# SEMANTIC TRANSFORMATIONS FOR XML QUERIES

By

**Dung Xuan Thi LE**
*BSc. (Swinburne), MSc. (La Trobe)*

THIS THESIS IS PRESENTED FOR

THE DEGREE OF DOCTOR OF PHILOSOPHY

DEPARTMENT OF COMPUTING, FACULTY OF SCIENCE

2012

# Statement of Authorship

I certify that the work in this thesis titled "Semantic Transformations for XML Queries" has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree to any other university or institution other than Macquarie University.

I also certify that the thesis is an original piece of research and it has been written by me. Any help and assistance I have received in my research work and the preparation of thesis itself have been appropriately acknowledged.

In addition, I certify that all information resources and literature used are indicated in the thesis.

Dung Xuan Thi LE

# Acknowledgements

I am deeply indebted to Prof. Mehmet Orgun for his supervision, support and advice which have been a great inspiration. The completion of this thesis would not have been possible without his patience, input and evaluation of my earlier efforts. I would also like to thank my co-supervisor Dr. Peter Busch who has also evaluated my drafts. I would like to extend my acknowledgement and thanks to my former supervisors Dr. Wenny Rahayu, Dr Eric Pardede and Dr. Stephane Bressan. I especially wish to acknowledge the efforts of Dr. Eric Pardede; this thesis would not have been possible without his dedicated time to evaluate my draft, friendship, support and encouragement. I thank Dr Michael Hill for his time to check and proofread part of my work.

I would also like to thank Dr. Saqib Ali, Oldooz Dianat, Dr. Paul Conilione, Dr. Regina Ryan, Lachlan Williams, Mary Law, Trairat Em-Oj and Elizabeth Le for their friendship, help and support.

I am grateful to the families of Dinh, Kelly, Bob and Bill and to the Vincent De Paul Society for their love and generosity. Without their kindness, I would not be here in Australia today.

I would also like to thank the directors of GMS and MMM Technology, my much earlier employers, who gave me many opportunities to gain invaluable experience during my period of employment. In particular, I thank the Ang family for their kindness, generosity and encouragement of my further study.

This thesis is specially dedicated to a very special person who has had a great influence on what I do and who I have become. His constant support, encouragement and love have made me believe in a better tomorrow.

Saving the most important for last, my heartfelt thanks go to my parents for their unconditional love and sacrifices which have inspired me to make them proud. I thank my dear brothers and sisters for their continuous support and love.

# List of External Publications & Statement of Contributions

This thesis contains material that has been published and prepared for submission as follows:

1. *Le, D. X. T., Stéphane Bressan, David Taniar, J. Wenny Rahayu, Eric Pardede: Using semantics for XPath query transformation. IJWGS 6(1):58-94 (2010).*

2. *Le, D. X. T., Bressan, S., Taniar, D., Rahayu, W., "Semantic XPath Query Transformation: Opportunities and Performance". In the 12th International Conference on Database Systems for Advanced Applications (DASFAA), 2007, pp. 994-1000. (Ranked in A category in ERA-2010).*

3. *Le, D. X. T., Bressan, S., Pardede, E., Rahayu, W., Taniar, D., "A Utilization of Schema Constraints to Transform Predicates in XPath Query". In the 21st International Conference on Database and Expert Systems Applications (DEXA), 2010, pp. 331-339. (Ranked in B category in ERA-2010)*

The first paper is the consolidation and revision of work reported in the second and third papers; they resulted from my PhD research at La Trobe University. The co-authors J. Wenny Rahayu and Eric Pardede were my PhD supervisors there. Stéphane Bressan and David Taniar were my external PhD advisors. I collaborated with Stéphane Bressan on background of semantic transformations and experimental design. I collaborated with David Taniar on the area of result presentations. My contributions to these three papers as a whole are 80%.

4. *Le, D. X. T., Bressan, S., Pardede, E., Rahayu, W., Taniar, D., "Semantic Transformation Approach with Schema Constraints for XPath Query Axes". In the 11th International Conference on Web Information System Engineering (WISE), 2010, pp. 456-470. (Ranked in A category in ERA-2010).*

The fourth paper above resulted from my PhD research at La Trobe University. The co-authors J. Wenny Rahayu and Eric Pardede were my PhD supervisors there. Stéphane Bressan and David Taniar were my external PhD advisors with whom I collaborated on the analysis of the experimental result of my research. My contributions to this paper as a whole are 85%.

5. *Le, D. X. T., Pardede, E., "On Using Semantic Transformation Algorithms for XML Safe Update". In the 8th International Conference on Information Systems Technology and its Applications (UNISCON), 2009, 276-278. (Ranked in B category in ERA-2010).*

6. *Le, D. X. T., Pardede, E., "Towards Performance Efficiency in Safe XML Update". In the $8^{th}$ International Conference on Web Information Systems Engineering (WISE), 2007, pp. 563-572. (Ranked in A category in ERA-2010).*

The fifth and sixth papers above resulted from the research collaboration that combines my research on *semantic transformations* and the *safe update* research conducted by Dr Eric Pardede to propose semantic safe updates for XML queries. My contributions to these papers including the proposed techniques, implementation and analysis of the results as a whole are 65%.

7. *Le, D. X. T., Orgun, M., "Evaluation of Semantic XPath Queries for XML-Enabled Databases". (In Preparation for Journal Submission).*

The seventh paper above is prepared for a journal submission. It has resulted from the work reported in chapters 8 and 9 of this thesis. My contribution to this paper is 80%.

# Abstract

The ever-increasing adoption of XML has created a need to ensure that XML query languages perform efficiently. *Query optimization* and *trans*formation for XML query languages, both syntactically and semantically, have received much attention from research communities in recent years. However, due to the fast progress of the application of XML data management solutions, XML-Enabled Database Management Systems still face several challenges. Among these challenges is query processing, especially the processing of XML queries specified with XPath axes and redundancies that may exist in predicates used in XML queries.

Semantic query optimization utilizes constraints in XML schemas to directly optimize a given query with a set of optimization rules. Due to the current complexity of the XML data structure which is enabled by rich semantics in XML Schemas, semantic query optimization should be performed in a more systematic manner. For a complete solution, this research proposes a series of semantic transformations to transform given XML queries to semantically equivalent, but more efficient, XML queries for optimization purposes, by using the semantics provided in XML Schemas. The proposed semantic transformations are grouped into three categories: (1) *Semantic Path Transformations*, (2) *Semantic Transformations for XPath Queries Specified with Predicates,* and (3) *Semantic Transformations for XPath Queries Specified with XPath Axes*.

After a semantic transformation is applied to an XML query, the equivalent semantic XML query can be processed more efficiently by an XML data management system and returns the same result set.

The proposed semantic transformations are then translated into a series of algorithms which are implemented and empirically evaluated for their efficiency and effectiveness. The experimental studies were carried out by using both real data sets (DBLP) and Benchmark data sets (Michigan) to illustrate that the majority of semantic transformations achieved significantly improved performance in XML

query processing; this also enabled the research presented here to identify semantic transformations as optimization devices.

# Table of Contents

# List of Figures