

IMPROVING PERFORMANCE OF HINDI TO ENGLISH PBMT USING EFFICIENT PHRASE MANAGEMENT TECHNIQUES

SUSMITA GUPTA



DEPARTMENT OF MATHEMATICS
INDIAN INSTITUTE OF TECHNOLOGY DELHI
OCTOBER 2018

©Indian Institute of Technology Delhi (IITD), New Delhi, 2018

IMPROVING PERFORMANCE OF HINDI TO ENGLISH PBMT USING EFFICIENT PHRASE MANAGEMENT TECHNIQUES

by

SUSMITA GUPTA

Department of Mathematics

Submitted

in fulfillment of the requirements of the degree of doctor of philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI
OCTOBER 2018

CERTIFICATE

This is to certify that the thesis entitled *Improving Performance of Hindi to English PBMT Using Efficient Phrase Management Techniques* submitted by *Mrs. Susmita Gupta* to the Indian Institute of Technology Delhi, for the award of the Degree of the **Doctor of Philosophy**, is a record of the original bona fide research work carried out by him 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.

Prof. Niladri Chatterjee
Department of Mathematics
Indian Institute of Technology Delhi

ACKNOWLEDGEMENTS

I am deeply indebted to *Mother Saraswati*; the Goddess of Knowledge. Due to *Her* divine grace I have been able to carry out this research work.

This thesis would not have been possible without the 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 **Professor Niladri Chatterjee** for teaching me basics of the research. I am grateful for his unfailing support throughout this work. He was great in letting me do what I wanted, and encouraged me when the things were not going fine. He has been extremely patient with my failings, and his trust in my abilities helped me overcome many difficulties.

I thank IIT Delhi for providing me necessary facilities to carry out my research work. I thank Head of Department of Mathematics, DRC Chairperson and my SRC members for their help at various stages of my research work.

I am thankful to manager of Cadence Design Systems Ms Reenee Raizada Tayal, for granting permission to join Ph.D. in IIT Delhi. I am especially thankful to my colleagues for their constant support during my research work.

Most importantly I would like to thank many people in my personal life for giving their unwavering support throughout, especially during this research work. No word can be ever enough to thank my family for their wholehearted support for carrying out research work without bothering for the day-to-day problems. Finally, I thank to all my elders for their blessings for successfully carrying out this research work.

Last but not least, I am very much thankful to the anonymous examiners and reviewers for their valuable comments and suggestions for improvising this thesis.

Susmita Gupta

ABSTRACT

Despite the arrival of Neural Machine Translation (NMT) systems traditional Phrase-based Statistical Machine Translation (PBSMT) systems are still relevant particularly for the resource poor languages. The demand on huge volume of data that is an essential requirement for NMT systems still leads researchers on Machine Translation involving such languages to use resort to PBSMT as the preferred translation paradigm. The focus of the present research is to improve the quality of PBSMT systems by using Natural Language based techniques. The schemes developed have been tested for Hindi to English MT. Translations involving resource-poor languages (e.g. Hindi) are often constrained by limited parallel corpora which is an essential knowledge source for any SMT system. To overcome this limitation, we have prescribed four different techniques at different stages of the PBSMT system. The four schemes are:

1. Efficient Pruning of the Phrases. The scheme proposed is aimed at pruning the phrases that are not pertinent for the translation task at hand. This ensures good quality translation at a much smaller time.
2. Dynamic Phrase Table Generation. The scheme proposed here uses a Phrase Table created dynamically, for a given the input sentence. This helps in generating only the useful phrases for the given translation task.
3. Improved Similarity Measurement of Source Language Phrases with special attention to Hindi. The scheme proposed helps in finding similar phrases which can help in generating translation for the given input sentence.
4. Developing Efficient Recombination Rules. As the translation for a given input is generated phrase by phrase, it is important that the generated translations follow the linguistic rules of target language. The proposed scheme aims at achieving this goal.

The combined effect of the proposed scheme is measured and improvement in translation quality and/or the time requirement while maintaining the quality is achieved. The proposed scheme should be suitable for translation involving Indian languages other than Hindi as well.

सार

तंत्रिका मशीन अनुवाद (एनएमटी) सिस्टम के आगमन के बावजूद परंपरागत वाक्यांश आधारित सांख्यिकीय मशीन अनुवाद (पीबीएसएमटी) सिस्टम अभी भी संसाधन गरीब भाषाओं के लिए प्रासंगिक हैं। एनएमटी सिस्टम के लिए डेटा (तथ्य) की भारी मात्रा में मांग की आवश्यक अभी भी मशीन अनुवाद पर शोधकर्ताओं को दूसरे तरीके पर ले जाती है, जिसमें पीबीएसएमटी का उपयोग पसंदीदा अनुवाद प्रतिमान के रूप में करने के लिए किया जाता है। वर्तमान शोध का ध्यान प्राकृतिक भाषा आधारित तकनीकों का उपयोग करके पीबीएसएमटी सिस्टम की गुणवत्ता में सुधार करना है। विकसित योजनाओं का हिंदी से अंग्रेजी एमटी के लिए परीक्षण किया गया है। संसाधन-गरीब भाषाओं (जैसे हिंदी) को शामिल करने वाले अनुवाद अक्सर समान समांतर निगम द्वारा बाधित होते हैं जो कि किसी भी एसएमटी प्रणाली के लिए एक आवश्यक ज्ञान स्रोत है। इस सीमा को दूर करने के लिए, हमने पीबीएसएमटी प्रणाली के विभिन्न चरणों में चार अलग-अलग तकनीकों को निर्धारित किया है। चार योजनाएं हैं:

1. वाक्यांशों का कुशल कटौती। प्रस्तावित योजना का उद्देश्य उन वाक्यांशों को छेड़छाड़ करना है जो अनुवाद कार्य के लिए प्रासंगिक नहीं हैं। यह बहुत कम समय पर अच्छी गुणवत्ता वाले अनुवाद सुनिश्चित करता है।
2. गतिशील वाक्यांश तालिका जनरेशन। यहां प्रस्तावित योजना इनपुट वाक्य के लिए, गतिशील रूप से बनाई गई वाक्यांश तालिका का उपयोग करती है। यह दिए गए अनुवाद कार्य के लिए केवल उपयोगी वाक्यांश उत्पन्न करने में मदद करता है।
3. हिंदी पर विशेष ध्यान देने के साथ स्रोत भाषा वाक्यांशों के बेहतर समानता मापन। प्रस्तावित योजना समान वाक्यांश खोजने में मदद करती है जो दिए गए इनपुट वाक्य के लिए अनुवाद उत्पन्न करने में मदद कर सकती हैं।
4. कुशल पुनर्संरचना नियम विकसित करना। चूंकि दिए गए इनपुट के लिए अनुवाद वाक्यांश द्वारा वाक्यांश उत्पन्न होता है, इसलिए यह महत्वपूर्ण है कि जेनरेट किए गए अनुवाद लक्ष्य भाषा के भाषाई नियमों का पालन करें। प्रस्तावित योजना का लक्ष्य इस लक्ष्य को प्राप्त करना है।

प्रस्तावित योजना के संयुक्त प्रभाव को गुणवत्ता बनाए रखने के दौरान अनुवाद गुणवत्ता में सुधार और / या समय की आवश्यकता में सुधार किया जाता है। प्रस्तावित योजना हिंदी के अलावा अन्य भाषाओं में अनुवाद के लिए उपयुक्त होना चाहिए।

Table of Contents

CERTIFICATE.....	i
ACKNOWLEDGEMENTS.....	iii
ABSTRACT.....	v
List of Figures	xiii
List of Tables.....	xv
Chapter 1 Introduction	1
1.1 Brief Description of Different MT Paradigms	2
1.1.1 Rule Based Machine Translation.....	3
1.1.2 Example Based Machine Translation	6
1.1.3 Statistical Machine Translation.....	8
1.1.4 Neural Machine Translation	9
1.2 Objectives of the Thesis	11
1.3 Organization of the Thesis	13
Chapter 2 Existing MT systems for Indian Languages.....	19
2.1 Existing MT systems for Indian Languages.....	19
2.1.1 ANGLABHARTI.....	19

2.1.2	ANUBHARATI	20
2.1.3	ANUBHARATI-II by Indian Institute of Technology, Kanpur (2004).....	20
2.1.4	ANUSAARAKA.....	21
2.1.5	MANTRA (Machine Assisted Translation Tool)	21
2.1.6	SHIVA MT System for English to Hindi	22
2.1.7	SHAKTI MT System for English to Hindi, Marathi and Telugu	22
2.1.8	SAMPARK	23
2.1.9	ANUBAAD.....	24
2.1.10	Web Based Hindi to Punjabi Machine Translator	25
2.1.11	VAASAANUBAADA.....	25
2.1.12	Google Translate.....	25
2.1.13	Babel Fish	26
2.1.14	Bing Translator	26
2.2	Difference between Hindi and English sentence structure.....	26
2.2.1	Word Order	27
2.2.2	Order of main verb and auxiliary verbs	27
2.2.3	Presence of subject in a sentence	28
2.2.4	Preposition Vs Postposition	28

2.2.5	Order of Verb and Adverb in a sentence	28
2.3	Translation Outputs and Analysis	29
2.4	Concluding Remarks	33
Chapter 3 Phrase Identification and Pruning Technique		35
3.1	Background	38
3.1.1	Statistics Based	39
3.1.2	Significance Based.....	41
3.1.3	Entropy Based.....	42
3.2	Proposed Phrase Table Pruning Technique.....	42
3.3	Experiments.....	48
3.3.1	Stratified Phrase pair reduction.....	51
3.3.2	Phrase length based reduction.....	52
3.3.3	Entropy based reduction	54
3.3.4	Source language (Hindi) Syntax Based reduction	55
3.3.5	Target language (English) Marker based reduction.....	56
3.3.6	Combination of Syntactic and Marker based technique for reduction	56
3.4	Results and Analysis	61
3.5	Concluding Remarks	64

Chapter 4 Dynamic Phrase Table Creation.....	66
4.1 Background	68
4.2 Proposed Dynamic Phrase Table Generation.....	70
4.3 Experiments.....	75
4.3.1 Setup 1: Without Feedback Loop	77
4.3.2 Setup 2: With Feedback Loop.....	81
4.4 Results and Analysis	85
4.5 Concluding Remarks	87
Chapter 5 Similarity Identification of Hindi Phrases.....	89
5.1 Background	89
5.1.1 String Based and Word Based Similarity	89
5.1.2 Tree Based Similarity	91
5.2 Proposed Similarity Identification Technique for Hindi Phrases.....	91
5.2.1 Hindi WordNet Based Similarity.....	93
5.2.2 Hindi Syntactic Phrase Similarity.....	95
5.3 Experiments.....	99
5.3.1 Golden Standard Set creation for similarity measurement	99
5.3.2 Evaluating the proposed similarity technique.....	100

5.4	Results and Analysis	102
5.5	Concluding Remarks	103
Chapter 6 Recombination of Translated Phrases		104
6.1	Background	106
6.2	Approaches used in the proposed translation technique	108
6.2.1	Phrase-based Approach.....	108
6.2.2	Rule-based Approach.....	108
6.2.3	Statistical Approach	109
6.3	Proposed Recombination Technique.....	110
6.3.1	Parsing and Sentence Type identification of the Hindi sentence.....	111
6.3.2	Noun Case-ending identification	113
6.3.3	Phrase generation	114
6.3.4	Phrase translation	115
6.3.5	Rule based Phrase Recombination.....	117
6.4	Experiments.....	122
6.5	Results and Analysis	124
6.6	Concluding Remarks	124
Chapter 7 Conclusion and Future Work		126

7.1	Proposed Translator 1.....	127
7.2	Proposed Translator 2.....	129
7.3	Result Analysis and Concluding Remarks.....	130
7.4	Future work.....	132
	References.....	134
	About the Author.....	144

List of Figures

Figure 1.1 Rule Based Machine Translation Approaches	3
Figure 1.2 Word alignment between Hindi and English Sentence	7
Figure 1.3 Word replacement by adaptation module	7
Figure 3.1 An example parallel sentence pair and word alignment.....	36
Figure 3.2 MOSES sample output	46
Figure 3.3 Redundant Parallel Phrases from MOSES	46
Figure 3.4 Output of Hindi Parser for Syntactic Phrases.....	47
Figure 3.5 Example Phrases retained after pruning scheme	48
Figure 3.6 Experiments for Phrase Reduction Techniques.....	51
Figure 4.1 Letter Based suffix array over the word 'Translation'	67
Figure 4.2 POS Structure output generated by shallow parser	72
Figure 4.3 Detailed POS structure output generated by LTRC Shallow Parser	73
Figure 4.4 Translation flow using the Proposed Dynamic Phrase Table.....	75
Figure 5.1 Algorithm for Similarity Measurement for any two phrases.....	92
Figure 5.2 Algorithm for calculating WordNet based Semantic Distance	94
Figure 5.3 Noun Phrase Similarity Measurements	97
Figure 6.1 Overall Translation Flow.....	110
Figure 6.2 Parse output from the LTRC Shallow Parser	112
Figure 6.3 Parse output for Negative Hindi Sentence.....	113

Figure 6.4 Parse output for Interrogative Hindi Sentence	113
Figure 7.1 Proposed Translator 1	127
Figure 7.2 Proposed Translator 2.....	129

List of Tables

Table 2.1 Translated outputs for Hindi sentences from Google	29
Table 2.2 Translated outputs for Hindi sentences from BabelFish.....	30
Table 2.3 Translated outputs for Hindi sentences from Bing	30
Table 3.1 Possible phrases using word alignment	37
Table 3.2 A comparison of the number of parallel sentence pairs and corresponding phrase pairs	38
Table 3.3 Two-by-two contingency table for a phrase pair (f, e)	41
Table 3.4 Average Bleu Score variation with Stratified Phrase Pair Reduction	52
Table 3.5 Average Bleu Score variation with Phrase Length Based Reduction for UMC002	53
Table 3.6 Average Bleu Score variation with Phrase Length Based Reduction for IIT Bombay Parallel Corpus	54
Table 3.7 Average Bleu Score variation with Entropy Based Phrase Pair Reduction.....	55
Table 3.8 Average Bleu Score variation with No Pruning, Syntax Based, Marker Based and Combined Scheme for UMC002.....	58
Table 3.9 Average Bleu Score variation with No Pruning, Syntax Based, Marker Based and Combined Scheme for IIT Bombay Parallel Corpus Set 1	59
Table 3.10 Average Bleu Score variation with No Pruning, Syntax Based, Marker Based and Combined Scheme for IIT Bombay Parallel Corpus Set 2	59
Table 3.11 Average Bleu Score variation with No Pruning, Syntax Based, Marker Based and Combined Scheme for IIT Bombay Parallel Corpus Set 3	60
Table 3.12 Average Bleu Score variation with No Pruning, Syntax Based, Marker Based and Combined Scheme for IIT Bombay Parallel Corpus Average.....	60

Table 3.13 BLEU Score, Phrase Table size and Average Translation time comparison using various techniques for UMC002 Corpus.....	61
Table 3.14 BLEU Score, Phrase Table size and Average Translation time comparison using various techniques for IIT Bombay Parallel Corpus Average	62
Table 4.1 BLEU Score and Runtime comparison with POS structure overlap percentage for UMC002.....	79
Table 4.2 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus for Set 1.....	79
Table 4.3 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus for Set 2.....	80
Table 4.4 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus for Set 3.....	80
Table 4.5 Average BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus for Set	81
Table 4.6 BLEU Score and Runtime comparison with POS structure overlap percentage Using Setup 2 for UMC002	82
Table 4.7 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus Set 1 Using Setup 2.....	83
Table 4.8 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus Set 2 Using Setup 2.....	83
Table 4.9 BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus Set 3 Using Setup 2.....	84
Table 4.10 Average BLEU Score and Runtime comparison with POS structure overlap percentage for IIT Bombay Parallel Corpus Using Setup	84
Table 4.11 BLEU Score and Runtime comparison with Proposed and Baseline Techniques for UMC002	86
Table 4.12 BLEU Score and Runtime comparison with Proposed and Baseline Techniques for IIT Bombay Parallel Corpus	86
Table 5.1 Weights for Similarity Measurement and their effect on average Precision for UMC002	101

Table 5.2 Weights for Similarity Measurement and their effect on average Precision for Set 1 of IIT Bombay Parallel Corpus.....	101
Table 5.3 Weights for Similarity Measurement and their effect on average Precision for Set 2 of IIT Bombay Parallel Corpus.....	102
Table 5.4 Weights for Similarity Measurement and their effect on average Precision for Set 3 of IIT Bombay Parallel Corpus.....	102
Table 6.1 Translation of Hindi to English using Google Translator.....	105
Table 6.2 Phrase translation with MOSES	116
Table 6.3 Phrase translation with Google.....	116
Table 6.4 Negative Phrase translation by Google translator.....	117
Table 6.5 Recombination rules for Simple Positive sentences.....	119
Table 6.6 Recombination rules for Interrogative sentences.....	120
Table 6.7 Recombination rules for Negative sentences.....	120
Table 6.8 BLEU score for 5000 test sentences.....	123
Table 7.1 A comparison of the BLEU score and the translation time for MOSES PBSMT system with the proposed translators for UMC002 Corpus.....	131
Table 7.2 A comparison of the BLEU score and the translation time for MOSES PBSMT system with the proposed translators for IIT Bombay Parallel Corpus.....	131