FACULTY OF ENGINEERING AND INFORMATION
TECHNOLOGY

# Topic-based Analysis for Technology Intelligence

**Hongshu Chen**

A thesis submitted for the Degree of

Doctor of Philosophy

U|T|S|

University of Technology, Sydney
January, 2016

# CERTIFICATE OF AUTHORSHIP/ORIGINALITY

This thesis is the result of a research candidature conducted jointly with another University as part of a collaborative Doctoral degree. I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as part of the collaborative doctoral degree and/or fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Candidate

_____

# ACKNOWLEDGEMENTS

First and foremost I would like to express my sincere gratitude to my supervisors, Prof. *Jie Lu* and Prof. *Guangquan Zhang*, who not only offered me a unique opportunity to study in UTS, but also have provided tremendous guidance and support for my research and life in the past four years. I still remember the first day I met Prof. *Jie Lu* in her office, she knew that I was so nervous to step into a new area of research, thus told me, "doing research is just like climbing up a mountain, the higher you go, the more beautiful scenery you will see". Her enthusiasm for research and the excellent example she has provided as a successful researcher were so motivational, has become my life treasure. I also would like to thank Prof. *Guangquan Zhang* for all his continuous guidance, contribution of ideas, time, and discussions during my study. Both of my supervisors have supported me academically and emotionally through the road to finish my PhD study. Their respectful personality and precise academic attitude have benefited me so much. Without their excellent supervision and continuous encouragement, this thesis could not have been finished on time. I really appreciate all the kind help from them.

As a dual-degree PhD student, I would like to express my earnest thanks to my supervisor in the Beijing Institute of Technology (BIT), Professor *Donghua Zhu*, who gave me an enormous amount of help and support since the year 2009. I have significantly benefited from his invaluable suggestions and support. Without his kind help and professional supervision, this research could not have been accomplished either.

My experience at the Decision Systems and e-Service Intelligent (DeSI) Lab in the centre for Quantum Computation and Intelligent Systems (QCIS), has been more than amazing. The great lab members of mine have contributed immensely to my professional

and personal time at UTS. I would like to thank all the friends that provide me with kind support during my PhD study. Especially, I would like to acknowledge Dr. *Dianhuang Wu*, Mr. *Junyu Xuan*, Mr. *Fan Dong* and Mr. *Yi Zhang* for their help and advice on my experiments. In addition, I am grateful to all members of the Knowledge Management and Data Analysis lab, School of Management & Economics, BIT. Thanks for all their support and help.

I want to dedicate this thesis to my grandfather, who has passed away in year 2005, just before I received my college acceptance letter. He has always told me how much he treasured me, trusted me and loved me, and how important for a person to be educated, which shaped my values and made me the person that I am today.

Last but not least, I would like to thank my family for their continuous support during my PhD study. I have amazing parents, who have provided me with unconditional support and encouragement all these years.  My family is my energy source and the harbour of my heart.

I love both of you, Mom and Dad.

*Hongshu Chen*

7th January 2016
At Sydney

# ABSTRACT

Since the past several decades, scientific literature, patents and other semi-structured technology indicators have been generating and accumulating at a very rapid rate. Their growth provides a wealth of information regarding technology development in both the public and private domain. However, it has also caused increasingly severe information overload problems whereby researchers, analysts and decision makers are not able to read, summarize and understand massive technical documents and records manually. The concept and tools of technology intelligence aims to handle this issue.

In the current technology intelligence research, one of the big challenges is that, the frameworks and applications of existing technology intelligence conducted semantic content analysis and temporal trend estimation separately, lacking a comprehensive perspective on trend analysis of the detailed content within an area. In addition, existing research of technology intelligence is mainly constructed on the fundamentals of semantic properties of the semi-structured technology indicators; however, single keywords and their ranking alone, are too general or ambiguous to represent complex concepts and their corresponding temporal patterns. Thirdly, systematic post-processing, forecasting and evaluation on both content analysis and trend identification outputs are still in great demand, for diverse and flexible technological decision support and opportunity discovery.

This research aims to handle these three challenges in both theoretical and practical aspects. It first quantitatively defines and presents temporal characteristics and semantic properties of typical semi-structured technology indicators. Then this thesis proposes a framework of topic-based technology intelligence, with three main functionalities, including data-driven trend identification, topic discovery and comprehensive topic

evaluation, to synthetically process and analyse technological publication count sequence, textual data and metadata of target technology indicators. To achieve the three functionalities, this research proposes an empirical technology trend analysis method to extract temporal trend turning points and trend segments, which help with producing a more reasonable time-based measure; a topic-based technological forecasting method to first discover and characterize the semantic knowledge underlying in massive textual data of technology indicators, meanwhile estimating the future trends of the discovered topics; a comprehensive topic evaluation method that links metadata and discovered topics, to provide integrated landscape and technological insight in depth. In order to demonstrate the proposed topic-based technology intelligence framework and all the related methods, this research presents case studies with both patents and scientific literature. Experimental results on Australian patents, United States patents and scientific papers from Web of Science database, showed that the proposed framework and methods are well-suited in dealing with semi-structured technology indicators analysis, and can provide valuable topic-based knowledge to facilitate further technological decision making or opportunity discovery with good performance.

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES