



Neural Machine Translation of Sinhala to English

D. G. P. Hansadi
(Reg. No.: MS23002142)

A THESIS
SUBMITTED TO
SRI LANKA INSTITUTE OF INFORMATION TECHNOLOGY
IN PARTIAL FULFILMENT OF THE REQUIREMENTS
FOR THE DEGREE OF
MASTER OF SCIENCE IN INFORMATION TECHNOLOGY

December 2024

I certify that I have read this thesis and that in my opinion it is fully adequate, in scope and in quality, as a thesis for the degree of Master of Science.

Dr. Dilshan De Silva

Approved for MSc. Research Project:

MSc. Programme Co-ordinator, SLIIT

Approved for MSc:

Head of Graduate Studies, FoC, SLIIT

DECLARATION

This is to certify that the work is entirely my own and not of any other person, unless explicitly acknowledged (including citation of published and unpublished sources). The work has not previously been submitted in any form to the Sri Lanka Institute of Information Technology or to any other institution for assessment for any other purpose.

Sign: .....

D.G.P. Hansadi

Date:16/01/2025.....

ABSTRACT

Neural Machine Translation of Sinhala to English

D. G. P. Hansadi

MSc. in Information Technology

Supervisor: Dr. Dilshan De Silva

December 2024

This addresses the challenges of translating Sinhala in the form of Singlish, a unique and informal variation of Sinhala written in English letters, into standard English, aiming to bridge the communication gap for Sinhala-speaking individuals. The research explores the application of advanced transformer architectures such as T5, mBART, and MarianMT, each chosen for their proven capabilities in handling multilingual and low-resource language tasks. The study involves training and fine-tuning these models on a curated Singlish and English parallel sentences including dataset and evaluates their performance across several key metrics, including Quantitative and Qualitative analysis. The system developed as part of this research integrates practical features like speech-to-text for Singlish input, text-to-speech for English output, and translation history for user convenience, making it a comprehensive tool for real-time translation needs. The thesis not only aims to produce an effective translation system but also contributes valuable insights into optimizing transformer models for low-resource languages, offering benchmarks and techniques for future research in neural machine translation. By combining cutting-edge technology with a user-centric design, this research seeks to enhance accessibility for Sinhala speakers and sets the stage for advancing machine translation for underrepresented languages and dialects.

ACKNOWLEDGEMENT

I would like to express my deepest gratitude to my supervisor, Dr. Dilshan De Silva, for his exceptional guidance, unwavering support, and encouragement throughout the course of this research. I am truly grateful for his patience and for always being available to discuss ideas, provide direction, and offer constructive criticism. His mentorship has not only guided me academically but has also been a source of great inspiration, motivating me to persevere through challenges and refine my work to a higher standard.

TABLE OF CONTENTS

DECLARATION	iii
ABSTRACT.....	iv
Neural Machine Translation of Sinhala to English.....	iv
ACKNOWLEDGEMENT	v
List of Figures.....	viii
List of Tables	ix
Chapter 1 Introduction	10
1.1 Background.....	10
1.2 Evolution of Machine Translation and the Rise of Neural Approaches	10
1.3 Problem Statement.....	11
1.4 Research Objectives.....	12
1.5 Significance of the Study	14
1.6 Review of Existing Literature	19
1.7 Gap Identification	32
1.8 Conceptual Framework.....	34
1.8.1 Research Problem	34
1.8.2 Characteristics of Input Data: Singlish-English Sentence Pairs.....	34
1.8.3 Data Preprocessing:.....	34
1.8.4 Model Selection and Comparison	35
1.8.5 Model Training and Fine-Tuning Process.....	35
1.8.6 System Deployment and Real-Time Translation	36
1.8.7 Evaluation and Validation.....	36
1.8.8 Implications and Contributions	36
Chapter 2 Research Methodology.....	38
2.1 NMT.....	38
2.2 Research Design.....	45
2.3 Sampling Strategy and Dataset Preparation.....	48
2.4 System Development and Model Architecture	50
2.5 Comparative Analysis of Applied Transformer Architectures	55
2.6 Connecting the model with the Translation Web Application Using Flask.....	57
2.7 Testing and Deployment	60
Chapter 3 Results	62
3.1 Qualitative Analysis Results	62
3.2 Quantitative Analysis Results	65

3.3 Performance Evaluation among three selected Transformers Architectures	66
3.4 Optimization of T5-base with BPE.....	67
Chapter 4 Discussion	69
Chapter 5 Conclusion.....	75
References	78

List of Figures

Figure 2. 1 High-level Diagram of Translation Application.....	46
Figure 3. 1 Example Translation outputs	63
Figure 3. 2 Translating Various Forms of Sentences.....	63
Figure 3. 3 Comparison 1 – Comparing T5-Base, Marian-MT and mBart and Achieving Same Outputs... ...	64
Figure 3. 4 Comparison 2 – Comparing T5-Base, Marian-MT and mBart and Achieving Correct Outputs ...	64
Figure 3. 5 Comparison 3 – Comparing T5-Base, Marian-MT and mBart and Achieving Outputs with Slights Changes.....	65
Figure 3. 6 Comparison 4 – Comparing T5-Base, Marian-MT and mBart and Found Some Performance Differences Among Models	65

List of Tables

Table 3. 1 Performance Evaluation.....	67
Table 3. 2 Performance Evaluation of T5-base after BPE Integration.....	68