# Zulu-English Cross-Language Information Retrieval: an analysis of errors

## By

## Johannes Gerhardus Nel

*Submitted in fulfilment of the requirements for the degree of*

*Magister Informationis Scientia*

At the

Department of Information Science
*Faculty of Engineering, Built Environment and Information Technology*

University of Pretoria, Pretoria

*Supervisor: Dr Erica Cosijn*
*November 2003*

©

> *Who knows for what we live, struggle, and die? Wise men write many books, in words too hard to understand. But this, the purpose of our lives, the end of our entire struggle, is beyond all human wisdom...*

**Alan Paton**

## Abstract

This study is an analysis of the published results of Zulu to English Cross Language Information Retrieval experiments. In these studies, approximate string matching was used to match the running text of the natural language queries to the base forms of the dictionary entries and this aspect was found to be reasonably successful. However, on a conceptual level, Zulu-English CLIR was very poor, in part due to the lack of technical terminology in Zulu. These studies have shown that unique translation problems are to be encountered if the language pairs used contain disparate vocabularies and increases the complexity of CLIR.

In support of the analysis in this study, a detailed overview of the Zulu language as one of the indigenous languages of South Africa is provided, together with a discussion on the linguistic structure of Zulu. Furthermore, the rationale for Zulu-English CLIR is given within the context of the digitisation of indigenous knowledge.

The manual analysis of the query set used in the original experiments shows that two main problems can be identified, namely dictionary problems and translation problems. The dictionary problems are further subdivided into issues relating to orthography, borrowed words and proper names while the translation problems consist of issues relating to paraphrasing, word inflection, homonyms, affixes and mismatched word forms. These two primary problems have further been categorised into problems related either to matching issues, or to culture-related issues.

It is envisioned that interactive CLIR, which will enable a user to select the best-translated keys and add his/her own source keys, might be an ideal solution for Cross-Cultural Information Retrieval. It is certainly worth investigating the following proposed solutions: query expansion, applying metadata to describe the content, applying normalisation, and improving dictionary coverage to manage untranslatable terms.

I declare that

**Zulu-English Cross Language Information Retrieval: an analysis of errors**

is my own work, and that all sources applied or quoted have been indicated and acknowledged by means of complete references.

JG Nel

*What you put on paper and how you put it there reveal your standards of excellence, your knowledge, and the quality of your thinking more eloquently than anything else about you.*

**Leedy**

## Acknowledgements

I would like to express my sincere gratitude to:

- Dr Erica Cosijn, for her invaluable advice, encouragement, motivation, exceptional leadership and insight;
- Prof. Theo Bothma, for the opportunity to study at the Department of Information Science;
- Prof. Hannes Britz (for everything);
- Dr Magdalena Nel, for her motivation and setting an example to follow;
- My parents, for their love, their faith in me, their support and financial assistance;
- My friends, for their prayers and support;
- Clare Hassett, for her wonderful encouragement and faith in me;
- Hilda Kriel, mentor and dear friend;
- Mareli Esterhuysen, for editing and proofreading of the thesis;
- Audrey, Mzangwa, Pansy, Violet and others for their assistance in the translation of the texts; and
- The Lord, for the talent He has given me.

**Dedicated to my grandfather**

**JG Nel, sen.**

**1929–1997**

"Keep your face to the sunshine and you cannot see the shadows."

Helen Keller

# Table of Contents

*We have to do the best we can. This is our sacred human responsibility.*

**Albert Einstein**

# List of figures

# List of tables