*Dedicated to*

*Maa*

*Late Lilavati Patowary*

*&*

*Deuta*

*Mr. Kumud Chandra Patowary*

# Declaration

I, Pallabi Patowary, hereby declare that the thesis entitled *"Disease Biomarker Identification for Esophageal Squamous Cell Carcinoma using Gene Expression Data Analysis"* submitted to the Department of Computer Science and Engineering under the School of Engineering, Tezpur University, in partial fulfilment of the requirements for the award of the degree of Doctor of Philosophy is based on bonafide work carried out by me. The results embodied in this thesis have not been submitted in part or in full, to any other university or institute for award of any degree or diploma.

**(Pallabi Patowary)**

# Tezpur University

## Certificate

This is to certify that the thesis entitled *"Disease Biomarker Identification for Esophageal Squamous Cell Carcinoma using Gene Expression Data Analysis"* submitted to the Department of Computer Science and Engineering under the School of Engineering, Tezpur University in partial fulfilment of the award of the degree of Doctor of Philosophy in Computer Science and Engineering is a record of research work carried out by *Mrs. Pallabi Patowary* under my personal supervision and guidance.

All helps received by her from various sources have been duly acknowledged. No part of this thesis has been submitted elsewhere for award of any other degree.

Signature of Supervisor

(Dhruba Kumar Bhattacharyya)
Designation: Professor
School: Engineering
Department: Computer Science and Engineering
Tezpur University
Assam, India-784028

# Certificate

This is to certify that the thesis entitled *"Disease Biomarker Identification for Esophageal Squamous Cell Carcinoma using Gene Expression Data Analysis"* submitted by *Mrs. Pallabi Patowary* to Tezpur University in the Department of Computer Science and Engineering under the School of Engineering in partial fulfilment of the requirements for the award of the degree of Doctor of Philosophy in Computer Science and Engineering has been examined by us on ................................... and found to be satisfactory.

The Committee recommends for award of the degree of Doctor of Philosophy.

Signature of Principal Supervisor                  Signature of External Examiner
Date:                                              Date:

# Acknowledgment

*"Our greatest glory is not in never falling, but in rising every time we fall."*

It will be injustice to the success of this thesis if due acknowledgement is not expressed to the people who stood by me through the thick and thin of this journey. First and Foremost, I would like to thank my supervisor Prof. Dhruba Kumar Bhattacharrya for his constant support, trust, valuable feedback, encouragement and innumerable advices. He gave freedom to pursue my ideas and work at my own pace, and was always available to discuss various problems on the way. His encouragement and guidance have provided a good basis for completion of my research work.

With great pleasure I would like to express my gratitude to Prof. Sarat Saharia, Prof. Utpal Sharma, Dr. Rosy Sarmah, and Dr. Pankaj Barah who were part of my doctoral committee at Tezpur University for their valuable suggestions and guidance. I would also like to thank Prof. Jugal Kumar Kalita, University of Colorado, USA, for all his help, support and valuable guidance. I extend my sincere thanks to other faculty members specially Prof. Bhogeswar Borah, Prof. Smriti Kumar Sinha. Dr. Sanjib Deka and Dr. Arindam Karmakar and the non-teaching staff of the Department, for their generous help in various ways towards the completion of my research work.

My hard work alone would not have resulted this success without the support from my husband (Dr. Madhurjya Pratim Das) and son (Rutvik Taksheel Das).

Guidance and support from seniors, friends and other Departments of Tezpur University helped me travelled this painful journey comfortably. I would like to acknowledge my sincere gratitude to all of them specially to Prathana, Munmi, Trishna, Satya, Hussain Da, Partha, Upasana, Koyel Di, Manaswita Ba, and Nabonita.

Special thanks to all my family members for their blessings, love, and wishes

that sustained me this far. A very special thanks to Dada, Bow, Dangor Baa, Biju Baa and Saru Baa. Finally, I would like to thank all those who have directly or indirectly helped me in different capacities to complete my research work.

Last but not the least, I thank the almighty for everything.

**Pallabi Patowary**

# List of Figures

# List of Tables

# List of Algorithms

# Glossary of Terms

| | |
|---|---|
| BC | Bicluster |
| BP | Biological Process |
| CC | Cellular Component |
| DCA | Differential co-expression analysis |
| DEA | Differential expression analysis |
| DEGs | Differentially Expressed genes |
| DE | Differentially Expressed |
| ESCC | Esophageal Squamous Cell Carcinoma |
| FIR | Fractionated Irradiation |
| G | Gene |
| GCN | Gene Co-expression Network |
| GED | Gene Expression Data |
| GEO | Gene Expression Omnibus |
| GO | Gene-Ontology |
| GRN | Gene regulatory network |
| H | Hub gene |
| MF | Molecular Function |
| N | Normal |
| NMRS | Normalized mean residue similarity |
| PCC | Pearson Correlation Coefficient |
| PPI | Protein-Protein Interaction Network |
| RNA-seq | RNA sequencing |
| S | Sample |
| scRNA-Seq | single-cell RNA sequencing |
| SNMRS | Scaling-and-Shifting Normalized Mean Residue Similarity |
| T | Tumor |
| TF | Transcription Factor |