# DECLARATION

I, **Amrita Bhattacharjee**, bearing Registration Number.-Ph.D/2054/12 dated 12/09/2012, hereby declare that the subject matter of the thesis entitled **"Computational Characteristics of Words over Formal Languages"** 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.

(AMRITA BHATTACHARJEE)
Research Scholar
(Ph.D. Registration No.-Ph.D/2054/12)

Place:
Date:

# Acknowledgements

The endless thanks goes to Almighty for all the blessings He has showered onto me, which has enabled me to write this last note in my present study. I am grateful to Him for giving me the mental and physical strength needed to complete this assignment.

It gives me immense pleasure to express my heartiest sense of gratitude towards my supervisor Dr. Bipul Syam Purkayastha for his invaluable advice, constant encouragement, guidance and support without which this work would not have been possible. From him I have learnt all the necessary basic concepts and ideas for making this Mathematics oriented Computer Science study possible. He has given me the right amount of freedom throughout the present study.

I thank Kh. Raju Singha, Mr. Abhijit Paul, Mr. Arindam Dey, Md. Joynal Abedin and some other fellow scholars for their continuous support.

I also convey my thankfulness to the Head of the Department, all the faculty members and staff in the Department of Computer Science, Assam University, Silchar for their help on various occassions.

I am thankful to several anonymous referees of various journals and conferences for their valuable suggestions and comments which have immensely enhanced the quality of my work.

I thank all my family members, specially to my husband, son and daughter, for their support and encouragement. Finally it's time to thank and pay gratitude to my parents and elders for their blessings, and active co-operation.

— Amrita Bhattacharjee

## *Dedication*

This thesis is dedicated to my beloved parents

*Smti Archana Bhattacharjee and*

*Shri Lakshmikanta Bhattacharjya,*

parents-in-law

*Smti Nilima Chaudhuri and Late Nikhiles Chaudhuri,*

husband

*Biplab Chaudhuri,*

son and daughter

*Biprajyoti and Adrija*

*For their endless love, support, encouragement and sacrifice ,*

*without whom none of my success would be possible.*

# Contents

# List of Symbols

| | |
|---|---|
| $\emptyset$ | The empty set |
| $\mathbb{N}$ | The set of natural numbers |
| $\mathbb{Z}$ | The set of integers |
| FSA | Finite state automata |
| NFSA | Non-deterministic finite state automata |
| DFSA | Deterministic finite state automata |
| $\Sigma$ | An ordered alphabet |
| $\Sigma^*$ | The set of words formed from $\Sigma$ |
| $|w|_{a_i}$ | The number of occurrences of $a_i$ in a word $w \in \Sigma^*$ |
| $\Psi_{M_n}(\zeta)$ | $n \times n$ Parikh matrix over an word $\zeta$ |
| $m_{ij}$ | The element in the $i^{th}$ row and $j^{th}$ column of a matrix |
| $\frown_r$ | Ratio property |
| $\frown_{wr}$ | Weak ratio property |
| NLP | Natural language processing |
| $R(\zeta)$ | M-ambiguity reduction factor of the word $\zeta$ |
| $d_S(\alpha, \beta)$ | Stepping distance between the words $\alpha$ and $\beta$ |
| CFG | Context-free grammar |
| i.e. | That is |

# List of Tables

# List of Figures