*Dedicated to*

*My family*

# Declaration

I certify that

- The work contained in the dissertation is original and has been done by myself under the general supervision of my supervisors.

- The work has not been submitted to any other Institute for any degree or diploma.

- I have followed the guidelines provided by Tezpur University in writing the thesis.

- I have conformed to the norms and guidelines given in the Ethical Code of Conduct of the university.

- Whenever I have used materials (data, theoretical analysis, and text) from other sources, I have given due credit to them by citing them in the text of the dissertation and giving their details in the references.

*parthajit borah*

**Parthajit Borah**

**Department of Computer Science & Engineering**
**Tezpur University**
**Napaam, Tezpur- 784028, Assam, India.**

**Dr. Dhruba Kr Bhattacharyya**
**Professor**

Phone::03712-275353
E-Mail : dkb@tezu.ernet.in

# Certificate

This is to certify that the thesis entitled **"Detection of Malware and Malware-based Attacks using AI Approaches"** submitted to Tezpur University in the Department of Computer Science and Engineering under the School of Engineering in partial fulfillment of the award of the degree of Doctor of Philosophy in Computer Science and Engineering is a record of research work carried out by **Parthajit Borah** under my supervision and guidance.

All help he received from various sources has been duly acknowledged. No part of this thesis has been submitted elsewhere for award of any other degree.

Signature of Supervisor
(Dhruba Kr Bhattacharyya)
Professor
Department of Computer Science and Engineering
Tezpur University
Assam, India-784028

# Certificate

This is to certify that the thesis entitled **"Detection of Malware and Malware-based Attacks using AI Approaches "** submitted to Tezpur University in the Department of Computer Science and Engineering under the School of Engineering in partial fulfillment of the award of the degree of Doctor of Philosophy in Computer Science and Engineering is a record of research work carried out by **Parthajit Borah** under my supervision and guidance.

All help he received from various sources has been duly acknowledged. No part of this thesis has been submitted elsewhere for award of any other degree.

Signature of Co-Supervisor
(Jugal Kalita)
Professor
Department of Computer Science
University of Colorado
Colorado Springs, CO 80918, USA

# Certificate

This is to certify that the thesis entitled **"Detection of Malware and Malware-based Attacks using AI Approaches"** submitted by **Mr. Partha-jit Borah** to Tezpur University in the Department of Computer Science and Engineering under the School of Engineering in partial fulfillment of the requirements for the award of the degree of Doctor of Philosophy in Computer Science and Engineering has been examined by us on .......................08/10/2024.......................... and found to be satisfactory.

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

(D. K. Bhattacharyya)

Prof. D.K. Bhattacharyya
Deptt. Of Computer Sc. & Engg.
Tezpur University.

Signature of Principal Supervisor

# Acknowledgment

It is a great pleasure to express my gratitude to all those who provided guidance and support in the successful completion of my doctoral program at Tezpur University. Achieving this milestone would not have been possible on my own. Many people extended their supportive hands, contributing to the success of my work. I would like to sincerely thank everyone who supported and assisted me throughout this journey at Tezpur University.

I am deeply indebted to my supervisor, Prof. Dhruba Kr. Bhattacharyya, for his invaluable guidance, patience, and encouragement throughout my PhD journey. His expertise and insights have been instrumental in shaping my research and helping me overcome numerous challenges. His unwavering support and belief in my abilities have been a constant source of motivation, and I am truly grateful for his mentorship and dedication.

I would also like to extend my heartfelt gratitude to my co-supervisor, Prof. Jugal Kalita, for his continuous support and constructive feedback. His detailed and insightful suggestions have greatly improved the quality of my work. His availability and willingness to assist at every stage of my research have been incredibly helpful. I deeply appreciate his commitment and contributions to my academic and professional growth.

Additionally, I am profoundly thankful to all the lab members and friends who have been a part of this journey. Their camaraderie, encouragement, and support have made this challenging path more enjoyable and manageable. I am grateful for the insightful discussions, collaborative spirit, and moments of shared laughter and perseverance.

I would like to express my sincere appreciation to the faculty of the Department of Computer Science and Engineering at Tezpur University. Their assistance and support in various administrative and academic matters have been invaluable, and I am truly thankful for their contributions.

I sincerely thank the members of my thesis review committee and the anonymous reviewers for their precious comments and feedback. Their constructive criticism and valuable suggestions have significantly improved the quality of my work.

A very special mention goes to my parents, who have constantly encouraged and supported me in every walk of my life. Their unwavering faith and love have been my foundation, and I will always be grateful for their sacrifices and guidance. I extend my heartfelt thanks to all my family members for their endless support, understanding, and encouragement throughout this journey.

I would like to thank the Almighty for providing me with the strength, patience, and perseverance needed to complete this journey. Without His blessings, none of this would have been possible.

Finally, I would like to thank those who have directly or indirectly helped me complete my research work in different capacities.

**Parthajit Borah**

# List of Figures

# List of Tables

# Glossary of Terms

| | |
|---|---|
| API | Application Programming Interface |
| MaaS | Malware-as-a-Service |
| CUDA | Compute Unified Device Architecture |
| CPU | Central Processing Unit |
| GPU | Graphics Processing Unit |
| FCG | Function Call Graph |
| CNN | Convolutional Neural Network |
| GNN | Graph Neural Network |
| GCN | Graph Convolutional Network |
| GAT | Graph Attention Network |
| GIN | Graph Isomorphism Network |
| GraphSAGE | Graph Sample and AggregatE |
| VGG | Visual Geometry Group |
| ResNet | Residual Network |
| SWaT | Secure Water Treatment |
| MAC | Message Authentication Code |
| LAN | Local Area Network |
| AP | Access Point |
| TUANDROMD | Tezpur University Android Malware Dataset |
| TUMALWD | Tezpur University Windows Malware Dataset |