# Sinhala - Tamil Statistical Machine Translation (SMT) for Official Documents

Farook Fathima Farhath

(168034C)

Degree of Master of Philosophy in Computer Science and Engineering

Department of Computer Science And Engineering

University of Moratuwa

Sri Lanka

Oct 2018

# Sinhala - Tamil Statistical Machine Translation (SMT) for Official Documents

Farook Fathima Farhath

(168034C)

Thesis submitted in partial fulfillment of the requirements for the
Degree of Master of Philosophy in Computer Science and Engineering

Department of Computer Science And Engineering

University of Moratuwa
Sri Lanka

Oct 2018

# Declaration

I, Farook Fathima Farhath, 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.

Also, I hereby grant to University of Moratuwa the non-exclusive right to reproduce and distribute my dissertation, in whole or in part in print, electronic or other medium. I retain the right to use this content in whole or part in future works (such as articles or books).

Signed:

Date:

The above candidate has carried out research for the MPhil thesis under my supervision.

Name of Supervisor: Prof. Sanath Jayasena

Signature of supervisor:                                        Date:

Name of Supervisor: Dr. Surangika Ranathunga

Signature of supervisor:                                        Date:

## *Abstract*

Sinhala and Tamil are declared to be the official languages of Sri Lanka. This requires each government related dissemination/communication to be done in both the languages. Even though the requirement for translation is higher, the number of available human translators is limited. One feasible option to boost the productivity would be assisting the human translators with machine translation output. Here the machine translation output is given to translators to work on by post editing, rather than translating from the scratch. However, Sinhala - Tamil pair does not have any well-performing machine translation system. Therefore, the focus of this research is to develop a machine translation system for short official government documents.

This thesis presents two main contributions towards building 'Si-Ta', the first domain-adapted machine translation system for Sinhala - Tamil. The first contribution is building the baseline translation system. The second is implementing data pre-processing techniques to improve the translation quality of the baseline system.

The baseline system was built using Moses, a phrase-based statistical translation system. This was the feasible option with the available resources.

To improve the quality of the translation, three main approaches were explored. They are: (a) domain adaptation, (b) integration of terminology, dictionary, and name lists, and (c) addressing out-of-vocabulary (OOV) problem using word-embedding-based paraphrasing.

In order to adapt the system for the domain of official government documents, different language model design techniques and a data filtration technique were experimented. Under terminology integration, experiments were carried out to evaluate the effect of incorporating bilingual terminology lists to the system. Moreover, a novel data augmentation technique was experimented to generate parallel data using bilingual lists and available parallel data. Further, open domain dictionary entries, as well as a list of person names and addresses were integrated and evaluated. In addition, word-embedding-based paraphrasing was used along with a novel heuristic-based filtering to address the out-of-vocabulary issue.

All the above-mentioned approaches gave an improvement over the baseline, apart from data filtering technique. Yet, all these scores were above the scores of already available machine translation systems for this language pair. Though our techniques/approaches were evaluated only on Sinhala - Tamil pair, they are feasible to be applied to other low-resourced, highly inflectional language pairs.

**Keywords:** Machine Translation, Sinhala, Tamil, Domain Adaptation, Terminology Integration, Out-of-vocabulary

## *Acknowledgements*

# Contents

# List of Figures

# List of Tables