# DECLARATION

I, Prajadhip Sinha, Registration No.: Ph.D/1780/11 Dated 22/09/2011 hereby declare that the subject matter of the thesis entitled **"Part-of-Speech tagging with Computational Approach: A Special Reference to Nepali Language"** is the record of works done by me and that the contents of the thesis did not form the basis for award of any other degree to me or to anybody else to the best of my knowledge. The thesis has not been submitted in any other University/Institute. This thesis is being submitted to Assam University for the degree of Doctor of Philosophy in Computer Science.

**Place:**                                                          (*Prajadhip Sinha*)

**Date:**                                                           Research Scholar

# ACKNOWLEDGEMENT

Upasana, Dulal, Runa, Suma, Partha, Pratap, Uday and Prachi for their continuous support, care and encouragement.

Finally I express my special appreciation and acknowledgement to my wife Barnali(Pinky), for her constant support, co-operation and sacrifice throughout my research work. Last but not the least; I thank all my well-wishers who directly or indirectly contributed for the completion of this thesis.

**Date:**                                                                 *Prajadhip Sinha*

# DEDICATION

To my Parents

Mr. & Mrs Krishnadas and Pramila Sinha

And

To my Wife and Daughter

Barnali & Prachi Sinha

This humble work is a sign of my love to

You!

# Content

## CHAPTER 1

## CHAPTER 2

## CHAPTER 3

# CHAPTER 4

## FOUNDATIONAL CONSIDERATIONS…………………………32

# CHAPTER 5

## PROPOSED POS TAGGER…………………………………..57

# CHAPTER 6

# CHAPTER 7

# LIST OF FIGURES

___

**FIGURE NAME**                                          **PAGE NO.**

# LIST OF TABLES

# LIST OF ABBREVIATIONS

NLP    Natural Language Processing

IR     Information Retrieval

IE     Information Extraction

MT    Machine Translation

BIS    Bureau of Indian Standards

CL     Computational Linguistics

CRF    Conditional Random Field

GUI    Graphical User Interface

HMM    Hidden Markov Model

HON    Honorificity

ILPOST   Indian Language Part of Speech Tagset

MSD    Morpho Syntactic Descriptions

MWE    Multi Word Expression

POST    Part of Speech Tagging

TBL    Transformation Based Learning

WSD    Word Sense Disambiguation