# A CORPUS-BASED INVESTIGATION OF THE INTERLANGUAGE OF UNIVERSITY STUDENTS IN EAST MALAYSIA



## RESEARCH MANAGEMENT INSTITUTE
## UNIVERSITI TEKNOLOGI MARA MALAYSIA
## 40450 SHAH ALAM,
## SELANGOR, MALAYSIA

### PREPARED BY:

### PUAN LILLY METOM
### PROF. MADYA DR. SIMON BOTLEY
### @ FAIZAL HAKIM

### 2011

# TABLE OF CONTENTS

# ABSTRACT

This report describes the culmination of the learner corpus project called CALES (Corpus-based Archive of Learner English in Sarawak). The original two phases of the project collected 356,000 words of learner writing in the form of argumentative essays written by students taking English proficiency courses in UiTM's Sarawak Branch Campus (Botley et al, 2005, 2007). This new project has increased this total to over 480,000 words of digital text, and has collected more essays from four different institutions in order to further expand and enrich the corpus.

The project follows the methodological principles laid down by the International Corpus of Learner English (ICLE) project in Belgium (Granger et. al., 2002). The data was digitised and analysed in order to investigate different types of language error. A number of observations were made concerning some prominent error categories in the data, and their pedagogical implications were explored.

It is hoped that these findings will further contribute to our understanding of the way in which Malaysian learners of English actually perform in their writing. Also, it is hoped that the outcomes of this project will form a foundation for a larger-scale understanding of the interlanguage (Selinker, 1972) of Malaysian EFL learners at university level, as well as providing a data resource for future research in this area.

# CHAPTER 1: INTRODUCTION

## 1.1 Research Background and Problem Statement

Educators in EFL (English as a Foreign Language) are all too aware of the errors, or performance features[1] frequently found in writing produced by students of English. However, EFL educators are often unable to make full use of the information revealed by such features in order to help students to improve their language performance. One reason for this is a lack of reliable and clear examples taken from real student texts. Such examples could then be used as a source of reference to help teachers predict the features students may display in their writing and speech, and then do something about them in a systematic and principled manner.

Many EFL educators tend to rely upon their professional experience and linguistic intuitions to predict what kinds of features will be displayed in the writing of a particular non-native-speaker group. For instance, it is widely known that Malaysian learners of English regularly under-use the definite article, and turn non-countable nouns onto countable ones (*a staff*, rather than *a member of staff*).

Errors such as these may be traced back to the L1 which in most cases in Malaysia is Bahasa Melayu, a language which does not have a system of definite and indefinite articles, and in which the notion of countability is somewhat different to that in English (see Botley, Haykal and Monaliza, 2005 for a recent discussion of this issue).

---

[1] Here, we prefer the term 'performance features', because common terms such as 'errors' or 'mistakes' can be considered judgemental and prescriptive. Furthermore, see the section on definition of terms below.