**INFORMATICS INSTITUTE OF TECHNOLOGY**
In collaboration with
**University of Westminster, UK**


**BSc. (Hons) in Computer Science**


Final Year, Project 2022/2023


# BliSS

**(Hate Speech Detection in Sinhala-English Code-Mixed Language
Deep Learning Approach)**


A dissertation by

R.S.W. B Chathurangani Ranathunga
2019923 / w1790948


Supervised By

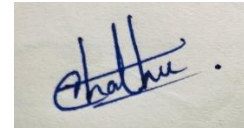Mr. Rathesan Sivagananalingam

# Declaration

This research project is wholly original, with no previous submissions for academic credit to any other organization. All materials and information that were taken from other sources have been properly cited, with proper acknowledgment of all of the sources. I acknowledge that any mistakes or omissions in the research project are entirely my responsibility, and I am aware that any type of academic misconduct will not be permitted and might result in disciplinary action. I have read and comprehended the rules and regulations pertaining to academic honesty and research ethics.

Student Full Name: Rajakaruna Seneviratne Wijekoon Bandaranayakalage Chathurangani Ranathunga

Student Registration Number: w1790948 / 2019923

Date : 02$^{nd}$ May 2023          Signature:

# Abstract

Hate speech is a growing issue in today's culture, and being able to identify it quickly could help lessen its negative effects. In current society, people use social media to express their thoughts without any doubts. People try to take revenge on social media. Especially in the YouTube video comment section, there are so many hate comments. Most people use their native language to type using English characters.

The study investigates the use of a deep learning technique called BERT for identifying hate speech in code-mixed Sinhala and English. This study intends to solve the difficulties of recognizing hate speech in code-mixed languages, which are common in multilingual nations like Sri Lanka. Hate speech identification is a crucial problem in natural language processing.

The paper compares different machine learning models with the performance of BERT models that were trained on a dataset of hate speech in a code-mixed Sinhala-English language. The results show that BERT models, with an accuracy of 92.3%, outperform other models in identifying hate speech in code-mixed languages. The research advances the creation of tools for the detection of hate speech in multilingual settings and shows the potential of deep learning techniques in this field.

Key Words : Deep Learning, Machine Learning, BERT,  Hate Speech Detection, Code-Mixed language, Sinhala, English