# - Language Assistant -

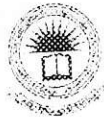## Towards Automatically Translating Text from English to Tamil in a Controlled Environment

*By*

## S . Mohanarajah

### ( 93 / M.Sc.Com.Sc / 32 )

31405

*Dissertation Submitted in*
*Partial Fulfilment of the Requirement for the*
## M.Sc in Computer Science
*at the*
### Department of Statistics and Computer Science
### University of Colombo, Sri Lanka

*June* — **1998**

# Abstract

Automatic translation is the process of automatically converting a text in one natural language to another using computers. The resultant translated text should tally with the original text in meaning, effect and style. Since this automated translation process involves handling a huge amount of dictionary information, and complex grammar and word formation rules, the vocabulary and sentence structures of the input text may be restricted based on domain specific knowledge to get a better translation.

Sri Lanka is a multi-ethnic country. Communication between the ethnic groups is essential for co-operation in the interest of national development. Translation is the only viable solution for communication between communities which use different languages. Sri Lanka requires a six-way translation tool between Tamil, English and Sinhala. The major objective of this study, as a first step for the overall motivation, is to build a model for a Machine Translation system that could translate text in English to Tamil in a controlled environment.

The proposed system is built on transformer based architecture. A code for information interchange in Tamil (SCIIT) is designed and used in this project. The system developed in this project works well for the English sentence structures found in elementary level textbooks. This system may be used readily for demonstration purposes, however, it needs further enhancements for it to be used in a real translation task.

# Contents

Abstract

Acknowledgment

Declaration

Contents

## CHAPTER III

LINGUISTIC ANALYSIS ON ENGLISH

## CHAPTER IV

LINGUISTIC ANALYSIS ON TAMIL

## CHAPTER V

## CHAPTER VI

## CHAPTER VII

LIMITATIONS AND FUTURE EXTENSIONS

## CHAPTER VIII

SUMMARY & CONCLUSION