# DECLARATION

I, **Shri Yumnam Bablu Singh**, bearing Registration Number **PHD/1162/2010** dated 15/09/2010, hereby declare that the work presented in this thesis entitled **"WEB ENABLED MULTILINGUAL MANIPURI DICTIONARY"** is the record of work done by me under the guidance of **Prof. Bipul Syam Purkayastha (Supervisor)**, HOD, Department of Computer Science, Assam University, Silchar, Assam and **Prof. Chungkham Yashawanta Singh (Co-Supervisor),** HOD, Linguistics Department, Manipur University, Canchipur, Manipur. I further declare to the best of my knowledge that no part of the thesis has been submitted for the award of any degree to any other University or Institution.

Place: ……………

Date: ……………

**(Yumnam Bablu Singh)**

Research Scholar

*To my parents*

*Mr. Yumnam Ibobi Singh & Mrs. Y. Ibebi Devi*

# ACKNOWLEDGEMENTS

*I would like to thank my friends Kh. Raju Singha, Churjit Thiyam, Poireiton S, Shantikumar N., Krishnabati Devi and well-wisher who have helped me directly or indirectly in completion this research work.*

*I would like to dedicate the successful completion of the thesis work to my parents Yumnam Ibobi Singh and Yumnam Ibebi Devi, two brothers Rishikanta and Misra singh and sister Surbala for supporting me financially, for always encouraging me with positive thought even when I lost my hope as well as and for their unlimited help rendered to me.*

**Yumnam Bablu Singh**

Date………………

Place………………

# CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| ABL | Ablative |
| ACC | Accumulative |
| AIATSIS | Australian Institute of Aboriginal and Torres Strait Islander Studies |
| API | Application Programming Interfaces |
| ASCII | American Standard Code for Information Interchange |
| ASP | Aspect |
| ASP.NET | Active Server Page dot net |
| AMD | Anglo Manipuri Dictionary |
| BE | Be Verb |
| BN | Benefactive |
| BNC | British National Corpus |
| BIS | Bureau of Indian Standards |
| CAB | Case Ablative |
| CAUS | Causative |
| CFILT | Center for Indian Language Technology |
| CIIIL | Central Institute of India languages |
| CL | Computational Linguistics |
| CN | Case |
| CQL | Corpus Query Language |
| CRF | Conditional Random Field |
| CON | Concession |
| COP | Copula |
| CC_CCD | Conjunction_Co-ordinator |
| CC_CCS | Conjunction_Subordinator |
| CC_CCS_UT | Conjunction_Quotative |
| C-DAC | Centre for Development of Advanced Computing |
| DA | Deitic |
| DFA | Deterministic Finite Automata |
| DIST | Distal |
| DIT | Department of Information Technology |

| | |
|---|---|
| DM_DMD | Demonstrative Deictic |
| DM_DMR | Demonstrative Relative |
| DM_DMQ | Demonstrative_Wh-word |
| ETL | Electronic Text Library |
| EMPH | Emphatic |
| EVI | Evidential |
| FEMD | Friends English Manipuri Dictionary |
| GNU | General Public License |
| GUI | Graphical User Interface |
| IE | Information Extraction |
| ILCI | Indian Languages Corpora Initiatives |
| INTRA | Interrogative |
| IR | Information Retrieval |
| IPA | International Phonetic Alphabet |
| JJ | Adjective |
| OALD | Oxford Advanced Learner's Dictionary of Current English |
| OCR | Optical Character Recognition |
| PDE | Determiner |
| POS | Part of Speech |
| PPF | Post Position |
| PSP | Postposition |
| RB | Adverb |
| RD_SYM | Residuals Symbol |
| RD_PUNC | Residuals Punctuation |
| RD_UNK | Residuals Unknown |
| RD_RDF | Residuals_ForeignWord |
| RP_RPD | Particles Default |
| RP_CL | Particles Classifier |
| RP_INJ | Particles Interjection |
| RP_INTF | Particles Intensifier |
| RP | Received Pronunciation |
| RP_NEG | Particles Negation |
| SAMD | Student Anglo Manipuri Dictionary |

| | |
|---|---|
| SGML | Standard Generalized Markup Language |
| SL | Source Language |
| SQL | Structured Query language |
| TTF | Thin Film Transistor |
| TG | Together |
| TL | Target Language |
| LSI | Linguistics Survey of India |
| LU | Language Unit |
| MEPLA | Manipuri Equivalents of Parliamentary, legal and administrative Terms |
| MESAT | Manipuri Equivalents of Scientific and Technical Terms |
| MSD | Morphosyntactic Descriptions |
| MMD | Multilingual Manipuri Dictionary |
| MRD | Machine Readable Dictionary |
| MT | Machine Translation |
| MWE | Multi Word Expression |
| NEG | Negation |
| NFA | Non- Deterministic Finite Automata |
| NLP | Natural Language Processing |
| NMZ | Nominalizer |
| OED | Oxford English Dictionary |
| SOV | Subject Object Verb |
| SVO | Subject Verb Object |
| TB | Tibeto-Burman |
| URL | Uniform Resource Locator |
| WSD | Word Sense Disambiguation |
| WX | Indian Language in ASCII format |
| V_CV | Cunjunc verb |
| V_CV_C | Compound verb |
| V_CV_SPC | Conjunc verb with spa |