.3	78.2
HMM with Sub divided Database	Ligature with mapped dots	82.1	72.4

Where data set 'A' contain good strokes and data set 'B' contain poor stroke.

Table 7.4: Comparison with previous systems

System Name	Technology Used	Recognition Ligature	Dataset size
OLUCR	Back Propagation Neural Network	93%	240 Ligatures
Online Urdu Handwritten Recognition	Tree Based dictionary search	96.2	49 Ligatures
Proposed	HMM with Sub Classes	87.4 & 74	1500 Ligatures

Nasta'liq and Naskh font. Finally the results are also compared with other systems developed for Urdu. As most of the extracted features are font independent therefore it works for both fonts. The result in table 7.4 shows that HMM divided into sub classes (based on the shape of the ligature) improves the results. Several experiments have been performed to measure performance. Most of the errors are due to imperfect feature extraction that leads to incorrect class.

7.2 HMM and Fuzzy Logics: Hybrid Approach

We combined HMM with fuzzy logics for classification of handwritten Urdu script. This hybrid approach uses fuzzy rules as preprocessing of features to normalize the input and at post processing step to improve the results respectively. HMM is used as main classifier for the recognition of basic shape by putting it into fuzzy rules (as inner and outer shell) as shown in figure 7.5. Firstly, the preprocessing, feature extraction, feature

purification and clustering are performed through fuzzy rules as inner shell. Secondly the futures matrix is forward to HMM's for further classification of primary stroke. Finally this classified output is again purified through fuzzy rule outer shell by modeling the language concept on to the recognized output. We put circular shield around HMM to make it performance better than classical HMM. As the fuzzy logic is the powerful tool to classify the irregular and complex pattern. Language properties are observed deeply to increase the beauty of fuzzy logic in the recognition of Urdu script shown in figure 7.2. Fuzzy rule are used in three steps, preprocessing, feature extraction and at the end for post processing in classification correction and stroke mapping. Basically the proposed approach uses two classifier Hidden Markov Models (HMMs) and fuzzy rules in three layers. Fuzzy rules are used at layer 1 and layer 3 i.e. preprocessing and post processing whereas the HMM is used at layer 2 as main classifier. The output of layer 1 is fed to the HMM as the input of layer 2 for classification. Then classified result is again screened through fuzzy rules at layer 3. Finally mapping of diacritical marks on to the associated primary strokes by using fuzzy rules which are developed by critically analyzing the Urdu language structure. The whole process is described in 7. 2.

The classification phase is divided into three sub phases are

1. Phase I: Clustering
2. Phase II: Classifier I-HMM's
3. Phase III: Classifier II- Fuzzy Rules

Table 7.5: Vector Quantization

1	Starting Loop and ending with hedge i.e. ص بص .	7	Starting with long down and ending with long up. کا طا
2	Starting with Jeem and ending with ray بر جر	8	Ending with long left to right direction سے
3	Starting with Jeem and ending with hedge i.e. ج جی .	9	Starting with jeem and ending with long down جم
4	Starting with small down and ending with long up i.e. سہل چہل	10	Starting and ending with long up. ل , لا
5	Ending with ray ر , د , د	11	Starting with jeem and ending with left to right hedge ج ج
6	Starting with loop and ending with right to left direction ق	12	Starting with ayen and ending with long right to left direction ع ع .

7.2.1 Phase I: Clustering

The fuzzy inspired algorithm for classification into sub classes is based on the basic shape of the stroke guessed from the starting and ending of the stroke. Clustering is the pre-step for HMM classification of primary strokes and it is subdividing the strokes into subclasses according to their structure similarity. The division into subclasses can also be handled with the help of secondary strokes marks number, but here we are dealing with Multilanguage by treating diacritics separately and secondly number of dot may differ for the same shape in other languages like Arabic and Persian. The primary stroke is sub classified divided into 62 classes to increase the recognition rate and reduce the computational complexity.

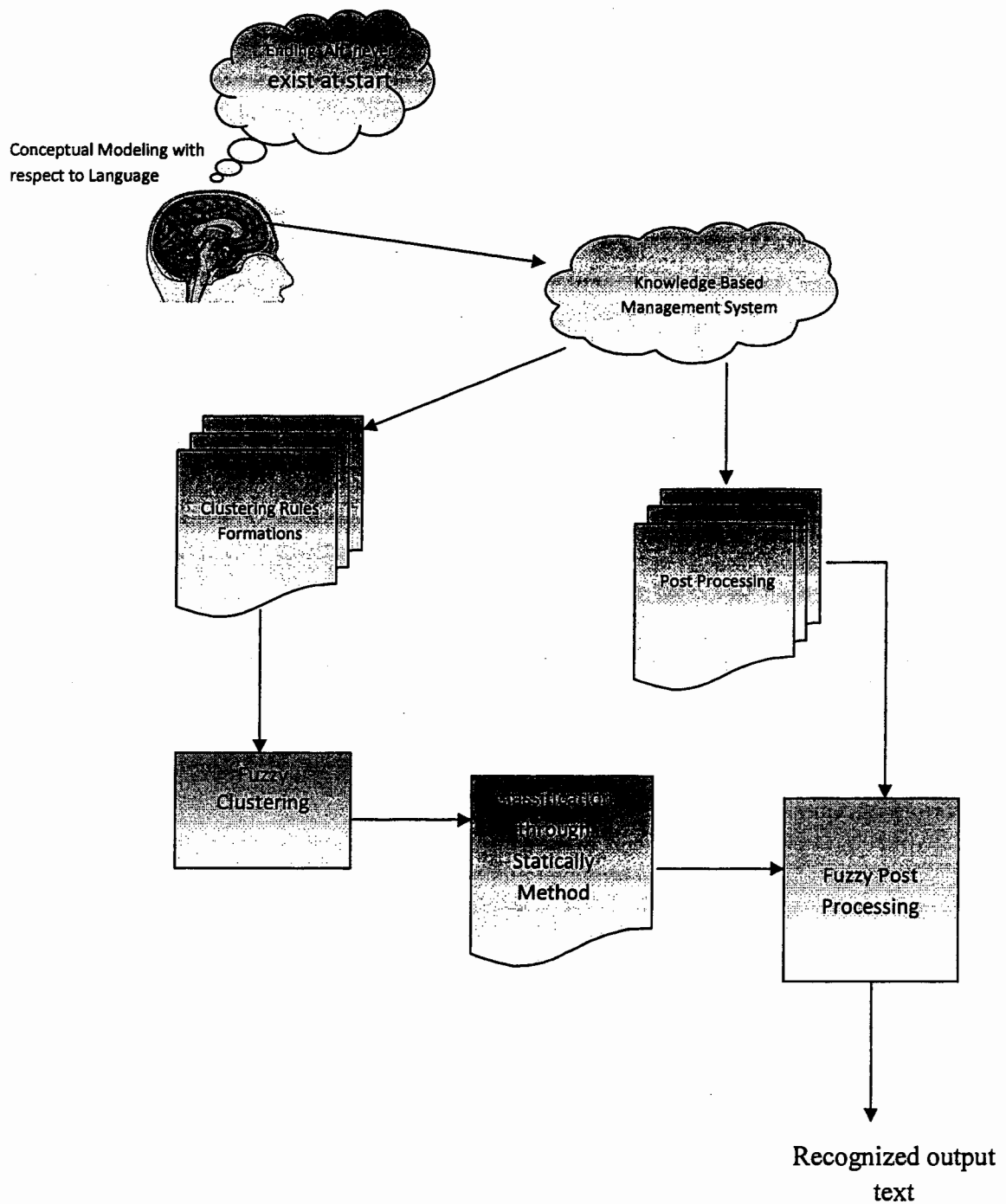


Figure 7.3: Fuzzy modeling on Statistical model with respect to language modeling.

The subdivision of strokes is based on the fuzzy rule, computed from the beginning pattern and ending pattern of the stroke i.e. class 1: بیم، بہجم، بانجم and class 2: $\text{خب، خت، خپ، خٹ، خٹ، چب}$. The variations in shape at end and starting (small, big, curve, loop, up, down) is computed by fuzzy triangular membership function and strokes are divided into classes based on the output of these fuzzy terms. For better clustering, maximum dependency attributes based on rough set can also be used. The variation at starting and ending are observed with respect to the height and width of the stroke.

7.2.2 Phase II: Classifier I-HMM's

The problem with fuzzy rule based method is absence of training and secondly it is impossible to form large set of rules that can be able to model all the possible set of shapes. The purpose of hybrid approach is to uses HMM's in folding of fuzzy rules method. Basically Fuzzy rules are applied before forwarding the input to HMM's i.e. classification into subclasses and again fuzzy inference rules are applied to purify the output obtained through HMM's. The extracted features are the set of basic shapes within strokes that uniquely define the structure of the stroke.

We used a right to left, no skip state and fixed state discrete Hidden Markov Models (HMM's) for primary stroke recognition and HMM's is built for each primary ligature. HMM database contains 62 sub databases whereas each sub database represents different basic shapes shown in figure 7.6. Quantization is required to convert the feature vector sequence into discrete symbols for observation sequence. Linguistically features FM_i i.e. loop; cups are quantized to the discrete value O_i . Eight directional features up to length γ are quantized from 0 to 8 and three kinds of loops upward, downward and oval are

discredited to 9, 10 and 11 respectively and upward and downward cups are quantized to 12 and 1 and similarly other feature matrix is quantized.

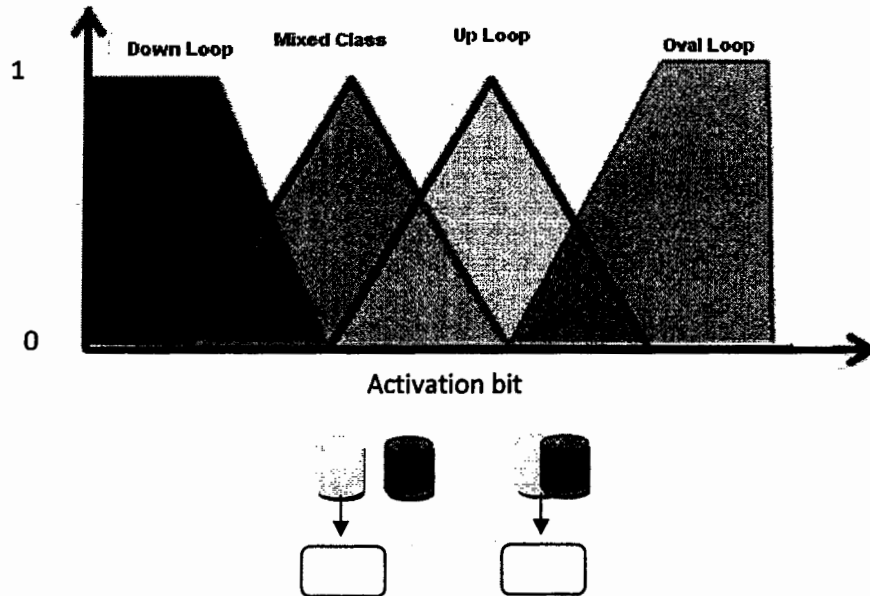


Figure 7.4: (a) Fuzzy member function for classification (b) Proper classification (c)

Mixed dataset classification

Figure 7.5 describe the layered hybrid structure member with the shell of fuzzy logic. The outer layer is responsible of classification into sub patterns. This sub patterns HMM database consist of small set of HMM's and used classification using HMM. This classification into sub pattern reduces the computation complexity due to less computation required for small number of classes. And increases the results due to small number of HMM's. Thus it is more convenient to find the best probabilistic match from this small dataset instead of large dataset.

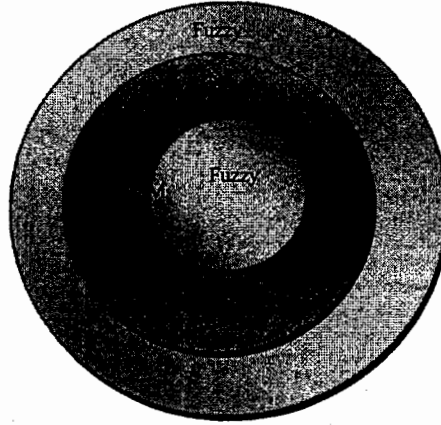


Figure 7.5: Layered Fuzzy with HMM

7.2.3 Phase III: Classifier II- Fuzzy Rules

In some cases HMM fails to recognize the very closely similar shapes. Thus we used fuzzy rule based method on the output probabilities of the HMM's to clarify the conflicted shapes better. In such cases rule based method is more powerful due to human reasoning involved in development of rules. Fuzzy rule-based method utilizes the linguistic variables. Fuzzy purification method is applied on to the output probabilities from the statistical method. If the difference between three maximum probabilities is less than α then Fuzzy inference is applied on to the probabilities. The probabilities of the strokes ا، ج are very close in some cases, thus we involved reasoning to resolve this conflict.

IF $P_m - P_{m-1} < \alpha \quad || \quad P_m - P_{m-2} < \alpha$

THEN USE Fuzzy Rules.

IF Strokes starts with diagonally right to left downward

Then it's

7.3 Post Processing

As secondary strokes are necessary to differentiate between the similar shapes in inter languages i.e. Arabic, Urdu and Persian, etc. and in intra language. Urdu script based languages contains zero or more secondary strokes correspond to one primary stroke. The mapping of secondary strokes on to the primary stroke is not easy task, because the dots may not appear exactly above or below the character in the stroke shown in figure 7.6.

The mapping of the secondary stroke is very critical especially when related and surrounded character are very small or related stroke may contain more than one stroke. During the segmentation process position of each stroke is computed and similarly position of each feature is also recorded during the feature extraction phase. To map the secondary stroke, dictionaries of possible valid words for all languages is used. The word dictionary contains 14150 words and is divided into sub dictionary based languages i.e. Urdu, Arabic, and Persian. The dictionary for Urdu consists of 10500 valid ligatures, whereas dictionary for Arabic and Persian contains 2200 and 1450 ligatures. The procedure of forming the valid ligature/ by mapping the diacritical marks is shown in figure 7.6. These valid ligatures are further used to form the valid word. The ligature dictionary contain 1800 ligature which are further used by word dictionary to form the valid words.

Table 7.6: Rule for Secondary strokes handling

Rule for Recognition of Secondary Strokes
<p>D_i is the directional feature where $i=1:8$ and D_1 the left to right side directional feature</p> <p>Rule 1: IF $SIZE(C(P_i)) < \gamma$ THEN Single Dot</p> <p>Rule 2: IF $SIZE(C(P_i)) < 3\gamma$ AND $C(P_i(x_l, y_l)) < C(P_i(x_f, y_f))$ THEN Double Dot</p> <p>Rule 3:: IF $C(P_i(x_l, y_l)) > C(P_i(x_f, y_f))$ AND D_6 at start D_6 at middle, D_6 at end THEN Kaf Line</p> <p>Rule 4: IF $C(P_i)$ contain loop then Tawn</p> <p>Rule5:: IF $C(P_i(x_l, y_l)) < C(P_i(x_f, y_f))$ AND D_8 at start D_8 at middle, D_8 at end THEN Kaf Line</p> <p>Rule 6: IF $C(P_i(x_l, y_l)) > C(P_i(x_f, y_f))$ AND $D_7 \parallel D_6$ at start D_5 at middle, $D_7 \parallel D_6$ at end THEN Madda.</p>

$$D_{Urdu} = \{\text{Valid set of Urdu words}\}$$

$$D_{Arabic} = \{\text{Valid set of Arabic words}\}$$

$$D_{Persian} = \{\text{Valid set of Persian words}\}$$

For all $SS_j \in PS_i$ Where $j=0$ to n and PS_i is the current primary stroke.

$$PS_i \{ \text{حب، حٹ، حپ، حٹ، جب، جٹ، جپ، جٹ، جٹ، جٹ} \}$$

$$\{ \text{خب، خٹ، خپ، خٹ، خٹ، جب، جٹ، جپ، جٹ، جٹ، جٹ} \}$$

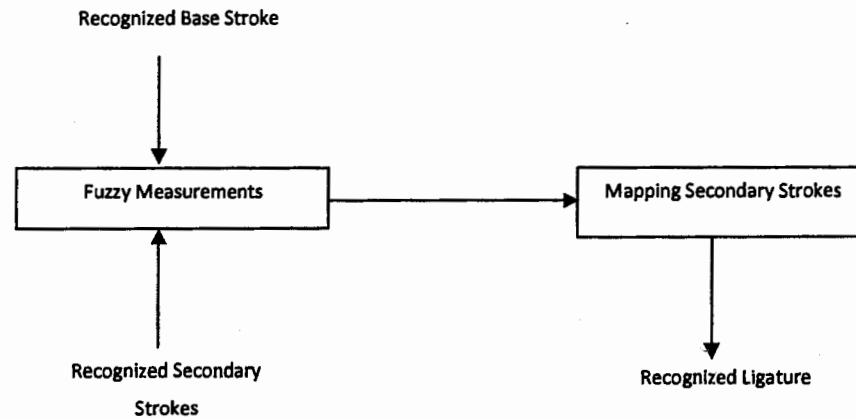


Figure 7.6: Stroke Mapping.

7.4 Final Remarks

For testing purpose we took Nasta'liq samples from 10 skilled users due to unavailability of Nasta'liq database whereas for Naskh style we used IFNENIT database for Arabic. The testing phase is divided into two parts i.e. Nasta'liq and Naskh and finally the results are also compared with previous work for Urdu script. As most of the extracted features are time variant and font independent thus it performs well for both Nasta'liq and Naskh. The result in table 5 shows that HMM divided into sub classes based on the basic ligature shape improved the results and furthermore the recognition rate is improved by using the Hybrid approach. The improvement in result using Hybrid approach is due to the involvement of human reasoning using fuzzy classifier in both features extraction and post processing phases. Several experiments are performed in order to measure this effect and average performance is increased by using hybrid approach. The most of the errors that was due to imperfect classification into subclasses is also reduced by the introduction

of fuzzy classification into sub data set. Moreover the system can recognize both Nasta'liq and Naskh font due independency of features on font. As the Nasta'liq style has more complex shapes than Naskh style [25] and Nasta'liq covers almost all shapes in Naskh style. Due to this reason the proposed system also works for Persian, Arabic and other languages written in Nasta'liq or Naskh style. We developed vocabulary of 10500, 2200 and 1450 words for Urdu, Arabic and Persian respectively shown in table 5. The computational complexity is very less as compared to conventional HMM due to the reduction in dataset for comparisons by defining the sub dataset. For example, we trained 1800 ligatures and there are 67 classes whereas class-I contain 30 ligatures. In conventional comparison is performed on 1800 HMM's whereas in proposed method only 30 HMM's are used as a trained dataset for comparison.

Table 7.7: Comparison of recognition rate WRT to Methodology

HMM Structure	Input Type	Sample Data A	Sample Data B
Classical HMM	Primary Stroke	81.2	68.7
Classical HMM	Ligature with mapped dots	75.4	63.2%
HMM with Sub divided Database	Primary Stroke	85.3	78.2
HMM with Sub divided Database	Ligature with mapped dots	81.1	70.4
Hybrid Classifier (HMM + Fuzzy Logics)	Primary Stroke	87.4	80.1
Hybrid Classifier	Ligature with mapped dots	87.3	78.6

Table 7.8: Comparison of recognition rate with previous System

System Name	Classifier Used	Recognition Ligature Nasta'liq and Naskh	Dataset size
OLUCR	Back Propagation Neural Network	93%	240 Ligatures
Online Urdu Handwritten Recognition	Tree Based dictionary search	96.2	49 Ligatures
Proposed Hybrid Approach	HMM + Fuzzy Logics	87.6 & 74.1	1800 Ligatures D _{Urdu} =9500, D _{Arabic} =1800, , D _{Persian} =1450

In this chapter we introduce a HMM based method for recognition of online Cursive handwritten Urdu Nasta'liq and Naskh Script. The system is currently trained for 1500 ligatures without segmenting the ligature into characters. The HMM's are divided into subclasses to improve the recognition rate and decrease the complexity. For testing purpose input was taken from 15 experienced users. We extracted 32 directional and structural features by critical examining the structure of the language and preprocessing on features matrix is applied to reduce the feature matrix from unnecessary features. We build right to left topology based HMMs for ghost characters which have 15 states and 32 observation symbols. The states are fixed and no skip state is allowed. The system provides 82.1% and 72.4% accuracy for both handwritten Nasta'liq and Naskh font respectively. K-mean clustering is used instead of HMM to recognize the secondary strokes and seven structural features are extracted. Finally the secondary strokes are mapped onto the primary stroke using position information through grammar. In future

we will perform clustering on dot is required for correct segmentation from the primary strokes.

CHAPTER 8

COMPARATIVE ANALYSIS AND CONCLUSION

This chapter analyzes the presented approaches i.e. ghost character theory, fuzzy preprocessing and feature extraction and finally compares the segmentation free approach with the biologically inspired (segmentation base) approach. The presented approaches are analyzed on both Nasta'liq and Naskh fonts. Some experiments were also performed for multi-language character recognition system. The fusion of fuzzy logic, fuzzy rules and human biological vision modeling provides a better solution to deal with complex character recognition problem. The fuzzy has the lack of learning and context knowledge whereas it provides flexibility in solving problems with large variation of data. A case study has been presented and has been divided into three parts. The first part presents the effect of ghost character theory for multi-language character recognition system and the second part analyzes the fuzzy based bio-inspired preprocessing and feature extraction

whereas bio-inspired segmentation based and segmentation free approaches have been analyzed in the third part. Finally, conclusion is drawn and future work is presented.

8.1. Ghost Character Based Recognition.

Normally both Naskh and Nasta'liq are followed by Arabic script based languages. Nasta'liq is mostly followed for Urdu, Persian, Punjabi, Sindhi, etc. whereas Naskh is mostly followed for Arabic, Pashto, etc. All Arabic script based languages are using almost similar writing script whereas they differ only by a few numbers of characters. Almost all ghost characters are the same in all Arabic script based languages. The mapping of diacritical marks and dictionary mapping is dependent upon the language selection. Each language has its own grammar, thus the ligature formation based on diacritical marks and word formation based on the ligatures which are according to the selected language. As every language has its own writing rules, ligatures and words but the basic shapes are the same therefore the recognition of basic shapes does not need any word formation rules, dictionary, etc. It is only dependent on the writing style used i.e. Nasta'liq or Naskh. The ligature formation from recognized ghost characters and recognized diacritical marks, and word formation from recognized ligatures required language modeling because the rules are dependent upon the language. The ghost character recognition theory is very successful for large datasets in the sense that Multilanguage recognition on Arabic script based languages can be developed by using rules of each language. Figure 8.1 shows the result of simple ligature formation while figure 8.4 shows the results for complex ligature formation for different languages.

The following section describes the dictionary formation of different Arabic script based language and it is evident that the number of ligatures is substantially reduced by using the Ghost character theory.

Dictionary D= (Urdu, Arabic, Persian, Punjabi, Pashto, Sindhi)

Ligature Dictionary for Urdu = [$L_1 \{ \dots \}, L_2 \{ \dots \}$

$L_1 \{ \text{حَب، حَت، حَب، حَت، حَب، حَت، حَب، حَت، حَب، حَت، حَب، حَت} \}$

$\{ \text{خَب، خَت، خَب، خَت، خَب، خَت، خَب، خَت، خَب، خَت، خَب، خَت} \}$

$\dots \dots \dots L_n \{ \dots \}]$

Ligature Dictionary for Arabic = [$L_1 \{ \dots \}, L_2 \{ \dots \}$

$L_1 \{ \text{جَ، جِ، جُ، حَ، حِ، حُ، خَ، خِ، خُ} \}$

$\dots \dots \dots L_m \{ \dots \}]$

The mapping of diacritical marks with respect to various languages dictionary with the same ghost ligature and same no of diacritical marks is shown in figure 8.1.

The major benefit of the proposed ghost character recognition theory is that the developed system works for all Arabic script based languages and the overall ligature set is decreased.

$\text{Ligature}_{\text{Multilanguage}} = \text{No of total ligatures by Arabic script based languages}$

$\text{Ligature}_{\text{Arabic}} = \text{No of total ligatures of Arabic}$

$\text{Ligature}_{\text{Urdu}} = \text{No of total ligatures of Urdu}$

$\text{Ligature}_{\text{Persian}} = \text{No of total ligatures of Persian}$

$\text{Ligature}_{\text{Punjabi}} = \text{No of total ligatures of Punjabi}$

$\text{Ligature}_{\text{other Arabic script based languages}} = \text{No of total ligatures of other Arabic script based languages like Pashto, Sindhi, etc.}$

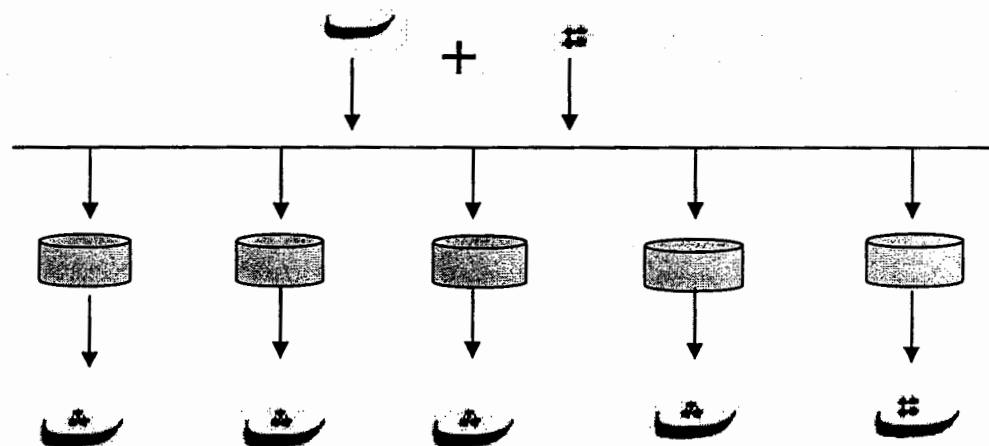


Figure 8.1(a): Combination of diacritical marks with respect to languages

Figure 8.2 shows the 2nd character for all languages. The basic shape is same for all language while difference in diacritical marks. From these 11 characters, 5 belong to Urdu and 3 belong to Arabic.

As we know

$\text{Ligature}_{\text{Ghost ligature of Urdu}} \lll \text{Ligature}_{\text{Urdu}}$

Thus we can say,

$\text{Ligature}_{\text{Multilanguage (Ghost)}} \sim < \text{Ligature}_{\text{Arabic}} + \text{Ligature}_{\text{Urdu}} + \text{Ligature}_{\text{Persian}} + \text{Ligature}_{\text{Punjabi}} + \text{Ligature}_{\text{other Arabic script based languages}}$

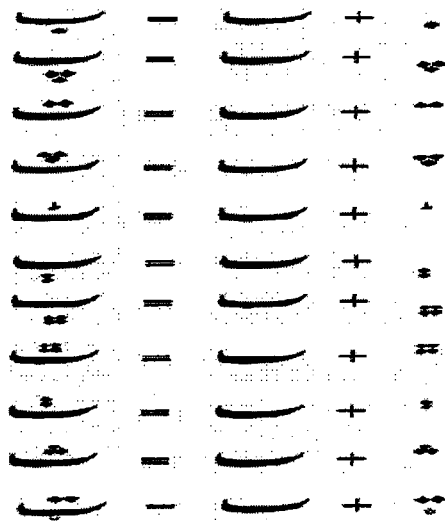


Figure 8.1(b): 'Bey' character in all languages

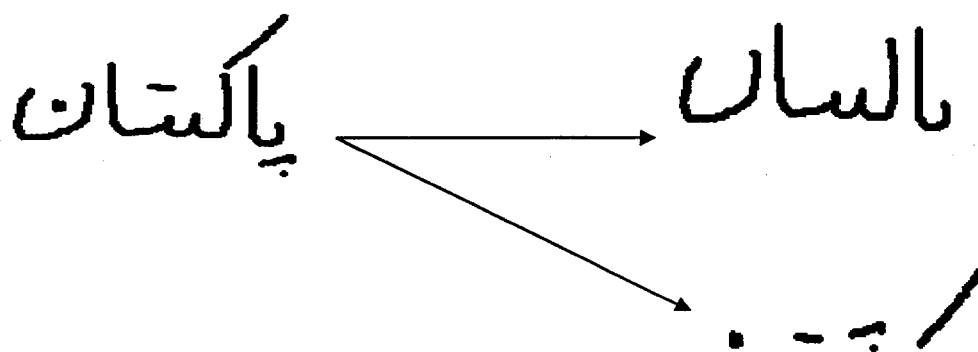


Figure 8.2 Combinations of diacritical marks with respect to languages

8.2 Feature Extraction and Fusion.

As successful character recognition methodology depends upon the particular choice of features used by the classifier. Thus it is the most critical part in any recognition problem. Due to the complexity and variation, direct recognition of handwritten stroke is almost impossible. The literature shows that directional and structural features i.e. loop, cusp,

etc. played important role in the classification. Fuzzy rules have been used to extract the unique and meaningful directional and structural features and shape defining patterns i.e. loops, cusp, endpoints, starts points, etc. Further post processing is also applied on extracted features to remove the unnecessary and noisy features by modeling the language rules and fusion of directional and structural features. Where and what information of each salient pattern is extracted and used for feature fusion. The features are two close points are combined to form new features. The approach was based on both directional and structural features. We fused both directional and structural features and it provides accuracy.

Fusion of structural and directional features is performed based on relative position information of each feature. Z is the feature and x,y are the feature coordinates. Where and what information of each salient pattern is extracted and used for feature fusion. The features are two close points are combined to form new features based on the language rules. For example the directional feature `Start_Vertical_Down` is combined with structural feature `Loop` to form new feature for ligature ط. This new feature is more discriminant as compared to the previous two features because it utilizes the location information as well in the construction of new feature matrix. The new feature is more concise and discriminant.

8.3. Segmentation Free and Segmentation Based Approach

Character recognition of Arabic script based languages has some technical issues not existing in other languages; makes these languages more complicated and hence the low recognition accuracy. Segmentation based approach for handwritten Arabic script based

languages incorporates a considerable overhead and have very less accuracy. The analysis of the Arabic script is complicated due to segmentation problem of cursive script and use of complimentary characters. We have presented both segmentation and segmentation free approaches. In segmentation free approach, HMM is trained for each ghost ligature. A Fuzzy logic is used in the classifier while fuzzy rules have been used at post processing step to improve the results. The division in classes reduced the time complexity and increased the efficiency. Generally the system works for both handwritten Nasta'liq and Naskh fonts of Arabic script and provided 87.4% and 80.1% classification result for the two fonts respectively tested on 1800 ligatures obtained from 15 trained users. The ligature based approach is limited due to the training data used and required a considerable computational overhead during training. We have tested the segmentation free system on 10500, 2200 and 1450 valid ligatures for Urdu, Arabic and Persian respectively as shown in table 5.

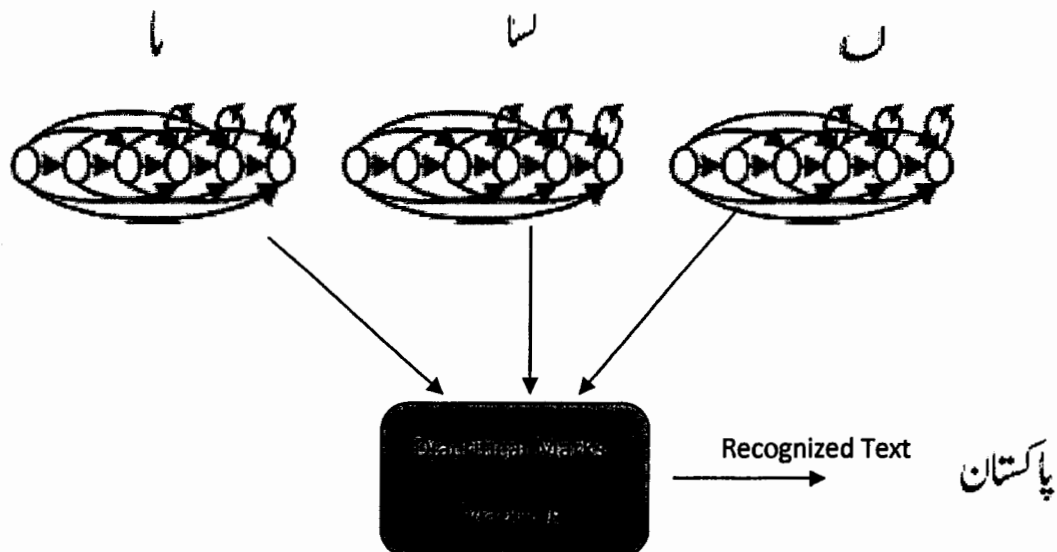


Figure 8.3: Combination of diacritical marks with respect to languages

The computational complexity during classification is less as compared to conventional HMM due to the reduction in dataset for comparisons. For example, we trained 1800 ligatures and there are 67 classes whereas class-I contain 30 ligatures. In conventional HMM, classification is performed by single HMM for the 1800 ligatures whereas in proposed method only 30 ligatures. The ligature based approach is suitable for limited words i.e. city name, verbs, etc.

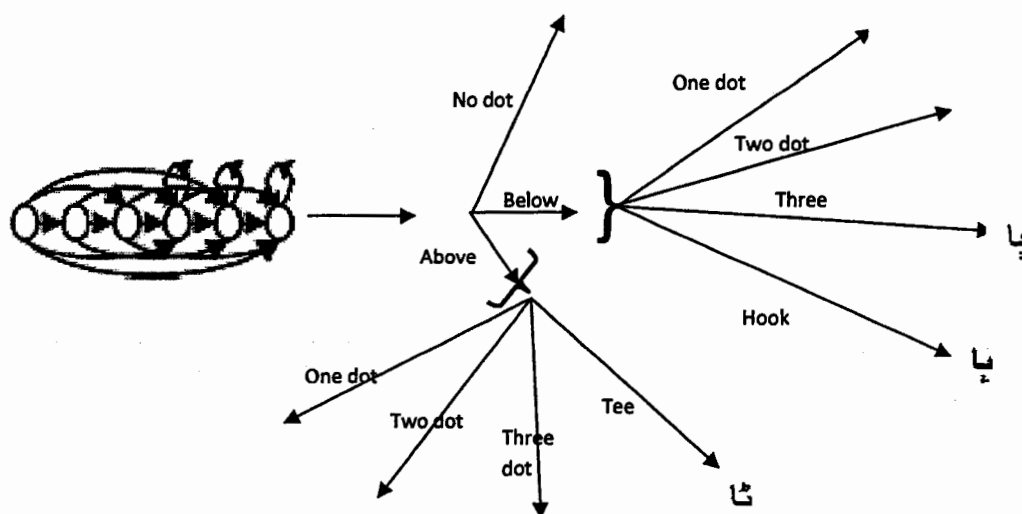


Figure 8.4: Combination of diacritical marks with respect to languages

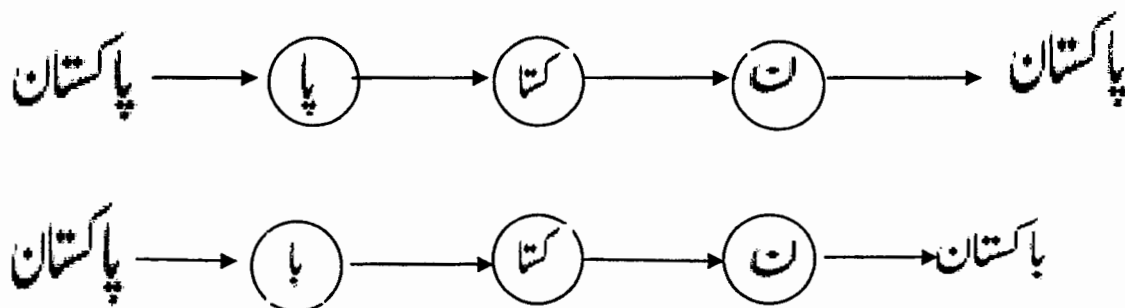


Figure 8.5: (a) Ligature formation by combination of diacritical marks and ghost character (b) Word level recognition HMM for same word in Urdu and Arabic.

Level-0

پاکستان

Level-1 Diacritical Marks
Separation, Noise
Reduction

پاکستان

Level-2: Baseline, Slant
etc. estimation

پاکستان

Level-3: Simple
Feature Extraction

پاکستان

Level-4: complex Feature
Formation

پاکستان

Level-5: Character
Level Recognition

پاکستان

Forming the character and mapping the diacritical marks

Level-6: Context Level
Computation

پاکستان

Figure 8.6: Bio-inspired character recognition process of word Pakistan

The segmentation based approach is difficult and less accurate due to the complexity of the Arabic script especially the text written in Nasta'liq style. On the other segment based

approach can treat the full script. We presented biologically inspired segmentation based approach in which the segmentation and recognition is done simultaneously using

Table 8.1: Comparison of Recognition Result of different system

HMM Structure	Input Type	Sample Data A	Sample Data B	Data Size
Classical HMM	Primary Stroke	81.2	68.7	1000 Ligature
Classical HMM	Ligature with mapped dots	75.4	63.2%	1000 Ligature
HMM with Sub divided Database	Primary Stroke	85.3	78.2	1500 Ligature
HMM with Sub divided Database	Ligature with mapped dots	81.1	70.4	1500 Ligature
Hybrid Classifier (HMM + Fuzzy Logics)	Primary Stroke	87.4	76.1	1800 Ligatures
Hybrid Classifier	Ligature with mapped dots	87.6	74.6	1800 Ligatures
Bio-Inspired Character Recog.	Primary Strokes	86.2	77.4	Full Urdu Text
Bio-Inspired Character Recog.	Ligatures with Diacritical marks	85.2	73.9	Full Urdu Text
Bio-Inspired for Word Recog.	Word	89.5	75.1	2000 Words

the fuzzy logics. The uncertainties in character recognition system have different effect depending on the type of problem whereas it is difficult to classify the complex problem into different level of uncertainty. Vast variation and inconsistencies makes handwritten Arabic script based languages character recognition much more complex than any other language. Thus a careful observation preformed and the complex handwritten recognition is divide into small uncertainty levels as similar the human ventral stream divide the handwritten problem into different uncertainty levels. The layered fuzzy approach is used

that lead to recognized results by deducing the high end patterns from small patterns. These high end patterns are used for the creation of further high end pattern. The integration human biological vision with fuzzy logic provided speed, accuracy and efficiency.

The biological visual perception is studied on 50 people, native readers of Arabic, Urdu, Punjabi, Sindhi and Pashto written in two different scripts Naskh and Nasta'liq to visualize the recognition level of diacritical marks. The concept of ghost character recognition is built on this study and this concept helps to build biology inspired Multilanguage character recognition system. The character recognition system is divided into layers according to human visual perception. The linguistics rules are added at each level to add the help of language behavior. The presented biological inspired provide 86.2% accuracy for Arabic, Urdu, Persian and Punjabi written in only Nasta'liq style. The result can be increase by adding the context clue at level 6.

8.4 Conclusion

Although the segmentation free approach provided good results as compared to segmentation based approach but the segmentation free approach is suitable for small data set i.e. cities name, country name, verbs, etc. Whereas the segmentation based approach can work for full Urdu text. Thus multi-language system is developed using the ghost character recognition theory and language rules by extracting the diacritical marks from the handwritten text before recognition and their association to the recognized base shape. Using this concept, our proposed system works successfully for multiple languages based on Arabic character script. The bio-inspired fuzzy approach is tested for multiple languages written in both Nasta'liq and Naskh and it provides 86.2%, 74.1% accuracy respectively on input from ten trained user. The main focus of the thesis was the fuzzy based bio-inspired based character recognition by processing each step with the help of fuzzy logic to add human reasoning.

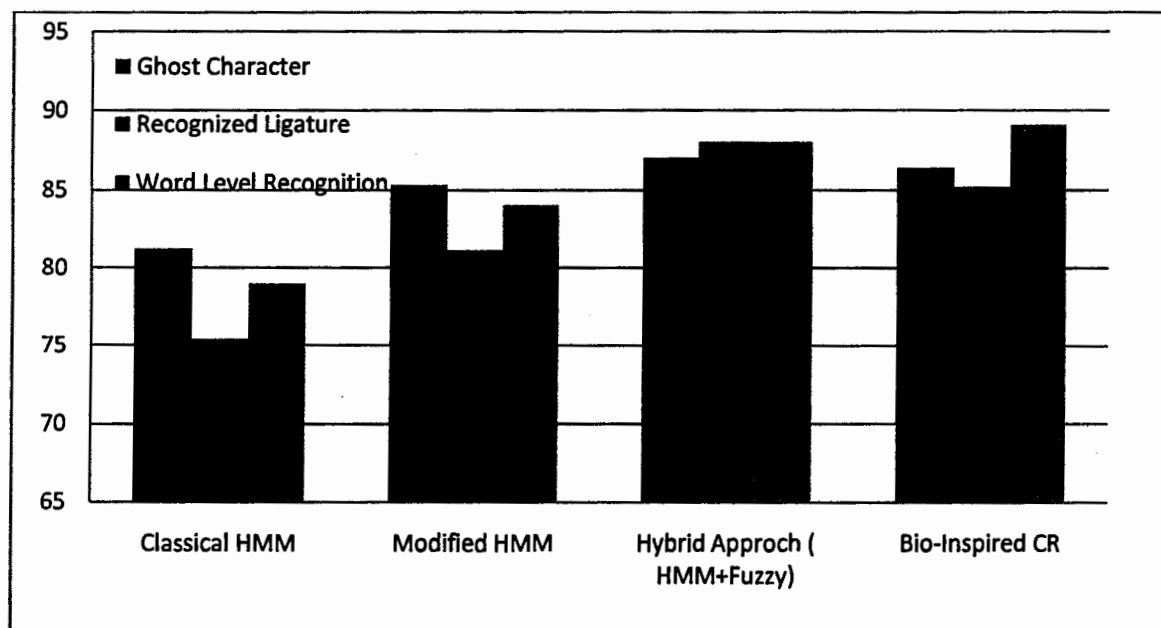


Figure 8.7. Recognition rate

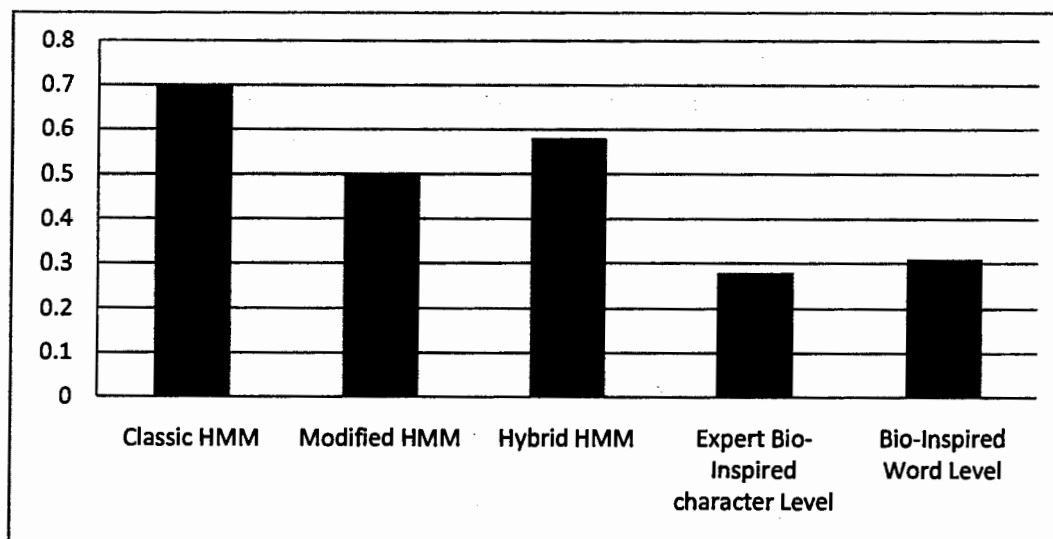


Figure 8.8. Computational Complexity

Table 8.2: Comparison with existing systems

Online Character Recognizer	Approach	Classifier	Data Set	Writing Style	Recog. Rate
Biologically Inspired Fuzzy Based Expert System	Simultaneous Segmentation	Fuzzy Logic Fuzzy Rules Bio-Inspired Methodology	Full Data set (Signature Level)	Nasta'liq Naskh	86.2
Hybrid Approach HMM and Fuzzy Logics [11]	Ligature Based Approach	HMM and Fuzzy Logics	14150 Words	Nasta'liq Naskh	87.6
Arabic Character Recognition		Genetic Algorithm And Visual Encoding	200 Words	Naskh	95
Online Urdu Character Recognition [14]	Ligature Based Approach	Backpropagation Neural Network	240 Ligatures	Nasta'liq	93
Online Arabic Handwritten recog With templates [16]		Template Based Matching	Full Data set	Naskh	68.2
Bio-Inspired Handwritten Farsi [9]		KNN ANN SVM	Farsi Digits (MNIST)	Naskh	81.3 (KNN) 94.65 (ANN) 98.75 (SVM)
Arabic Handwritten using Matching Algorithm [13]	Ligature Based	Matching Algorithm and Decision Tree	Arabic Alphabets	Naskh	98.3
Recognition Based Segmentation of Arabic [20]	Recognition based Simultaneous	HMM	Arabic	Naskh	88.8
Rule Based Urdu Online Character Recognition [21]	Simultaneous Segmentation based Recognition	Fuzzy Logics	Urdu Full Dataset	Nasta'liq and Naskh	78
Online Arabic Characters by Gibbs Modeling [23]		Gibbs Modeling of Class Conditional Densities	Arabic Characters	Naskh	84.85 (direct Bayes) 90.19 (Indirect Bayes)
Arabic Character Recognition using HMM [24]		HMM	Arabic full data set	Naskh	78.25

8.5 Future Work

The presented work is not focused on post processing phase. In future, we will focus on post processing steps (level-6) and add the context knowledge. For example the presented work has lack of word level dictionary, In future; work can be carried out on dictionary to improve the recognition rate as a post processing step.

Secondly, probabilistic based multi solution from lower level to final level can be tried instead of exact solution during the feature extraction and classification as shown in figure 8.9.

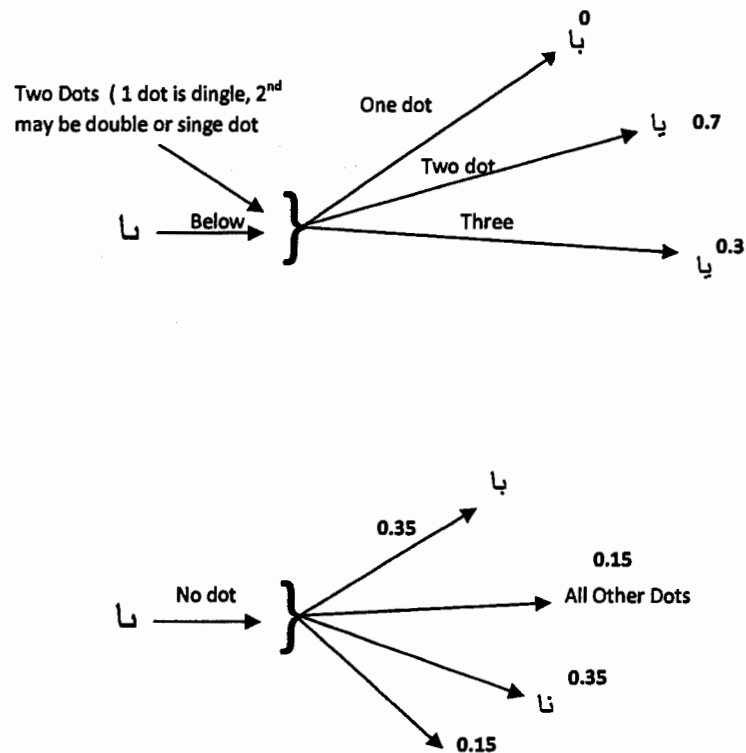


Figure 8.9: Probabilistic based modeling for future work

We will extract few close solutions and then these close solutions can be further modeled in layered manner. Some work on this is performed in chapter 6 by using hybrid approach: Fuzzy and HMM, the output from HMM can be further analyzed using fuzzy logic but in this case only one solution can be considered whereas in future we will consider several top solutions and these solution will be further reduced to one solution by performing the probabilistic based modeling till word formation. For example the diacritical mapping is discussed in figure 8.7.

REFERENCES

- Abuhaiba, M.J.J. Holt, and S. Datta, (1998) "Recognition of Off-Line Cursive Handwriting," Computer Vision and Image Understanding, vol. 71, pp. 19-38.
- Adeed S.A, Higgins C., Elliman D. (2002), " Recognition of Offline Handwritten Arabic Word Using Hidden Markov Model Approach" , 16th International Conference on Pattern Recognition (ICPR'02) - Volume 3
- Ahmad Z., Orakzai J.K., Shamsher I., Adnan A, (2007) "Urdu Nastaleeq Optical Character Recognition" World Academy of Science, Engineering and Technology 32, 2007.
- Ahn J.H., Lee J, Jo J., Choi Y, Lee Y, (2009)"Online Character Recognition using Elastic Curvature Matching", Seventh International Conference on Advances in Pattern Recognition.
- Al-Absi H.R.A., Abdullah A.B. (2009) "A Proposed Biologically Inspired Model for Object Recognition" IVIC 2009, LNCS 5857, pp. 213-222,
- Al-Badr. B., Mahmoud, S.A. (1995), "Survey and Bibliography of Arabic Optical Text Recognition", Elsevier Science, Signal Processing, Vol. 41, pp. 49-77.
- Al-Badr B. and R. Haralick, (1998) "A Segmentation-Free Approach to Text Recognition with Application to Arabic Text," International Journal Document Analysis and Recognition, vol. 1, pp. 147-166,.
- Al-Badr B. and R. Haralick, (1995) "Segmentation-Free Word Recognition with Application to Arabic," Proc. International Conference Document Analysis and Recognition, pp. 355-359.
- Al-Ghoneim K. (2001) "Sub Stroke Approach to HMM-based On-line Kanji Handwriting Recognition" , International Conference on Document Analysis and Recognition, 2001
- Al-Habian, G. Assaleh, K., (2007) "OnlineArabic handwriting recognition using continuous Gaussian mixture HMMS " International Conference on Intelligent and Advanced Systems, 2007. .

Ramy A.H.M., Sulem L.L., Mokbel C. (2009), Combining slanted-frame classifiers for improved HMM based Arabic and writing recognition, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, Volume: 31 Issue:7

Al-Hamad H.A., Zitar R. A. (2010) "Development of an efficient neural-based segmentation technique for Arabic handwriting recognition" *Pattern Recognition* 43 2773–2798

Al-Khateeb J. H, Ren J., Ipson S.S, Jiang J. (2008), "Knowledge-based Baseline Detection and Optimal Thresholding for Words Segmentation in Efficient Pre-processing of Handwritten Arabic Text", *Fifth International Conference on Information Technology: New Generations*, 2008

AlKhateeb J.H., Jinchang Ren, Jianmin Jiang, Stan S Ipson, (2009), "Unconstrained Arabic Handwritten Word Feature Extraction: A Comparative Study", *Sixth International Conference on Information Technology: New Generations*,

Al-Rashaideh H., (2006) "Preprocessing Phase for Arabic Word Handwritten Recognition", *Information Transmission in Computer Networks Vol 6. No.1*, pp 11-19 .

Adeed S.A, Elliman D., Higgins C., "A Data Base for Arabic Handwritten Text Recognition Research", *The International Arab Journal of Information Technology*, Vol. 1, No. 1, January 2004

Alma'adeed, S., D. Elliman, and C.A. Higgins (2002), "A Data Base for Arabic Handwritten Text Recognition Research," *Proceeding Eighth International Workshop Frontiers in Handwriting Recognition*, pp. 485-489,.

Alama S.' adeed (2008). , Recognition of Offline Handwritten Arabic Word Using Hidden Markov Model Approach" , 16th International Conference on Pattern Recognition (ICPR'02) - Volume 3.

Al-Omari FA, Al-Jarrah O (2004) Handwritten Indian numerals recognition system using probabilistic neural networks. *Adv Eng Inform* 18(1):9–16. doi:10.1016/j.aei.2004.02.001.

Alonso J., W. Usrey and R. Reid (1996), Precisely correlated firing in cells of the lateral geniculate nucleus. *Nature*, 1996. 383: 815-819.

Aly W., Uchida S, Suzuki M (2007), "Identifying Subscripts and Superscripts in Mathematical Documents", *Mathematics in Computer Science*, 2008, pp 195-209

Amor N.B., Amara N.E.B (2006), "Multifont Arabic Characters Recognition Using Hough Transform and HMM/ANN Classification", *Journal of Multimedia*, Vol. 1, NO. 2, May 2006

Amara N. Ben and Bouslama F. (2003). Classification of Arabic script using Multiple sources of information: state of the art and perspective," *International Journal of Document Analysis and Recognition*. vol. 5. pp. 195-212.

- Anastas, P and Warner, J (2000) Green chemistry: theory and practice, Oxford University Press
- Anwar W, Wang X, Wang X.L. (2006), "A survey of automatic Urdu language processing" Proceedings of the Fifth International Conference on Machine Learning and Cybernetics, Dalian, 13-16 August 2006
- Baghshah M.S., Baghshah S.B., Kasaei S. (2006) "A novel fuzzy classifier using fuzzy LVQ to recognize online Persian handwriting", Information and Communication Technologies, 2006. ICTTA '06
- Baker, C.I., Liu, J., Wald, L.L., Kwong, K.K., Benner, T., Kanwisher, N., (2007). Visual word processing and experiential origins of functional selectivity in human extrastriate cortex. *Proc. Natl. Acad. Sci. U. S. A.* 104, 9087–9092
- Benouareth A., A.Ennaji, M. Sellami, (2008)"Arabic Handwritten Word Recognition Using HMMs with Explicit state Duration" *EURASIP Journal of Advance Signal Processing*, Vol. 2008
- Benouareth A, Ennaji A, Sellam M. (2008)"Semi-Continuous HMMs with Explicit State Duration Applied to Arabic Handwritten Word Recognition , *Pattern Recognition Letter* pp 1742-1752
- Biadisy F., J. El-Sana, N. Habash (2006), Online Arabic handwriting recognition using hidden Markov models, in: *Proceedings of the Tenth International Workshop on Frontiers in Handwriting Recognition*, pp. 85–90.
- Biederman I. (1987), "Recognition-by-components: A theory of human image Understanding" *Psychological Review*, 24: 115-147.
- Borji A., Hamdi M, F Mahmoudi, (2008), "Robust Handwritten Character Recognition with Features Inspired by Visual Ventral Stream" *Neural Processing Letters* Vol 28 pp 97-111.
- Bosner, R,(2006) "Patented biologically-inspired technological innovations:a twenty year view", *Journal of Bionic Engineering*, Vol 3 pp 39e41
- Bosner, R and Vincent, J Technology trajectories, innovation, and the growth of biomimetics, *Journal of Mechanical Engineering Science* ,2006, Vol 221 pp 1177-1180
- Boubaker H., M. Kherallah, A. M. Alimi, New Algorithm of Straight or Curved Baseline Detection for Short Arabic Handwritten writing, 10th International Conference on Document Analysis and Recognition.
- Bouchareb F., M Bedda, S. Ouchetati, (2006) "New Preprocessing Method for Handwritten Arabic Word", *Asian Journal of Information Technology* 2006, pp 609-613.

- Chan S, Tang S.W., Tang K.W., Lee W.K., Lo S.S, Kwong K.K. (2009), "Hierarchical coding of characters in the ventral and dorsal visual streams of Chinese language processing", *NeuroImage* 48 pp 423–435
- Chen, M.Y., Kundu, A. and S. N. Srihari, (1995) "Variable duration hidden Markov model and morphological segmentation for handwritten word recognition," *IEEE Transactions on Image Processing*, vol. 4, no. 12, pp. 1675–1688, 1995.
- Collins, M and Brebbia (2001), *C Design and nature II: comparing design in nature with science and engineering* Wessex, 2001,) Institute of Technology Press
- Deepu V., Sriganesh M. Ramakrishnan A. G. (2004), "Principal Component Analysis for Online Handwritten Character Recognition" 17th International Conference on (ICPR'04) Volume 2 - Volume 02.
- Dehghani A., F. Shabani, and P. Nava, (2001) "Off-Line Recognition of Isolated Persian Handwritten Characters Using Multiple Hidden Markov Models," *Proceeding International Conference Information Technology: Coding and Computing*, pp. 506-510, 2001.
- Deng W, Hu J, Guo J., Cai, W., Feng, D. (2010) "Emulating biological strategies for uncontrolled face recognition" *Pattern Recognition* 43 (2010) 2210–2223
- Durani A (2009), "Pakistani: Lingual Aspect of National Integration of Pakistan", www.nlait.gov.pk.
- Durani A (2008), "Urdu Informatics" Vol. 1, pp. 102-112, pp 8-15, National Language Authority Press
- Dong L, Ebroul Izquierdo, (2007) "A Biologically Inspired System for Classification of Natural Images", *IEEE transactions on circuits and systems for video technology*, vol. 17, no. 5, may 2007
- Douglas D. , Peucker T. (1973), Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. *The Canadian Cartographer*, 10(2):112–122, 1973.
- Eberhart, R., Simpson P., Dobbins R., "Computational Intelligence PC Tools" Ed. AP Professional.
- El-Hajj R., L. Likforman-Sulem, and C. Mokbel, (2005) "Arabic Handwriting Recognition Using Baseline Dependant Features and Hidden Markov Modeling," *Proceeding International Conference Document Analysis and Recognition*, pp. 893-897, 2005.
- Daifallah K., Zarka N, Jamous H., Recognition-Based Segmentation Algorithm for On-Line Arabic Handwriting, 2009 10th International Conference on Document Analysis and Recognition.

- Dong Le and Ebroul Izquierdo, A Biologically Inspired System for classification of Natural Images, IEEE Transactions on Circuits and System for Video and Technology, Vol. 17, No. 5, May 2007
- Douglas D. , Peucker T. Algorithms for the reduction of the number of points required to represent a digitized line or its caricature. The Canadian Cartographer, 10(2):112–122, 1973.
- Duan J, Li R, Hu Y, “A bio-inspired application of natural language processing: A case study in extracting multiword expression” Expert Systems with Applications 36 (2009) 4876–4883.
- Durani A., "Urdu Informatics" Vol. 1, pp. 102-112, pp 8-15, National Language Authority Press 2008.
- Durani A (2009), "Pakistani: Lingual Aspect of National Integration of Pakistan", www.nlait.gov.pk.
- Elanwar R. I., Simultaneous Segmentation and Recognition of Arabic Characters n an Unconstrained On-Line Cursive Handwritten Document International Journal of Computer and Information Science and Engineering Volume Number 4.
- Farooq F., V. Govindaraju, and M. Perrone,(2005) "Preprocessing methods for handwritten Arabic documents".ICDAR (Proceedings of the Eighth International Conference on Document Analysis and Recognition), pp: 267-271, 2005.
- Fahmy M..M. and S. Al Ali, (2001)“Automatic Recognition of Handwritten Arabic Characters Using Their Geometrical Features,” Studies in Informatics and Control Journal., vol. 10, 2001.
- Faradji F., Faez K., Nosrati M.S. (2007), “Online Farsi handwritten words recognition using a combination of 3 cascaded RBF veural networks” Internatioanl Conference on Intelligent and Advance System 2007
- Faradji F., Faez K., Nosrati M.S. (2007), “An HMM-based online recognition system for Farsi handwritten words” International Conference on Intelligent and Advance System 2007
- Farghaly A, Shaalan K, (2009)“Arabic natural language processing: challenge and solution” ACM Transactions on Asian Language Information Processing, Vol. 8, No. 4, 2009.
- Forbes, P The Gecko’s foot: bio-inspiration, (2005) engineering new materials and devices from nature Harper Collins
- Fujiswa, H, (2008)“Forty years of research in character and document recognition---an industrial perspective” Pattern Recognition 41 (2008), pp 2435 – 2446.

- Fujisawa H. and C.-L. Liu. (2003) Directional pattern matching for character recognition revisited. Proceedings of the 7th International Conference on Document Analysis and Recognition, Edinburgh, Scotland, 2003, pp. 794–798.
- Gillies A. M., (1992), “Cursive word recognition using hidden Markov models,” in Proceedings of the 5th U.S. Postal Service Advanced Technology Conference, pp. 557–562, Washington, DC, USA, November 1992
- Goodale M. and A. Milner (1992), Separate visual pathways for perception and action. Trends in Neurosciences, 1992. 15: 20-25
- Govindan, V. K., Shivaprasad, A. P. (1990), “Character Recognition, A review”, Pattern Recognition, Vol. 23, No. 7, pp. 671–683, 1990.
- Haghighi P. J., Nobile N., He C.L. Suen C.Y. , “A New Large-Scale Multi-purpose Handwritten Farsi Database” Lecture Notes in Computer Science, 2009, Volume 5627/2009, pp 278-286.
- Halavati R., M. Jamzad, M. Soleymani(2005), A novel approach to Persian online hand writing recognition, Transactions on Engineering, Computing and Technology 6 232–236.
- Hamada R. H. A. and A. B. Abdullah (2009) A Proposed Biologically Inspired Model for Object Recognition, Lecture Notes in Computer Science, Volume 5857, pp 213-222.
- Hamdani M, El Abed H. , Kherallah M, Alimi A.M. (2009), “Combining Multiple HMMs Using On-line and Off-line Features for Off-line Arabic Handwriting Recognition”, 10th International Conference on Document Analysis and Recognition
- Hanmandlu M., Murali Mohan K.R., Chakraborty S, Goyal S, Choudhury D.R. (2003) “Unconstrained handwritten character recognition based on fuzzy logic” Pattern Recognition 36 pp. 603 – 623
- Haraty R. and A. Hamid, (2002) “Segmenting Handwritten Arabic Text,” Proceeding. Int. Conf. Computer Science, Software Eng., Information Technology, e-Business, and Applications, 2002.
- Haraty R. and C. Ghaddar, (2004) “Arabic Text Recognition,” International Arab Journal Information Technology, vol. 1, pp. 156-163, 2004.
- Hathaway R.J, Bezdek J. C, Hu Y. (2000), “Generalized fuzzy c-means clustering strategies using Lp norm distances” *IEEE Trans. on Fuzzy Systems*, 2000 Vol. 8, pp. 576–582
- Haxby J. (1991), Dissociation of object and spatial visual processing pathways in human extrastriate cortex. PNAS, 1991. 88: 1621-1625.
- HelmsM, Vattam S S, Goel A K.(2009). “Biologically inspired design: process and products”, Design Studies 30 pp. 606-622

- Huang Y. Kaiqi Huang , Liangsheng Wang¹, Dacheng Tao, Tieniu Tan¹ and Xuelong Li, (2008), "Enhanced Biologically Inspired Model", IEEE, 2008
- Hung C., G. Kreiman, T. Poggio and J. DiCarlo, (2005), Fast Read-out of Object Identity from Macaque Inferior Temporal Cortex. *Science*, 2005. 310: 863-866.
- Husam A.A, Zitar R.A.(2010), Development of an efficient neural-based segmentation technique for Arabic handwriting recognition, *Pattern Recognition* 43 pp. 2773–2798
- Hussain, M., Khan, N. M. (2005), "Urdu Character Recognition Using Spatial Temporal Neural Network", INMIC 2005.
- Hussain S.A, Anwar F. Sajjad A. (2007) "Online Urdu Character Recognition System." MVA2007 IAPR Conference on Machine Vision Applications 2007
- Izadi S., Sue C.Y. (2008), "Online Writer-independent Character Recognition Using a Novel Relational Context Representation" 2008 Seventh International Conference on Machine Learning and Applications
- Khan, I. Haider, "Online Recognition of Multi-Stroke Handwritten Urdu Characters", International Conference on Image Analysis and Signal Processing (IASP), 2010.
- Khedher M.Z, Al-Talib G. (2007), "A fuzzy expert system for recognition of handwritten Arabic sub-words" 2007 IEEE
- Kherallah M, Bouri F., Alimi A.M.,(2009) "On-line Arabic handwriting recognition system based on visual encoding and genetic algorithm" *Engineering Applications of Artificial Intelligence* 22 pp. 153–170
- Kherallah M., Elbaati A, Abed H.E., Alimi A.M., "The On/Off (LMCA) Dual Arabic Handwriting Database" International Conference on Frontiers in Handwriting Recognition, 2008.
- Khorsheed M.S.(2003), "Recognising Handwritten Arabic Manuscripts Using a Single Hidden Markov Model," *Pattern Recognition Letters*, vol. 24, pp. 2235-2242,
- Kuo, W.J., Yeh, T.C., Duann, J.R., Wu, Y.T., Ho, L.T., Hung, D., Tzeng, O.J., Hsieh, J.C., A left-lateralized network for reading Chinese words: a 3 T fMRI study. *NeuroReport* 12, pp 3997–4001
- Lazard G., (1995) "The Rise of the New Persian Language" in Frye, R. N., *The Cambridge History of Iran*, , Vol. pp. 595–632, Cambridge: Cambridge University Press.
- Lee J.J, Kim J., Kim J.H. (2000) "Data driven design of HMM topology for online handwritten recognition" *Proceedings of the Seventh International Workshop on Frontiers in Handwriting Recognition*, September 11-13.
- Lesica N.A. and G.B. Stanley (2004), Encoding of natural scene movies by tonic and burst spikes in the lateral geniculate nucleus. *Journal of Neuroscience*, 24: 10731–40.

- Li Y, Dong M, Hua J,(2008) "Localized feature selection for clustering", Pattern Recognition Letters 29, pp 10–18
- Liu C.L., Suen C.Y. (2009), "A new benchmark on the recognition of handwritten Bangla and Farsi numeral characters" Pattern Recognition 42 pp. 3287 – 3295
- Lin C.T., G. Lee, (1996) "Neural fuzzy system: A neuro fuzzy synergism to intelligent system" Prentice Hall
- Lin C.W, Chen Y.H., Chen L.G, "Bio-inspired Unified model of Visual Segmentation system for CAPTCHA Character Recognition ", IEEE workshop on signal processing system,
- Lorigo L.M., V. Govindaraju, "Offline Arabic Character Recognition: A Survey" Pattern Recognition and Machine Intelligence (2006) Vol. 28 pp 712-724
- Gupta M.M., (1998) "Cognition, Perception and uncertainty in :M.M. Gupta(Ed.), Fuzzy Computing, 1998"
- Maddouri S.S.and Abed H.E. (2008). Baseline extraction: comparison of six methods on IFN/ENIT database. International conference of frontiers handwritten recognition.
- Mantas J., (1986)"An overview of character recognition methodologies", Pattern Recognition, Vol. 19, No. 6, pp. 425–430,.
- Malaviya A., L. Petrs,(1999), "Multilayered Handwritten Recognition Approach", Fuzzy Sets and System Vol 104, pp 219-227
- Malik, S. Khan, S.A., (2005) "Urdu Online Handwriting Recognition", Emerging Technologies, 2005. Proceedings of the IEEE Symposium on Volume, Issue, 17-18.
- N.Mezghani , Mitiche A., Cheriet M, "Bayes Classification of Online Arabic Characters by Gibbs Modeling of Class Conditional Densities" IEEE Transactions On Pattern Analysis And Machine Intelligence, Vol. 30, No. 7, July 2008.
- Mezghani N., Mitiche A., Cheriet M.(2005), "A new representation of shape and its use for high performance in online Arabic character recognition by an associative memory" International Journal of Document Analysis 7(4): pp 201–210.
- Miled H. and N.E. Ben Amara (2001), "Planar Markov Modeling for Arabic Writing Recognition Advancement State," Proc. International Conference Document Analysis and Recognition, pp. 69-73, 2001.
- Mishkin M. (1982), A memory system in the monkey. Philosophical Transaction of the Royal Society of London Series B, 1982. 298: 85.
- Mitoma H., Uchida S., Sakoe H. (2004), "Online Character Recognition Using Eigen-Deformations", 9th Int'l Workshop on Frontiers in Handwriting Recognition, 2004

- Mitsuru M., Akira N., (2001) "Sub Stroke Approach to HMM-based On-line Kanji Handwriting Recognition", Nakai, 2001
- Mohamed, M., Gader, P. (1996) "Handwritten word recognition using segmentation-free hidden Markov modeling and segmentation-based dynamic programming techniques," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 18, no. 5, pp. 548–554, 1996.
- Mohamad R.H., Sulem L.L, Mokbel C. (2009) Combining slanted-frame classifiers for improved HMM-based Arabic handwriting recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence
- Mozaffari S., Faez K., Margner V., El-Abed H., "Lexicon reduction using dots for off-line Farsi/Arabic handwritten word recognition" Pattern Recognition Letters 29 (2008) 724–734.
- Mumolo E. ,Massimiliano Nolich, (2006)"A Biological Inspired Robotic Auditory System Based on Binaural Perception and Motor Theory" American Assoc. Artificial Intelligence Technical Report, 2006
- Omer M.A.H. and Ma S.L, "Online Arabic Handwriting Character Recognition Using Matching Algorithm" The 2nd International Conference on Computer and Automation Engineering (ICCAE), Beijing, China , 2010.
- Papanek V, (1984) Design for the real world, human ecology and social change Pantheon Books, 1984 New York
- Parui S. K., Guin K., Bhattacharya U. , Chaudhuri B. B. (2008), "Online handwritten Bangla character recognition using HMM", IEEE 2008
- Pechwitz M. and V. Ma"rgner, (2003) "HMM Based Approach for Handwritten Arabic Word Recognition Using the IFN/ENIT-Database," Proc. International Conference Document Analysis and Recognition, pp. 890-894, 2003.
- Pechwitz M., Maddouri S. S., Märgner V., Ellouze N. , Amiri N., "IFN/ENIT-Database of Handwritten Arabic Words" , 7th Colloque International Francophone sur l'Ecrit et le Document , CIFED 2002, Oct. 21-23, 2002, Hammamet, Tunis, (2002).
- Potter M. , E. Levy (1969), Recognition memory for a rapid sequence of pictures. Journal of Experimental Psychology, 1969. 81: 10-15.
- Randa I. E, Mohsen A. R., and Samia A. M. (2007) "Simultaneous Segmentation and Recognition of Arabic Characters in an Unconstrained On-Line Cursive Handwritten Document World Academy of Science, Engineering and Technology 29 , p 288-29
- Razzak M.I., S.A.Hussain, M.Sher (2009) "Combining online and offline preprocessing for online Urdu character recognition" IMECS 09, HongKong
- Razzak M.I., S.A.Hussain, M.Sher (2009)"Numeral recognition for Urdu script in unconstrained enviornment" ICET, FAST 09

- Razzak M.I., S.A.Hussain,A.Belaid, M.Sher (2010) "Multifont numeral recognition for Urdu script based languages " International journal of research trend in engineering.
- Riesenhuber M. and Poggio T.,(1999), Hierarchical models of object recognition in cortex, Nature America Inc. <http://neurosci.nature.com>
- Reinagel P. (1999), D. Godwin, S. Sherman and C. Koch, Encoding of visual information by LGN bursts. Journal of Neurophysiology, 1999. 81: 2558-2569.
- Rieke F., D. Warland, R. van Steveninck and W. Bialek, Spikes. (1997), Cambridge, Massachusetts: The MIT Press.
- Sari T., L. Souici, and M. Sellami, (2002) "Off-Line Handwritten Arabic Character Segmentation Algorithm: ACSA," Proc. International Workshop Frontiers in Handwriting Recognition, pp. 452-457, 2002.
- Sattar S.A., Haque S., Pathan M.K. (2009), "Finite State Model for Urdu Nastalique Optical Character Recognition, IJCSNS International Journal of Computer Science and Network Security, VOL.9 No.9, September 2009.
- SlomanA., (1989), "On Designing a Visual System towards a Gibbsonian computation modal of vision" J. Exp. Theor, A.I. 1989.
- Sudo T. (2002), Thesis, "On-Line Recognition of Handwritten Text Based on Hidden Markov" . Model, School of Information Science Japan.
- Shafait F, HaA.san, Keysers, BreuelT. M., (2006) "Layout Analysis of Urdu Document Images" IUPR, 2006
- Shahzad N., Paulson B., Hammond T. (2009), "Urdu Qaeda: Recognition System for Isolated UrduCharacters" IUI 2009 Workshop on Sketch RecognitionFebruary 8, 2009, Sanibel Island, Florida.
- Sherman S., Tonic and burst firing: dual modes of thalamocortical relay. Trends in Neurosciences, 2001. 24: 122-126.
- Sternby J., Morwing J., Andersson J., Friberg C., (2009) "On-line Arabic handwriting recognition with templates" pp. Pattern Recognition 42 3278 - 3286
- Saabni R, Sana J E (2009), "Hierarchical online Arabic handwritten recognition" , 10th International conference on document analysis and recognition, pp 867-871.
- Safabakhsh R. and P. Adibi, (2005) "Nastaaligh Handwritten Word Recognition Using a Continuous-Density Variable-Duration HMM," The Arabian Journal Science and Engineering., vol. 30, pp. 95-118.
- Sagheer M.W., He C.L, Nobile N., Suen C.Y. , "A New Large Urdu Database for Off-Line Handwriting Recognition" Lecture Notes in Computer Science, 5716, pp. 538-546, 2009.

- Salah AA, Alpaydin E, Akarun L (2002) A selective attention-based method for visual pattern recognition with application to handwritten digit recognition and face recognition. *IEEE Trans Pattern Anal Mach Intell* 24(3):420–425. doi:10.1109/34.990146
- Sari T, L Soucici, M Sellami, (2002) "Off-line Arabic handwritten segmentation algorithm ACSA" 8th International Workshop on Frontiers in Handwriting Recognition.
- Sattar S.A, Haque S, Pathan M.K, Gee Q, (2008) "Implementation Challenges for Nastaliq Character Recognition" *CCIS* 20, pp. 279–285.
- Schenk J., S. Schwarzler, G. Rigoll, (2008), "PCA in Online Handwritten Recognition of Whiteboard Notes: A Novel VQ Design for Use with Discrete HMM's" in: *Proc. of Int. Conference on Frontiers in Handwriting Recognition* S. 544 – 549
- Soucici L., N. Farah, T. Sari, and M. Sellami, (2004), "Rule Based Neural Networks Construction for Handwritten Arabic City-Names Recognition," *Proceeding Artificial Intelligence: Methodology, Systems, and Applications*, pp. 331-340,.
- Sternby J., J Morwing, J Andersson, C Friberg, (2009) "On-line Arabic handwriting recognition with templates" *Pattern Recognition* 42 3278 – 3286
- Subrahmonia, J., (2000). Similarity measures for writer clustering. In L. R. B. Schomaker and L. G. Vuurpijl (Eds.). *Proceedings of the Seventh International Workshop on Frontiers in Handwriting Recognition*, pp. 541-546.
- Tappert, C. C., (1984), "Adaptive online handwriting recognition". *Proceedings of International Conference of Pattern Recognition*, pp. 1004-1007.
- Tanaka H., N. Iwayama, K. Akiyama, (2004), "Online Handwritten recognition technologies and its applications", *FUJITSU Science Technology Journal* pp.170-178.
- Thorpe S., D. Fize and C. Marlot, (1996), Speed of processing in the human visual system. *Nature*, 381: 520-522.
- Wandell B.A. (1999), "Computational Neuroimaging Of Human Visual Cortex" *Annu. Rev. Neurosci.* 22:145–73
- Wang W., Brakensiek A., Rigoll G. (2002), "Combination of multiple classifiers for handwritten word recognition", *IWFHR 2002*
- Weitzenfeld A. (2008) "From schemas to neural networks: A multi-level modelling approach to biologically-inspired autonomous robotic systems" *Robotics and Autonomous Systems* 56 , 177–197.
- Wu L., Predrag Neskovic and Leon N Cooper, (2006) "Biologically Inspired Hierarchical Model for Feature Extraction and Localization" *Proceedings of the 18th International Conference on Pattern Recognition (ICPR'06)*.

Zant T.V.D, Schomaker L, Haak K, (2008), "Handwritten-Word Spotting UsingBiologically Inspired Features" IEEE transactions on pattern analysis and machine intelligence, vol. 30, no. 11, November 2008

Zheng L, (2008) "Recognition for Arabic Character Based on Edge and BPNN", Proceedings of the World Congress on Engineering and Computer Science 2008

Zhou B.L.,(2000) "Bio-inspired study of structural materials", Materials Science and Engineering C 11, pp. 13–18

APPENDIX A

PUBLICATION LIST

Publication List**Journal Paper = 10****Total Impact Factor = 10.32**

1. Muhammad Imran Razzak, Fareeha Anwar, Syed Afaq Husain, Abdel Belaid, Muhammad Sher "A Hybrid Approach: HMM and Fuzzy Logics for Online Urdu Script Character Recognition". Knowledge-Based System: Elsevier, ISSN: 0950-7051, Vol. 23, 2010, pp.914-923 (Impact Factor 1.308).
2. Muhammad Imran Razzak, Syed Afaq Hussain, Muhammad Sher "Locally Baseline detection for Online Urdu Script based Languages Character Recognition" International Journal of Physical Sciences, ISSN 1992 – 1950, Vol. 5(6), July, 2010.(Impact Factor 0.554)
3. Muhammad Imran Razzak, Syed Afaq Hussain, Muhammad Sher "Handling Diacritical Marks for Arabic Script based Languages Online Character Recognition using Fuzzy C-mean Clustering and Relative Position" Journal of Information, ISSN:1343-4500 (Impact Factor 0.09)
4. Muhammad Imran Razzak, Muhammad Sher, Syed Afaq Hussain, "Effect of Ghost Character Theory on Arabic Script based Languages Character Recognition" Przegląd Elektrotechniczny, ISSN 0033-2097 (Impact Factor 0.196).

-
5. Muhammad Imran Razzak, Syed Afaq Hussain, Muhammad Sher
“Preprocessing for Arabic Script Based Languages” International Journal of
Innovative Computing, Information and Control (Impact Factor 2.92)
 6. Muhammad Imran Razzak, Syed Afaq Hussain, Muhammad Sher,
Biologically Inspired Urdu Script Based Languages Character Recognition,
International Journal of Innovative Computing, Information and
Control(Impact Factor 2.92)
 7. Muhammad Imran Razzak, Rubiyah Yosuf, Syed Afaq Hussain, Muhammad
Sher, Rule Based Online Urdu Character Recognition, ICIC-Express Letter:
An International Journal of Research and Survey Vol.4, No.2. (HEC
Approved).
 8. Muhammad Imran Razzak, Muhammad Sher, A. Belaid, Afaq Hussain Syed,
“Multifont Numerals Recognition for Urdu Script based Languages.
International Journal of Recent Trend in Engineering. Academy. Vol. 2, pp:
70-72. (HEC Approved)
 9. Muhammad Imran Razzak, Fareeha Anwar, Syed Afaq Hussain, Muhammad
Sher, A Fuzzy Expert System: Biologically Inspired Multilayered and
Multilanguage Character Recognition, Expert System: A Journal of
Knowledge Engineering. (Impact Factor 1.231) (Accepted).

10. Muhammad Imran Razzak, Rubiyah Yousf, Afaq Hussain Syed, Muhammad Sher "HMM Based Online Urdu Character Recognition". International Journal of Computer Mathematics. (Impact Factor 0.478) (Accepted)

Conference Paper

- 1 Muhammad Imran Razzak, Syed Afaq Hussain, Muhammad Sher, Numeral Recognition for Urdu Script Based Languages in Unconstrained Environment, International Conference on Emerging Technologies, Islamabad, Pakistan.
- 2 Muhammad Imran Razzak, Muhammad Sher, Afaq Hussain Syed, "Combining Offline and Online Preprocessing for Online Urdu Character Recognition". International Conference on Imaging Engineering, Hong Kong 18-20 March, 2009.

Book Chapter

1. Abdel Belaid, Muhammad Imran Razzak, Arabic Script Based Character Recognition, Handbook of Document Image Processing and Recognition, Publisher Springer Verlag.

