



**Handling Unknown Words in English to Nepali Statistical Machine
Translation Using Analogical Learning Approach**

A Dissertation

Submitted To

**Central Department of Computer Science and Information Technology
Tribhuvan University
Kirtipur, Kathmandu, Nepal**

**In Partial Fulfillment of the Requirements for the Degree of
Master of Science**

in

Computer Science and Information Technology

Submitted By

Pravakar Ghimire

CDCSIT, TU

(September, 2011)



**Handling Unknown Words in English to Nepali Statistical Machine
Translation Using Analogical Learning Approach**

A Dissertation

Submitted to

**Central Department of Computer Science and Information Technology
Tribhuvan University
Kirtipur, Kathmandu, Nepal**

**In Partial Fulfillment of the Requirements for the Degree of
Master of Science
in
Computer Science and Information Technology**

**Submitted By
Pravakar Ghimire**

**Supervisor
Prof. Dr. Shashidhar Ram Joshi**

**Co-Supervisor
Mr. Bikash Balami**



Handling Unknown Words in English to Nepali Statistical Machine Translation Using Analogical Learning Approach

Date :-

Recommendation

I hereby recommend that the dissertation prepared under my supervision by **Mr. Pravakar Ghimire** entitled “**Handling Unknown Words in English to Nepali Statistical Machine Translation Using Analogical Learning Approach**” be accepted as in fulfilling partial requirements for the degree of Master of Science.

Prof. Dr. Shashidhar Ram Joshi

Head of Department

Department of Electronics and Computer Engineering, Institute of Engineering,

Tribhuvan University

(Supervisor)



Handling Unknown Words in English to Nepali Statistical Machine Translation Using Analogical Learning Approach

Date :-

Recommendation

I hereby recommend that the dissertation prepared under my co-supervision by **Mr. Pravakar Ghimire** entitled “**Handling Unknown Words in English to Nepali Statistical Machine Translation Using Analogical Learning Approach**” be accepted as in fulfilling partial requirements for the degree of Master of Science.

Mr. Bikash Balami

Instructor

Central Department of Computer Science and Information Technology,

Tribhuvan University

(Co-Supervisor)



Tribhuvan University

Institute of Science and Technology

Central Department of Computer Science and Information Technology

We certify that we have read this dissertation work and in our opinion it is satisfactory on the scope and quality as a dissertation in the partial fulfillment for the requirement of Master of Science in Computer Science and Information Technology.

Evaluation Committee

Asso. Prof. Dr. Tanka Nath Dhamala

Head of Department

Central Department of Computer Science
and Information Technology

Tribhuvan University

Prof. Dr. Shashidhar Ram Joshi

Head of Department

Department of Electronics and
Computer Engineering,
Institute of Engineering,
Tribhuvan University

(Supervisor)

(External Examiner)

(Internal Examiner)

Acknowledgement

I consider it my pleasant duty to express my sincere gratitude to all the people who supported and encouraged me to complete this thesis work entitled “**Handling Unknown Words in English to Nepali Statistical Machine Translation Using Analogical Learning Approach**”.

First of all I would like to thank Tribhuvan University, Central Department of Computer Science and Information Technology for providing me this opportunity to perform this research work.

I express my sincere gratitude to my supervisor **Prof. Dr. Shashidhar Ram Joshi** for his keen supervise, able guidance and valuable suggestions. This thesis work would not be possible without his patience and support.

I'm greatly obliged to my co-supervisor **Mr. Bikash Balami** for his constant support. He was the one who was always available to deal with every obstacle that I faced during this thesis work with his insightful suggestions.

I'm also highly thankful to all the staffs and teachers of CDCSIT for providing me such a broad knowledge and enlightenment in two years of study period. Their motivation and support was really appreciable.

Special thanks to my parents and family for their unconditional love, constant support and motivation. They have always been a source of inspiration for whatever I have achieved so far.

Last but not the least, I would like to acknowledge and appreciate the direct and indirect support of my friends Mr. Dinesh kumar Khadka, Mr. Roshan Silwal and Mr. Satya Bahadur Maharjan and their willing co-operation to bring this thesis work in tangible form.

I have given my best effort to make this thesis work complete and error free but still if it contains some faults, suggestions regarding those mistakes will always be welcomed.

Pravakar Ghimire
(CDCSIT, TU)
September, 2011

Abstract

Unknown words are one of the great difficulties in the field of machine translation. In the process of translation, a system is most likely to encounter words that were not present in the available training data. While this is in part due to the segmentation issues, it is also often simply due to the lack of training data. In statistical machine translation to translate a sentence from one language to another we make use of a parallel corpus but it is not possible to a corpus to contain all the words from a whole language domain, hence the unknown word problem is obvious. In this thesis work an effective approach is used to translate those unknown words in English to Nepali statistical machine translation using word analogy. In this method the meaning of the unknown word is identified on the basis of other words presented in the corpus and the analogy between the prefixes and suffixes of those words with the unknown word.

Dedicated to,

My loving dad and mom

Table of Contents

Detail	Page no.
CHAPTER 1	
Introduction	1-5
1.1 Machine Translation	1
1.2 Approaches in Machine Translation	1
1.3 Statistical Machine Translation	2
1.4 Problems in Machine Translation	3
1.4.1 Ambiguity	3
1.4.2 Unknown Words	4
1.5 Proportional Analogies	5
CHAPTER 2	
Problem Definition	6-16
2.1 Background	6
2.2 Literature Review and Related Works	8
CHAPTER 3	
Implementation	17-20
3.1 Phases of Implementation	18
3.1.1 Identify unknown words	19
3.1.2 Solving analogical equation in source language	19
3.1.3 Translating the word triplets	20
3.1.4 Solving analogical equation in target language	20

CHAPTER 4

Testing and Analysis	21-26
4.1 Testing and Training Data	21
4.1.1 Training	21
4.1.2 Testing	25
4.2 Analysis	26
4.3 Verification and Validation	26

CHAPTER 5

Conclusion and Further Study	27-28
5.1 Conclusion	27
5.2 Further Study	27
Appendices	29-45
References	46-47

List of Tables

Detail	Page no.
Table 4.1 English parallel corpus used to train the system	21
Table 4.2 Nepali parallel corpus used to train the system	22
Table 4.3 Testing system with different unknown words	25

List of figures

Figure no.	Detail	Page no.
2.1	English to Nepali word level alignment including unknown word	7
2.2	Edit-distance tables computed between 'even' and 'usual' & 'even' and 'unevenly' along with synchronizations computed while solving the equation [even : usual = unevenly : ?]	15
3.1	Implementation model for handling unknown word	18
3.2	calculation of longest common subsequence between words 'mortal' and 'immortal'	19

Abbreviations Used

AI – Artificial Intelligence

AT – Automatic Translation

EBMT – Example Based Machine Translation

EM – Expectation Maximization

IBM – International Business Machine

INV – In Vocabulary

LCS – Longest Common Subsequence

MT – Machine Translation

NLP – Natural Language Processing

OOV – Out of Vocabulary

SMT – Statistical Machine Translation