# DUPLICATE BUG REPORT DETECTION USING PRE-TRAINED LANGUAGE MODELS

K.A.Udeshika Sewwandi

(199363T)

Degree of Master of Science in Computer Science

Department of Computer Science and Engineering

Faculty of Engineering

University of Moratuwa

Sri Lanka

July 2022

# DUPLICATE BUG REPORT DETECTION USING PRE-TRAINED LANGUAGE MODELS

Kahandawala Arachchige Udeshika Sewwandi

(199363T)

Thesis/Dissertation was submitted in partial fulfillment of the requirements for the degree MSc in Computer Science specializing in Data Science

Department of Computer Science and Engineering

Faculty of Engineering

University of Moratuwa

Sri Lanka

July 2022

# DECLARATION

I declare that this is my own work and this thesis/dissertation 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:…………………

The above candidate has carried out research for the PhD/MPhil/Masters thesis/dissertation under my supervision. I confirm that the declaration made above by the student is true and correct.

Name of the Supervisor: Dr. Surangika Ranathunga

Signature of the supervisor: ................................          Date:…………………

# ACKNOWLEDGEMENTS

# ABSTRACT

Software testing and defect reporting are significant factors of software development and maintenance. Defects are identified and reported in a bug tracking system like JIRA, or Bugzilla. Those reported defects are further triaged by an expert who has an understanding of the repository, system, and developers and assigns them to the developers to fix them. During this defect reporting there can be duplicate bugs reported and identifying duplicate bugs is a crucial task. Manual labeling of duplicate defects is time-consuming, may identify defects as duplicate bug reports, and also increases the cost of software maintenance. Therefore automated duplicate bug report detection is very significant. This research proposes a duplicate bug report classification methodology that leverages the Pre-trained language models BERT and XLNet with Multi-Layer Perceptron as the Deep Learning classifier for duplicate bug detection. We tested on publicly available datasets related to Eclipse, NetBeans, and OpenOffice bug reporting datasets. The selected models were shown to outperform the previously proposed systems for the same task. Among them, the approach used with BERT embeddings has shown the best results. Further experiments showed that BERT is capable of domain adaptation –meaning that even when the BERT was fine-tuned with different bug report datasets, it is still capable of detecting duplicate bugs in an unseen dataset. Finally, a multi-stage classification was done using a Convolutional Neural Network model and a BERT model using Eclipse and NetBeans datasets and a combined dataset of Eclipse and NetBeans. The approach used with the combined dataset has outperformed the baseline approach.

Keywords: Duplicate Bug detection, BERT, XLNet, MLP, CNN, Domain Adaptation, Multi-Stage Classification

# Table of Contents

# List of figures

# List of tables