# Developing an XML-based, exploitable linguistic database of the Hebrew text of Gen. 1:1-2:3

Thesis by

**Jan Hendrik Kroeze**

**(21345393)**

Submitted in fulfilment of the requirements for the degree

**Philosophiae Doctor (Information Technology)**

in the

**School of Information Technology**

in the

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

**University of Pretoria**

**Pretoria**

**Promoter: Prof. Dr. T.J.D. Bothma**

**Co-promoter: Dr. M.C. Matthee**                    **April 2008**

# Developing an XML-based, exploitable linguistic database of the Hebrew text of Gen. 1:1-2:3

## JH Kroeze

## Abstract

The thesis discusses a series of related techniques that prepare and transform raw linguistic data for advanced processing in order to unveil hidden grammatical patterns. A threedimensional array is identified as a suitable data structure to build a data cube to capture multidimensional linguistic data in a computer's temporary storage facility. It also enables online analytical processing, like slicing, to be executed on this data cube in order to reveal various subsets and presentations of the data. XML is investigated as a suitable mark-up language to permanently store such an exploitable databank of Biblical Hebrew linguistic data. This concept is illustrated by tagging a phonetic transcription of Genesis 1:1-2:3 on various linguistic levels and manipulating this databank. Transferring the data set between an XML file and a threedimensional array creates a stable environment allowing editing and advanced processing of the data in order to confirm existing knowledge or to mine for new, yet undiscovered, linguistic features. Two experiments are executed to demonstrate possible text-mining procedures. Finally, visualisation is discussed as a technique that enhances interaction between the human researcher and the computerised technologies supporting the process of knowledge creation. Although the data set is very small there are exciting indications that the compilation and analysis of aggregate linguistic data may assist linguists to perform rigorous research, for example regarding the definitions of semantic functions and the mapping of these functions onto the syntactic module.

# Summary

**Title:** Developing an XML-based, exploitable linguistic database of the Hebrew text of Gen. 1:1-2:3

**Candidate:** Jan Hendrik Kroeze (1958)

**Promoter:** Prof. dr. T.J.D. Bothma

**Co-promoter:** Dr. M.C. Matthee

**Departments:** Information Science; Informatics

**School:** School of Information Technology

**Degree:** Philosophiae Doctor (Information Technology)

The thesis touches on various sub-disciplines of computational linguistics, investigating the use of XML tagging to capture linguistic categories in the Hebrew text of Genesis 1:1-2:3 and to construct a threedimensional databank; using string processing algorithms to round-trip the data to and from the databank and computer program; using array processing to explore the semantic patterns hidden in the marked-up text; and using a graphical visualisation to investigate the mapping of semantic and syntactic functions. The thesis hopes to make a contribution by demonstrating the rigour enforced by the application of data-warehousing and data-mining concepts to a linguistic databank. It proposes a macro-structure that may be used in future to package and integrate multidimensional linguistic data.

Chapter 2 experiments with a threedimensional data structure, using Visual Basic 6, and finds that a threedimensional array could be used to represent inherently multidimensional linguistic data regarding Biblical Hebrew clauses. Various layers of linguistic knowledge can be integrated by stacking various modules of analysis onto each other.

Chapter 3 explores data warehousing and online analytical processing concepts to find ways to render meaningful subsets of linguistic data stored in a threedimensional

array. Concepts like slicing and dicing are adjusted to make them useful for the processing of linguistic data.

Chapter 4 tries to find a more elegant solution for the permanent storage of the databank using XML technology. Due to its flexibility XML is chosen to build a text-based databank. The experiment indicates that XML is indeed a very suitable mark-up technology that can be used to permanently store the linguistic data in a separate databank because it allows users to create their own tag sets which may simulate a multidimensional database structure.

Chapter 5 investigates round-tripping in order to satisfy the requirement of finding a stable platform for the data, while also allowing editing and advanced processing of the data. In addition, various viewing and searching functions are discussed. Create, update and delete functionalities are added to enable users to populate and edit the clause cube while it is in the array state and to save these updates both to the RAM and on permanent storage in XML format.

Chapter 6 focuses on the benefits of text data mining facilitated by the preceding technologies. Some data mining concepts are applied in two experiments by aggregating aspects of the semantic and syntactic modules tagged in Genesis 1:1-2:3. Computer-assisted explorations of the semantic and syntactic data captured in the XML database illustrate the rigour enforced by such a text-mining venture.

In Chapter 7 projects are suggested (one of which is implemented) that could use the XML-based data cube of Genesis 1:1-2:3 in visualisation ventures to clearly show linguistic patterns uncovered by means of a computer program. These techniques may be used to create user-friendly interfaces that may facilitate easier and more intuitive mining of linguistic data.

**Keywords**

The following keywords represent the most important aspects covered in the thesis:
- Threedimensional array
- Online analytical processing (OLAP)
- XML

- Round-tripping
- Database management
- Data warehousing
- Text-data mining
- Computational linguistics
- Visualisation
- Hebrew Bible

# Acknowledgements

I would like to thank the following people:

- My supervisors, Prof. Theo Bothma and Dr. Machdel Matthee, for your excellent supervision and feedback on the related papers and articles produced during the execution of the doctoral project

- The Head and Acting Head of the Department of Informatics, Prof. Carina de Villiers and Prof. Trish Alexander, the Dean of the EBIT Faculty, Prof. Roelf Sandenbergh, as well as the Department's management committee, for granting me various shorter and longer periods of study and research leave to work on the research project, and for subsidising various overseas trips to read the papers that developed parallel to the chapters of this thesis

- Mr. Danie Malan of the University of Pretoria's library who helped me with literature searches and to find sources that were not always easy to get hold of

- All the editors and peer reviewers, as well as other colleagues, who assisted me with advice, help and comments that helped to steer the research in the right direction

- Ms. Mariëtte Postma for her excellent language editing

- My wife, Irma, for allowing me and assisting me to embark, for the second time, on the lengthy process of writing a doctoral thesis, and for helping to get aspects regarding the philosophy of science in the introduction right

- My son, Jan, for writing the Java progam discussed in Chapter 7 as an example of the graphical visualisation of linguistic data

- My daughter, Christien, for creating the beautiful graphics of threedimensional data cubes, used as illustrations in Chapters 2 and 3

- My parents, Jan and Clasie Kroeze, and my mother-in-law, Babs Jansen van Rensenburg, as well as my other family, in-laws and friends for your continued support and encouragement

# Declaration of originality

I declare that this thesis, *Developing an XML-based, exploitable linguistic database of the Hebrew text of Gen. 1:1-2:3*, is my own work, that all the sources that I have used or quoted have been indicated and acknowledged by means of complete references, and that it has not been submitted for a degree at another university.

J.H. Kroeze

# Table of contents

**Databank:** Gen1_InputV15.xml

I **Addendum I: XML clause cube corrected** CD

AddendumI_XMLClauseCube_Corrected_20080411_Fin.pdf

**Databank:** Gen1_InputV15b.xml

J **Addendum J: XML clause cube style sheet** CD

AddendumJ_XMLClauseCube_StyleSheet_20080411_Fin.pdf

**Style sheet:** Gen1XMLdb03c.css

K **Addendum K: Viewing the clause cube in a browser** CD

AddendumK_ViewingClauseCubeInBrowser_20080411_Fin.pdf

L **Addendum L: Source code of Chapter 5** CD

AddendumL_SourceCode_Chapter5_20080411_Fin.pdf

**Program:** Gen1_XML_VB6_CRUD_Beta15_Ch5.exe

**Databank:** Gen1_InputV15_RT1.xml

M **Addendum M: Source code of Chapter 6** CD

AddendumM_SourceCode_Chapter6_20080411_Fin.pdf

**Program:** Gen1_XML_VB6_CRUD_SemFSynF_Ch6.exe

**Databank:** Gen1_InputV15_RT1.xml

N **Addendum N: Visualisation program (graphical topic map)** CD

**Program:** semantics.bat

# List of Figures