# DECLARATION

I, Md. Jaynal Abedin, Registration No.: Ph.D/1781/2011 Dated 22/09/2011 hereby declare that the subject matter of the thesis entitled "Automated Multiword Expressions Detection" 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:                                                    (Md. Jaynal Abedin)

Date:                                                     Research Scholar

i

*Dedicated to My Dearest Parents*

*Moulana Morthuza Hussain (father)*
*and*
*Late Reshma Khanom (mother)*

# ACKNOWLEDGEMENT

Date:                                                                                        MD. JAYNAL ABEDIN

# Contents                                      **Page No**

## List of Abbreviations

| | |
|---|---|
| AM | Association Measure |
| BIS | Bureau of Indian Standards |
| BNC | British National Corpus |
| CL | Computational Linguistics |
| D | Demonstrative |
| GA | Grammatical Annotation |
| GUI | GUI Graphical User Interface |
| HMM | Hidden Markov Model |
| IR | Information Retrieval |
| LM | Language Model |
| LT | Language Technology |
| LVC | Light Verb Construction |
| MWE | Multiword Expression |
| MWT | Multiword Term |
| NLP | Natural Language Processing |
| P | Precision |
| PV | Phrasal Verb |
| SA | Syntactic Annotation |
| SMT | Statistical Machine Translation |
| TnT | Trigrams'n'Tags |
| SVM | Support Vector Machine |
| TBL | Transformation Based Learning |
| VPC | Verb Particle Construction |
| WSD | Word Sense Disambiguation |

# List of Figures

# List of Graph

# List of Tables        Page No.