# EXPLOITING MULTILINGUAL CONTEXTUAL EMBEDDINGS FOR SINHALA TEXT CLASSIFICATION

G.V. Dhananjaya

218039D

Master of Science (Major Component Research)

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

May 2022

# EXPLOITING MULTILINGUAL CONTEXTUAL EMBEDDINGS FOR SINHALA TEXT CLASSIFICATION

G.V. Dhananjaya

218039D

Thesis submitted in partial fulfillment of the requirements for the degree
Master of Science (Major Component Research)

Department of Computer Science and Engineering
Faculty of Engineering

University of Moratuwa
Sri Lanka

May 2022

# DECLARATION

I declare that this is my own work and this Thesis does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or Institute of higher learning and to the best of my knowledge and belief it does not contain any material previously published or written by another person except where the acknowledgement is made in the text. I retain the right to use this content in whole or part in future works (such as articles or books).

Signature:                                        Date:  2023/04/20

The supervisors  should certify the Thesis with the following declaration.

The above candidate has carried out research for the Master of Science (Major Component Research) Thesis under our  supervision. We  confirm that the declaration made above by the student is true and correct.

Name of Supervisor: Dr. Surangika Ranathunga

Signature of the Supervisor:   Surangika      Digitally signed by          Date:
                               Ranathung      Surangika
                                              Ranathunga
                               a              Date: 2023.05.17
                                              00:19:18 +12'00'

Name of Supervisor: Prof. Sanath Jayasena

Signature of the Supervisor:                      Date:  22/05/2023

i

# DEDICATION

I dedicate this research work to all the individuals who supported me in academics, including my family, friends, and teachers.

## ACKNOWLEDGEMENT

# ABSTRACT

Language models that produce contextual representations (or embeddings) for text have been commonly used in Natural Language Processing (NLP) applications. Particularly, Transformer based, large pre-trained models are popular among NLP practitioners. Nevertheless, the existing research and the inclusion of low-resource languages (languages that primarily lack publicly available datasets and curated corpora) in these modern NLP paradigms are meager. Their performance for downstream NLP tasks lags compared to that of high-resource languages such as English. Training a monolingual Language model for a particular language is a straightforward approach in modern NLP but it is resource-consuming and could be unworkable for a low-resource language where even monolingual training data is insufficient. Multilingual models that can support an array of languages are an alternative to circumvent this issue. Yet, the representation of low-resource languages considerably lags in multilingual models as well.

In this work, our first aim is on evaluating the performance of existing Multilingual Language Models (MMLM) that support low-resource Sinhala and some available monolingual Sinhala models for an array of different text classification tasks. We also train our own monolingual model for Sinhala. From those experiments, we identify that the multilingual XLM-R model yields better results in many instances. Based on those results we propose a novel technique based on an explicit cross-lingual alignment of sentiment words using an augmentation method to improve the sentiment classification task. There, we improve the results of a multilingual XLM-R model for sentiment classification in Sinhala language. Along the way, we also test the aforementioned method on a few other Indic languages (Tamil, Bengali) to measure its robustness across languages.

**Keywords**: Multilingual language models, Multilingual embeddings, Text classification, Sentiment analysis, Low-resource languages, Sinhala language

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| Abbreviation | Description |
|---|---|
| AP | Auxiliary phrases |
| CL | Continual Learning |
| MLM | Masked Language Modeling |
| MMLM | Multilingual Language Models |
| NER | Named Entity Recognition |
| NLG | Natural Language Generation |
| NLI | Natural Language Inference |
| NLP | Natural Language Processing |
| NLU | Natural Language Understanding |
| NSP | Next Sentence Prediction |
| POS | Part-of-Speech tagging |
| RNN | Recurrent Neural Network |
| TLM | Translation Language Modelling |