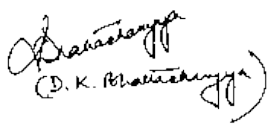*Dedicated to Maa and Deta...*

# Declaration

I, Upasana Sarmah, hereby declare that the thesis entitled *Detection of Web-based Attacks using Machine Learning Techniques* submitted to the Department of Computer Science and Engineering under the School of Engineering, Tezpur University, in partial fulfillment of the requirements for the award of the degree of Doctor of Philosophy is based on bona-fide work carried out by me. The results presented 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.

*Upasana Sarmah*

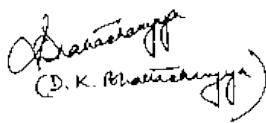**(Upasana Sarmah)**

# Tezpur University

## Certificate

This is to certify that the thesis entitled *"Detection of Web-based Attacks using Machine Learning Techniques"* submitted 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 is a record of research work carried out by **Upasana Sarmah** 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 reproduced elsewhere for award of any other degree.

(D. K. Bhattacharyya)

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

Signature of Research Supervisor
(Dr. Dhruba Kumar Bhattacharyya)
Designation: Professor
School: Engineering
Department: Computer Science and Engineering

# University of Colorado

**Certificate**

This is to certify that the thesis entitled ***"Detection of Web-based Attacks using Machine Learning Techniques"*** submitted 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 is a record of research work carried out by **Upasana Sarmah** under my co-supervision and guidance.

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

Signature of Research Co-Supervisor

(Dr. Jugal Kumar Kalita)

Designation: Professor

College of Engineering and Applied Science

Department of Computer Science

University of Colorado

Colorado Springs, CO 80918, USA

# Tezpur University

### Certificate

This is to certify that the thesis entitled *"Detection of Web-based Attacks using Machine Learning Techniques"* submitted by **Upasana Sarmah** 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 has been examined by us on _ _ _ _ _ 21/11/2024 _ _ _ _ _ _ and found to be satisfactory.

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

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

**Signature of Principal Supervisor**

**Date:** 21/11/2024

# Acknowledgment

The tenure of my research has been an intense learning experience which does not merely extend over research skills, but also inculcates a number of social skills on a personal level. It gives me immense pleasure to take this opportunity to express my deep sense of gratitude to my esteemed supervisor Prof. Dhruba Kumar Bhattacharyya. I am proud to be a student of such erudite, honest and considerate individual. His discipline and constant motivation was a key driver of my endeavor to achieve the goals of my research. He has been a constant source of inspiration throughout my PhD tenure. I sincerely appreciate his problem handling tactics, life long research experience and whole heartedly thank him for the trouble he took up to arrange facilities for my research works. This thesis is a result of his guidance, constant support, invaluable suggestions and encouragement. I also convey my heartiest thanks and gratitude to my co-supervisor Prof. Jugal Kumar Kalita, University of Colorado, USA, for all his help, support and guidance in shaping my Ph.D thesis upto this extent. It is my privilege to thank the authorities of Tezpur University and the Department of Computer Science and Engineering for providing me the facilities during the pursuit. Their co-operation and support will always be revered. I deeply acknowledge Dr. Sanjib Kr. Deka and Dr. Debojit Boro, members of my doctoral research committee for their valuable suggestions, inspirations and co-operations during the entire research tenure. I am also indebted to all the faculty members of the department, with special mention to Prof. Sarat Saharia for his relentless support during his tenure as Head of the Department. Everyday I feel blessed for the untiring moral support of my parents, Maa and Deta (who is not with me today but I know he is proud wherever he is). Maa and Deta have been and will always be my pillars of strength. This work is whole heartedly dedicated to them. I am thankful to my seniors Dr. Pooja Sharma, Dr. Nazrul Hoque, Dr. Prakash Chouhan, Dr. Hussain Ahmed Chowdhury, Dr. Monisha Devi and Dr. Pooja Dutta for giving me valuable insights whenever I was in need of it. Special mention to Dr. Ram Charan Baishya who helped me and showed me the way not only during my research tenure but also in life. I am grateful to my lab mates Khushboo, Prayakhi, Sushmita, Minhazur, and Rishang for their constant moral support. I owe a lot to my best friends Parthajit, Annushree and Suroshikha for being there for me whenever I needed that extra push. Words are less to express how much grateful I am to them. This note of acknowledgment can never be complete without a mention to a few more individuals, my constants Tanmaya

and Urma, my brothers Rimpu, Baba, Angshu and Pranjal for their constant love and appreciation. I am also grateful to all the research scholars of the department, and the office staff (Golap Da, Pronoti Baidew and Bobita Baidew) for all their help and support. Last but not the least, I thank the almighty for everything.

<div align="right">

*Upasana Sarmah*

**(Upasana Sarmah)**

</div>

# List of Tables

# List of Figures

# LIST OF ABBREVIATIONS

**OWASP**    Open Worldwide Application Security Project

**CVE**    Common Vulnerabilities and Exposures

**HTTP**    HyperText Transfer Protocol

**CI**    Critical Infrastructure

**IP**    Internet Protocol

**ZAP**    Zed Attack Proxy

**XSS**    Cross-site Scripting

**OSI**    Open Systems Interconnection

**CMIM**    Conditional Mutual Information Maximization

**MRMR**    Maximum Relevance and Minimum Redundancy

**MIFS**    Mutual Information-based Feature Selection

**GB**    Gradient Boosting

**XGB**    Extreme Gradient Boosting

**RF**    Random Forest

**EXTT**    Extra Trees

**Adaboost**    Adaptive Boosting

**RFE**    Recursive Feature Elimination

**URL**    Uniform Resource Locator