# Query Expansion for Document Retrieval Using Wordnet

## BY

## MOHD HAFIZ BIN HUSSIN
## BACHELOR OF COMPUTER SCIENCE (Hons)

## THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE REQUIREMENT FOR BACHELOR OF COMPUTER SCIENCE

## FACULTY OF COMPUTER AND MATHEMATICAL SCIENCES
## UNIVERSITI TEKNOLOGI MARA

## OCT 2010

# Acknowledgement

Firstly, I would like to express my gratitude towards Allah S.W.T that of his blessing that gives me capabilities, strength, and good health in order to finish up this project on time. In addition, this project successfully completed with my own effort and extra incentive. I would like to thank everybody who supported me with my project, especially my family for their moral support, love, advise and encouragement give me some confident through the progress of my research project.

Furthermore would like to special thank my thesis supervisor, Pn. Hayati Abd Rahman for encouraging, teaching, guiding, advising and supporting me to make sure that this project is successfully completed in time. Without the cooperation and her help, perhaps I have difficulties and I will not able to complete this project on time. Her patience in guiding me gives me confident to complete my project and I really appreciate it. I would also like to thank Dr. Noor Elaiza Binti Abd Khaled, our Final Year Project Coordinator for her guidance, support, encouragement, and criticism through the progress of my research project.

Last but not least, I would like to thank my friends, lecturers and everyone that is directly or indirectly involved in this project.

Thank you.

# ABSTRACT

This paper describes the experimentation conducted to test the effectiveness of query expansion. Several experiments generating queries extracted from WordNet. Results show that lexical expansion is not able to improve retrieval performance. Nevertheless, the experiments allow us to conclude that query expansion can benefit searching process which allows structured queries.

**KEYWORDS**

Information retrieval, query expansion, sports, WordNet.

# Table of Contents