



**INFORMATICS  
INSTITUTE OF  
TECHNOLOGY**

**INFORMATICS INSTITUTE OF TECHNOLOGY**

In Collaboration with

**UNIVERSITY OF WESTMINSTER**

**Sinhala Hate Speech Detection System**

A Project Proposal by  
Lahiru Aravinda

Supervised by

Mr. Achala Aponso

M.Sc. Advanced Software Engineering  
University of Westminster

**August 2022**

## **Abstract**

The volume of hate content on the internet is rising quickly along with the constant growth of user-generated content. Social networking platforms, discussion boards, and blogging platforms all allow users to express their opinions without many constraints. As a result, people start hating on others who disagree with them. This study focuses on recognizing writings that include hate speech and are written in Sinhala Unicode and Singlish, which are primarily utilized online by Sri Lankans. Due to the lack of Sinhala datasets, Facebook comments and other public forums were used to construct the dataset.

In this study, the ability to identify hate speech in Sinhala was tested using 5 machine learning algorithms. Additionally, their f1- score, memory, accuracy, and precision were assessed. Then, in order to build this hate speech detection system, Support Vector Machine (SVM), Multinomial Nave Bayes (MNB), XGBoost Classifier, Random Forest, and Logistic Regression classifiers were utilized. The best amount of 85% accuracy and f1-score return for the Random Forest method was achieved by feature extraction using the TFIDF transformer and count Vectorizer.

**Key Words:** Sinhala Hate Speech Detection, Machine Learning, Supervised learning, Data Science