

# Example-Based Parsing for Resource-Deficient Languages

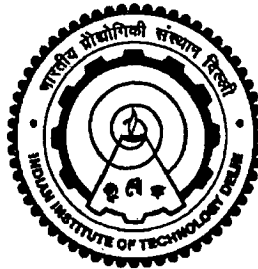
By

Shailly Goyal (Nee Kansal)

Department of Mathematics

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

*to the*



Indian Institute of Technology Delhi  
June 2007

I. I. T. DELHI.  
LIBRARY  
Acc. No. TH-3506

I.I.T. DELHI  
PROCESSED  
AGREENT

TH  
81 322  
G0Y-E

*To my dear husband . . .*

*a token for your love and inspiration*

# Certificate

This is to certify that the thesis entitled *Example-Based Parsing for Resource-Deficient Languages* submitted by *Ms. Shailly Goyal (Nee Kansal)* to the Indian Institute of Technology Delhi, for the award of the Degree of Doctor of Philosophy, is a record of the original bona fide research work carried out by her under my supervision and guidance. The thesis has reached the standards fulfilling the requirements of the regulations relating to the degree.

The results contained in this thesis have not been submitted in part or full to any other university or institute for the award of any degree or diploma.

New Delhi

June 2007



**Dr. Niladri Chatterjee**

**Supervisor**

# Acknowledgements

This thesis would not have been possible without support of many people. It is impossible to acknowledge all of them here, but there are few people who deserve a special mention.

I am forever indebted to my Ph.D. supervisor, Dr. Niladri Chatterjee, for teaching me the basics of research. I am grateful for his unfailing support throughout this work. He was great in letting me do what I wanted, gave consistently good advice, and encouraged me when the things were not going fine. He was extremely patient with my failings, and his trust in my abilities helped me overcome many difficulties.

I thank IIT Delhi for providing me teaching assistantship and the necessary facilities required to pursue my research smoothly. I am thankful to Prof. B. Chandra, Head of the Department and one of my SRC member, for her help. Special thanks are due to Prof. S. Kaushik, the other member of my SRC, for the constructive suggestions provided at various stages of this research. I am grateful to Dr. N. Nataraj for motivating me to join the Ph.D. Many thanks to my pre-Ph.D. course instructors - Prof. S. R. K. Iyenger, Prof. J. B. Srivastava and Dr. A. Tripathi - for helping me strengthen my basics. I express my regards to Prof. R. K. Sharma, DRC Chairperson, for his support. Thanks to non-teaching staff for their assistance.

I am very lucky to have many good friends and colleagues who made my experience in this department greatly enjoyable. Thanks to *Sary* and *Annie* for being

there to share the various highs and lows I faced during this research. I thank them also for the many fun moments I spent with them. I appreciate my seniors, Deepa, Sonia and Prabhakar, for making me comfortable in the department, and motivating in the initial days of my research. I will always remember Manju, Pratibha, Kanchan, Preeti, Pandeyji, Megha, Reshma, Geeta, Vaneeta, *Dinu* and all others for providing fun-filled and cheerful environment in the department. Work looked less mundane with all of them around.

I will never forget Krishna for his genuine concern. He always boosted and supported me selflessly during the whole process. I thank him also for countless impromptu academic and non-academic discussions which we had, from which I have learnt a lot about research, in particular, and life, in general. Just saying thanks will be too little a gratitude towards him.

Finally, and most importantly, thanks to many wonderful people in my personal life for giving their unwavering support throughout, especially during this work. No words can be ever enough to thank my husband, Nilesh, for his constant encouragement, understanding and endless patience. My Ph.D. had always been first preference for him, and he kept everything else on hold for this. I am grateful to my in-laws, Prof. R. K. Goyal and Mrs. Usha Goyal, for their whole-hearted support, and for providing very encouraging environment at home to pursue my research. I convey my earnest thanks to my sister-in-laws, Dr. Reena Mittal and Dr. Neena Jain, for their moral support. I would like to express my deepest gratitude towards *nanaji*, Mr. O. P. Mathuria, and *dadaji*, Mr. S. L. Goyal, for their blessings. Sincere thanks are due to my parents, Er. R. K. Kansal and Mrs. Manju Kansal, for their unending love and blessings. My siblings, Nidhi and Utkarsh, deserve a special mention for cheering me always.

Above all, I thank God for making this thesis possible.

New Delhi

  
Shailly Goyal (Nee Kansal)

# Abstract

Aim of the present research is to develop parsing schemes for natural language sentences. Parsed corpus is essential for various natural language processing (NLP) activities, but its availability cannot be guaranteed for most of the languages. Furthermore, development of parsed corpus or parser is not an easy task using traditional approaches, viz. rule-based and statistical. This is because success of these approaches almost invariably demands a huge amount of computational resources that are not typically available for most of the natural languages. We feel that example-based (EB) approaches can serve as suitable alternatives at this juncture. The major advantage of these approaches is that their demand on computational resources is much less in comparison with the traditional approaches, yet EB approaches are useful in developing robust techniques as is envisaged in many areas of artificial intelligence, NLP in particular. In this work we have pursued following two aspects of example-based parsing:

*Bilingual Parsing.* In this methodology a sentence is parsed using the parse of its parallel sentence. While projecting the syntactic relations from one language to another, we have considered similarities as well as dissimilarities between the two languages. Hence we are able to develop generalized schemes that can work on a wide variety of source-target language pairs. We have developed parsing schemes for simple as well as complex sentences.

*Monolingual Parsing.* In this scheme a sentence is parsed using the parse knowledge of examples of the same language. We have developed schemes for parsing sentences of a language by acquiring appropriate knowledge from a parsed example base of the same language. In this work we have devised ways to take care of various problems, such as unknown words, free word order property, morphological variations, effectively.

We have developed both these schemes in a generalized way with minimal dependence on linguistic knowledge so that the schemes developed can be used across a wide spectrum of languages. In this work we have done a thorough case-study on Hindi. For nitty-gritty details of the parsing schemes, where linguistic details are inevitable, we have considered English and Hindi as the source and target language, respectively. In this work we have chosen link grammar as the underlying grammar for representing the parse of sentences. One fundamental requirement therefore is a link grammar for languages under consideration. Since, no such grammar exists for Hindi, we have developed a link grammar for Hindi. For this task also we follow example-based approach.

*Development of Hindi Link Grammar.* Instead of developing the link grammar for Hindi from scratch, in this work we have made appropriate modifications in the English link grammar to suit the requirements of the Hindi grammar. We have shown how English links can be adapted for Hindi by taking care of its various grammatical nuances (e.g. free word order, noun and verb morphology, influence of subject/object on verb morphology) that make Hindi grammar distinctly different from English.

The parsing schemes developed in this work have been implemented, and tested on a reasonably sized example base. Still we have been able to demonstrate clearly the efficacy of these schemes. We feel that our research will pave the way for quick development of parsers for other languages.

# Contents

<b>Certificate</b>	<b>i</b>
<b>Acknowledgements</b>	<b>iii</b>
<b>Abstract</b>	<b>v</b>
<b>List of Figures</b>	<b>xi</b>
<b>List of Tables</b>	<b>xv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	2
1.2 The Proposed Approach . . . . .	5
1.2.1 Choosing an Appropriate Grammatical Formalism . . . . .	8
1.3 Description of the Work Done . . . . .	11
1.3.1 Organization of the Thesis . . . . .	12
1.3.2 Implementation and Evaluation . . . . .	16
<b>2 Hindi Link Grammar</b>	<b>19</b>
2.1 English Link Parser . . . . .	20
2.2 Word Order . . . . .	23

2.2.1	Effects on Links . . . . .	25
2.3	Noun Phrases . . . . .	26
2.3.1	English Noun Phrases and Links . . . . .	27
2.3.2	Hindi Noun Phrases and Links . . . . .	28
2.4	Verb Phrases . . . . .	34
2.4.1	Study of English and Hindi Verb Phrases . . . . .	35
2.4.2	Verb Phrase Links . . . . .	36
2.4.3	Complex and Compound Verbs . . . . .	40
2.5	Linking Subject/Object with Verb . . . . .	44
2.5.1	Dependence of Verb Morphology on Subject/Object . . . . .	44
2.5.2	Subject/Object-Verb Links . . . . .	46
2.6	Subordinate Clauses . . . . .	48
2.6.1	General Linking Scheme for Hindi Complex Sentences . . . . .	49
2.6.2	Noun Clause and its Links . . . . .	51
2.6.3	Adverb Clause and its Links . . . . .	54
2.6.4	Adjective Clause and its Links . . . . .	55
2.7	Cross-Linking . . . . .	59
2.7.1	Phrase Swapping . . . . .	65
2.8	Concluding Remarks . . . . .	67
<b>3</b>	<b>Bilingual Parsing: Projecting Clausal Links</b>	<b>71</b>
3.1	pDCA over Inadequacy of DCA . . . . .	73
3.1.1	The pseudo-Direct Correspondence Assumption . . . . .	77
3.2	The pseudo-Direct Projection Algorithm . . . . .	81
3.3	Parser for Hindi: A Case Study . . . . .	84
3.3.1	Identification of Phrases and Head Words . . . . .	85
3.3.2	The Rule Set $\mathcal{R}$ . . . . .	90
3.3.3	Morphological Analysis . . . . .	93
3.4	Experimental Setup and Results . . . . .	97

---

3.4.1	Illustrations . . . . .	99
3.4.2	Experimental Results . . . . .	102
3.5	pDPA for Compound and Complex Sentences . . . . .	105
3.5.1	Extending pDPA to Parse Compound Sentences . . . . .	105
3.5.2	Inadequacy of pDPA for Complex Sentences . . . . .	106
3.6	Concluding Remarks . . . . .	108
<b>4</b>	<b>Bilingual Parsing: Projecting Conjunctive Links</b>	<b>111</b>
4.1	The Overall Scheme . . . . .	113
4.1.1	Identifying Complex Sentence . . . . .	114
4.1.2	Identification of Clauses . . . . .	120
4.1.3	The Algorithm . . . . .	133
4.2	Projecting Links to Hindi Sentence . . . . .	134
4.2.1	ParseComplex1() . . . . .	136
4.2.2	ParseComplex2() . . . . .	138
4.3	Results and Discussions . . . . .	142
4.3.1	Examples . . . . .	143
4.3.2	Experimental Results . . . . .	151
4.4	Concluding Remarks . . . . .	154
<b>5</b>	<b>Monolingual Parsing</b>	<b>155</b>
5.1	Knowledge Elicitation . . . . .	158
5.1.1	EB Link Dictionary . . . . .	160
5.1.2	Phrase Templates . . . . .	165
5.2	The Parsing Algorithm . . . . .	172
5.2.1	An Illustration . . . . .	175
5.3	Handling Unknown Words . . . . .	177
5.3.1	Unknown Category Words . . . . .	178
5.3.2	New Words . . . . .	182

---

5.3.3	An Example . . . . .	186
5.4	Handling Free Word Order Languages . . . . .	187
5.5	Handling Inflectionally Rich Languages . . . . .	190
5.5.1	Stemming Various Words . . . . .	193
5.5.2	Assigning Links to New Words . . . . .	200
5.6	Overview of the Monolingual Parsing Scheme . . . . .	204
5.7	Results and Discussions . . . . .	205
5.7.1	English Results . . . . .	205
5.7.2	Hindi Results . . . . .	208
5.8	Concluding Remarks . . . . .	212
<b>6</b>	<b>Discussions and Conclusion</b>	<b>213</b>
6.1	Summarization of the Work . . . . .	217
6.2	Possible Extensions . . . . .	219
	<b>Bibliography</b>	<b>223</b>
<b>A</b>	<b>Hindi Links</b>	<b>239</b>
<b>B</b>	<b>Hindi Suffix Paradigms</b>	<b>243</b>
	<b>About the Author</b>	<b>253</b>