

# SEQUENTIAL LABELING FOR RECOGNITION OF IMAGE BASED PATTERNS

ANUPAMA RAY



DEPARTMENT OF ELECTRICAL ENGINEERING  
INDIAN INSTITUTE OF TECHNOLOGY DELHI

AUGUST 2017

©Indian Institute of Technology Delhi (IITD), New Delhi, 2017

# **SEQUENTIAL LABELING FOR RECOGNITION OF IMAGE BASED PATTERNS**

by

**ANUPAMA RAY**

DEPARTMENT OF ELECTRICAL ENGINEERING

Submitted

in fulfillment of the requirements of the degree of Doctor of Philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI

AUGUST 2017

For you Baba

# Certificate

This is to certify that the thesis titled **Sequential Labeling for Recognition of Image Based Patterns** being submitted by **Ms. Anupama Ray** to the **Department of Electrical Engineering**, Indian Institute of Technology Delhi, for the award of **Doctor of Philosophy** is a record of bona-fide research work carried out by her under my guidance and supervision. In my opinion, the thesis has reached the standards fulfilling the requirements of the regulations relating to the degree. The work presented in this thesis has not been submitted elsewhere, either in part or full, for the award of any other degree or diploma.

Professor Santanu Chaudhury  
Department of Electrical Engineering  
Indian Institute of Technology Delhi  
New Delhi - 110 016

# Acknowledgments

I would like to express my sincere gratitude to my advisor Prof. Santanu Chaudhury, who has been a huge inspiration in my life. His ideas and his zest for work has been a constant motivation. He has been very generous with his time and always encouraged me for better research. I take this opportunity to thank him for being so understanding and supporting me with all the financial help, without which PhD would be extremely difficult.

I am extremely grateful to the members of my SRC committee: Prof. P. K. Kalra, Dr. Sumantra DuttaRoy and Dr. Brejesh Lall, for their insightful suggestions and constant help which encouraged me to widen my research from various perspectives. I am grateful to the technical staff of Multimedia Lab who allowed me own this lab for the last four years. This lab has been home to me. I thank my seniors Ritu Mam, Meera Mam, and Nisha Mam for their help and advice throughout my PhD. I was lucky to have very friendly labmates like Manoj, Ronak, Deepika and Richa around me. I have been lucky to get the opportunity to team up with a few juniors of mine, some of whom who have been my co-authors. I specially mention Sai Rajeswar, who has seen the long sleepless nights with me in Multimedia Lab without any results. Thank you for being with me, tolerating me and supporting me to get throughout. I have been very fortunate to have friends like Peeyush and Chandni who stood by me through the good and bad. I would like to express my heartiest gratitude towards Karun who stood by me patiently through all my efforts, studied with me, read my papers my thesis, helped me prepare for coding interviews and constantly encouraged me to peruse my PhD. I can never thank you enough!!

This thesis is dedicated to my Baba, it was his dream which I tried to fulfill, I know I

couldn't do my best, but I tried really hard. My mother and sister have been my oxygen tank and no thanks could ever be sufficient. I thank my family for forming my vision, exposing me to the world and allowing me to live the way I wanted. Their infallible love and support has always been my strength. Their patience and sacrifices will remain my inspiration throughout my life.

Anupama Ray

# Abstract

Sequential Labeling is a structured output prediction task which requires labeling sequence of inputs by sequence of labels, and the sequences can be of arbitrary length. The label of each input in the input sequence is influenced by the label of the neighboring inputs. This thesis aims to extend the research in sequential labeling for recognition of 2D images. Since text recognition is a sequential labeling problem, we primarily focus on text recognition from images.

In terms of methodology, the first contribution of this work is a novel graph based sequential labeling framework using multiple hypotheses at input and Conditional Random Fields (CRF)/SVM-HMM as the probabilistic recognizer to capture associations and selecting the best path. We apply this framework for segmentation based OCR and demonstrate results on degraded and noisy documents of printed and handwritten text. The second contribution of this work aims to advance the state-of-art Long Short Term Memory (LSTM), a recurrent neural network by adding dimension to make it structurally suitable to 2D images. We construct a 2D LSTM and stack several layers to create a 2D deep Bidirectional LSTM (BLSTM) network. We use this 2D deep BLSTM architecture within a hypothesis-and-verify architecture along with a learned language model to create a complete script independent, segmentation free OCR architecture for multilingual documents. Our third contribution is an end-to-end trainable Convolutional LSTM (CLSTM) architecture for text in the wild. We propose a framework for iterative text localization using character recognition which shows improvements on state-of-the-art object detection, followed by text recognition using CLSTM. The fourth contribution of this work is a novel end-to-end trainable framework with Convolutional and Fixed point layers, to capture



label dependencies as well as spatial relationships with a faster and guaranteed convergence. We have focused on Indian language scripts since research on Indian scripts are not at par with the Latin scripts. All the frameworks developed are script independent and can be trained on any script.

# सार

अनुक्रमिक लेबल एक संरचित आउटपुट भविष्यवाणी कार्य है जिसके लिए लेबल के क्रम से इनपुट के लेबलिंग अनुक्रम की आवश्यकता होती है, और अनुक्रम मनमानी लंबाई का हो सकता है। इनपुट अनुक्रम में प्रत्येक इनपुट का लेबल पड़ोसी इनपुट के लेबल से प्रभावित होता है। इस थीसिस का उद्देश्य 2D छवियों की पहचान के लिए अनुक्रमिक लेबलिंग में अनुसंधान का विस्तार करना है। चूंकि पाठ पहचान एक अनुक्रमिक लेबलिंग समस्या है, इसलिए हम मुख्य रूप से छवियों से पाठ पहचान पर ध्यान केंद्रित करते हैं।

कार्यप्रणाली के संदर्भ में, इस काम का पहला योगदान एक नवीन ग्राफ आधारित अनुक्रमिक लेबलिंग ढांचे है जो इनपुट और कंडीशनल रैंडम फ़िल्ड (CRF) / SVM-HMM कई अनुमानों का उपयोग करता है जो एसोसिएशन पर कब्जा करने और सर्वोत्तम मार्ग का चयन करने के लिए संभाव्य पहचानकर्ता है। हम विभाजन के आधार पर ओसीआर के लिए इस ढांचे को लागू करते हैं और मुद्रित और हस्तलिखित पाठ के अपमानित और शोर दस्तावेजों पर परिणाम प्रदर्शित करते हैं।

इस काम का दूसरा योगदान 2D छवियों के लिए संरचनात्मक रूप से उपयुक्त बनाने के लिए आयाम जोड़कर एक लयबद्ध शॉर्ट टर्म मेमोरी (LSTM), एक पुनरावर्ती तंत्रिका नेटवर्क को आगे बढ़ाने का लक्ष्य है। हम एक 2D एलएसटीएम का निर्माण करते हैं और एक 2D गहरी द्विदिश एलएसटीएम (BLSTM) नेटवर्क बनाने के लिए कई स्तरों का ढेर लगाते हैं। हम इस 2D गहरी बीएलएसटीएम आर्किटेक्चर का उपयोग एक परिकल्पना और सत्यापित वास्तुकला के भीतर करते हैं, जिसमें सीखी भाषा मॉडल के साथ-साथ बहुभाषी दस्तावेजों के लिए एक पूर्ण स्क्रिप्ट स्वतंत्र, विभाजन मुक्त ओसीआर आर्किटेक्चर तैयार किया जाता है।

हमारा तीसरा योगदान प्राकृतिक दृश्य चित्रों में पाठ पहचान के लिए एंड-टू-एंड ट्रेनेबल कनवोल्यूशनल एलएसटीएम (CLSTM) वास्तुकला है। वर्णित पाठ स्थानीयकरण के लिए हम ढांचे को प्रस्तावित करते हैं, जो ऑब्जेक्ट का पता लगाने में नवीनतम तकनीकों में सुधार दिखाते हैं, इसके बाद पाठ पहचान के लिये हम सीएलएसटीएम का उपयोग करते हैं।

इस काम का चौथा योगदान एक परिकल्पनात्मक और निश्चित बिंदु परतों के साथ एक अन्तराल एंड-टू-एंड ट्रेनेबल फ्रेमवर्क है, जो कि लेबल निर्भरता को पकड़ने के साथ-साथ तेज और गारंटीकृत अभिसरण के साथ स्थानिक रिश्तों को भी प्राप्त करता है। हमने भारतीय भाषा लिपियों पर ध्यान केंद्रित किया है क्योंकि भारतीय लिपियों पर शोध लैटिन स्क्रिप्ट के अनुरूप नहीं हैं। विकसित सभी ढांचा स्क्रिप्ट स्वतंत्र हैं और किसी भी स्क्रिप्ट पर प्रशिक्षित किया जा सकता है।

# Contents

<b>Certificate</b>	<b>iii</b>
<b>Acknowledgments</b>	<b>v</b>
<b>Abstract</b>	<b>vii</b>
<b>List of Figures</b>	<b>xiii</b>
<b>List of Tables</b>	<b>xv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Introduction to Sequential Labeling . . . . .	1
1.2 Scope and Challenges . . . . .	3
1.3 Objectives . . . . .	4
1.4 Contributions . . . . .	4
1.5 Layout of thesis . . . . .	5
<b>2 Sequential Labeling and Applications: A Review</b>	<b>7</b>
2.1 Traditional Algorithms . . . . .	7
2.2 Sequential Learning Algorithms . . . . .	8
2.2.1 Hidden Markov Model . . . . .	8
2.2.2 Conditional Random Fields . . . . .	9
2.2.3 Structured SVM . . . . .	9

2.3	Deep Learning for Sequential Learning Applications . . . . .	9
2.3.1	Recurrent Neural Networks and variants . . . . .	10
2.4	Review of Applications . . . . .	11
2.4.1	Review of existing OCR frameworks . . . . .	11
2.4.2	Review of OCR techniques for Indic Scripts . . . . .	12
2.5	Related work in Scene Text Recognition . . . . .	15
2.6	Motivation . . . . .	17
<b>3</b>	<b>Multi-hypothesis architecture for Sequential Labeling</b>	<b>19</b>
3.1	Introduction . . . . .	19
3.2	Multi-hypotheses framework for OCR . . . . .	20
3.2.1	Generation of Multiple Preprocessing Hypotheses . . . . .	21
3.2.2	Recognition Architecture . . . . .	23
3.2.3	Construction of Multi-Hypotheses Tree . . . . .	23
3.2.4	Feature Extraction . . . . .	25
3.3	CRF Learning . . . . .	25
3.4	SVM-HMM Learning . . . . .	27
3.5	Experimentation Results . . . . .	28
3.6	Discussions . . . . .	30
<b>4</b>	<b>Sequential Labeling using Deep Recurrent Neural Networks</b>	<b>31</b>
4.1	Introduction . . . . .	31
4.2	Recurrent Neural Networks and vanilla LSTM . . . . .	32
4.3	Construction of 2D Deep BLSTM . . . . .	33
4.4	Segmentation free Multilingual OCR architecture . . . . .	36
4.5	Overview of Proposed OCR architecture . . . . .	37
4.6	Preprocessing Hypothesis Generation . . . . .	38
4.6.1	Data preprocessing . . . . .	40
4.6.2	Recognition Engine . . . . .	41

---

4.7	Language model based Verification . . . . .	41
4.7.1	Hypothesis Verification using Language Modeling . . . . .	42
4.8	Datasets . . . . .	43
4.8.1	IAM English offline handwritten dataset . . . . .	43
4.8.2	Printed Indian Datasets . . . . .	44
4.8.3	Proposed Hindi handwritten dataset . . . . .	44
4.9	Results and Discussions . . . . .	46
4.10	Discussions . . . . .	54
<b>5</b>	<b>Text Recognition in the Wild</b>	<b>55</b>
5.1	Introduction . . . . .	55
5.1.1	Scene Text Localization . . . . .	56
5.1.2	Motivation . . . . .	58
5.2	Overview and Contributions of Proposed Work . . . . .	59
5.3	Proposed Framework for Text Localization . . . . .	60
5.3.1	Region Proposals . . . . .	61
5.3.2	Filtering of Region Proposals . . . . .	61
5.3.3	Feature Selection . . . . .	64
5.3.4	Recursive Filtering using Character Recognition . . . . .	65
5.3.5	Textline formation . . . . .	66
5.4	End-to-end trainable Convolutional-LSTM network . . . . .	67
5.5	Text Detection Results . . . . .	70
5.5.1	ICDAR 2011 dataset . . . . .	70
5.5.2	Oriented Scene Text Dataset . . . . .	70
5.5.3	MSRA TD500 dataset . . . . .	71
5.6	Recognition Results . . . . .	73
5.6.1	ICDAR 2015 Incidental scene text dataset . . . . .	74
5.6.2	SVT dataset . . . . .	76

---

5.6.3	IIIT 5K word dataset . . . . .	77
5.7	Discussions . . . . .	78
<b>6</b>	<b>Convolutional Fixed Point Network</b>	<b>81</b>
6.1	Introduction . . . . .	81
6.2	CNN for structured labeling . . . . .	83
6.3	Overview of Proposed Network . . . . .	84
6.4	Convolutional Neural Networks . . . . .	86
6.5	Fixed Point Theorem . . . . .	87
6.5.1	Contraction Condition for convex loss . . . . .	89
6.6	Applications . . . . .	91
6.6.1	1D character recognition . . . . .	91
6.6.2	2D Object recognition . . . . .	93
6.6.3	2D semantic image segmentation . . . . .	94
6.6.4	Analysis of Speed of convergence . . . . .	95
6.7	Discussion . . . . .	96
<b>7</b>	<b>Conclusions</b>	<b>97</b>
7.1	Summary of Contributions . . . . .	98
7.2	Scope of Future Work . . . . .	100
	<b>Bibliography</b>	<b>101</b>
	<b>Publications</b>	<b>127</b>
	<b>Biography</b>	<b>129</b>

# List of Figures

1.1	Sequential labeling in POS Tagging. . . . .	2
1.2	Sequence alignment issues for an Oriya word due to one-to-many and many-to-one correspondences. . . . .	3
3.1	Illustration of multi-hypothesis sequence labeling application for 1D speech sequences(adapted from [1]). . . . .	20
3.2	Illustration of word-tree formed from mutiple input hypothesis. . . . .	20
3.3	Block Diagram of proposed framework . . . . .	23
3.4	Illustration of a word-tree in Indic script formed from multiple binarization schemes used as input hypotheses. . . . .	25
3.5	Sample image from Oriya dataset with output text file . . . . .	30
4.1	Single LSTM cell. . . . .	33
4.2	Block diagram of Deep Bidirectional Long Short Term Memory Architecture. . . . .	35
4.3	Block diagram of proposed framework. . . . .	38
4.4	Sample from proposed Devanagari Text dataset . . . . .	46
4.5	Sample image from Telugu dataset with output text file . . . . .	47
4.6	a-f show model training and testing for each script . . . . .	52
4.7	a,b show model training and testing for IAM and Devanagari (offline) hand-written documents . . . . .	53
4.8	Graph showing effect of size of training data on model training . . . . .	53

4.9	Graph showing training error curve for different LSTM architectures for same script and same training data size . . . . .	53
5.1	Examples of incidental scene text in natural images (from ICDAR 2015 end-to-end recognition Challenge [2]) . . . . .	57
5.2	Block Diagram of Proposed Text localization Framework. . . . .	59
5.3	Block Diagram of Proposed CNNLSTM Text Recognition Framework . . . . .	60
5.4	a, b show prior distribution of HOG features, Color variance for text and non-text respectively (top:text; bottom:non-text) . . . . .	62
5.5	a, b show prior distribution of Intensity distribution and Stroke width for text and non-text respectively (top:text; bottom:non-text) . . . . .	63
5.6	a, b show prior distribution of Distance variance and Entropy for text and non-text respectively (top:text; bottom:non-text) . . . . .	64
5.7	a-f show some challenging successful cases from our method . . . . .	72
5.8	a-f show some failure cases of our method . . . . .	73
5.9	Results of end-to-end recognition for ICDAR 2015. The green bounding box denotes groundtruth, while the red bounding box denotes detected text regions. The recognized text has been overlaid on the detected bounding box in red. . . . .	75
6.1	Block diagram of CFPN for Recognition of 1D sequences from images . . . . .	86
6.2	Variation of error rate with respect to number of epochs trained for both 1D and 2D sequence learning. . . . .	96



# List of Tables

2.1	Review of OCRs for Indian Languages . . . . .	12
3.1	Experimental Results on different datasets . . . . .	29
3.2	Comparison with state-of-art on Oriya dataset . . . . .	29
4.1	Results of proposed 2D Deep BLSTM network on printed datasets . . . . .	46
4.2	Test Error rate comparison of different LSTM architectures for several scripts .	49
4.3	Comparison of single versus multiple hypotheses . . . . .	50
4.4	Comparison of character error rate of printed scripts with state of art techniques	51
5.1	Network parameters of proposed CNNLSTM . . . . .	67
5.2	Text localization Results on ICDAR 2011 dataset . . . . .	70
5.3	Text localization Results on OSTD dataset . . . . .	71
5.4	Text localization Results on MSRA TD500 dataset . . . . .	71
5.5	Results of end-to-end recognition for ICDAR 2013 and 2015 . . . . .	76
5.6	Results on SVT full-image dataset . . . . .	77
5.7	Results on IIIT 5K word dataset . . . . .	78
6.1	Comparison of Error rates achieved by each method . . . . .	93
6.2	Experimental Results for Semantic Segmentation on PASCAL-Context . . . . .	95