UNIVERSITEIT VAN PRETORIA
UNIVERSITY OF PRETORIA
YUNIBESITHI YA PRETORIA

# The compilation of corpus-based Setswana dictionaries

# Fannie Sebolela

A thesis submitted in accordance with the requirements in the field of study: DLitt: African Languages at The University of Pretoria.

December 2009.

**Supervisor:** **Prof. D.J. Prinsloo**

# Summary

The aim of this thesis is to describe how corpus-based Setswana dictionaries should be compiled. The challenge to the modern Setswana lexicographer is to compile very practical descriptive and user-friendly dictionaries. A detailed evaluation of existing Setswana dictionaries will be performed in terms of the macrostructural and microstructural aspects:

- Coverage of frequently used words.
- Effective use of dictionary space.
- Use of standard dictionary conventions.
- Choice, ordering and composition of translation equivalent paradigms.

The focus will be on material collection and corpus building. Informants will be used to compile an oral corpus of 100,000 tokens. All ethical requirements such as informed consent requirements (See Appendix 1) will be honoured. Since the text corpus is an organic corpus, thus not a designed corpus aimed at balance and representativeness, the oral corpus will be constructed in the same way i.e. only basic selection criteria:

- Mother tongue speakers of Setswana.

- Adults (to be on a par with authors of the written sources in the text corpus). Age: ranging from 20-60 years.

- Male and female.

Critical analysis of all currently available Setswana dictionaries will be done with special reference to the dictionaries of Brown (1987) (SESD), Snyman, et al. (1990), Matumo (1993).(MSED), Kgasa (1976) (THAND) and Kgasa and Tsonope

(1995).(THAN) In all these cases the strategy would be in terms of the theoretical criteria and best practices in terms of a broad theoretical survey of core aspects of dictionary compilation. Finally, the study will be concluded with an analysis of corpus integrity and stability of Setswana corpora based on the model introduced by Prinsloo and De Schryver (2001a).

# Declaration

I declare that **The compilation of corpus-based Setswana dictionaries** is my own work and all sources that I have used or quoted have been indicated and acknowledged by means of complete references.

_____

Fannie Sebolela

# Acknowledgements

I am deeply indebted to Professor D.J. Prinsloo for his constructive criticism, purposeful guidance and patience without which this study would probably have been abandoned. May the Lord God preserve you.

My sincere gratitude also goes to Professor P.M. Sebate, Professor M. R. Malope, Dr R. Noormohamed, Ms. R. Ramagoshi, Ms. V. Nicolson and lastly, my dear sister Ms. C. C.L. Ngobeni for correcting and sharing some ideas in this research.

I am also grateful to Mr. M. Phaladi for his willingness to assist in the technical IT support.

Lastly, to my late wife Mmemme and our beloved children, Neo, Motumiseng, Omolemo and my grandson Motheo: You have been there for me from day one. This was a rough one, but finally I have completed the research.

# Table of contents

# List of figures

# List of tables