# A Novel Quartet-Based Method for Inferring Evolutionary Trees from Molecular Data

A thesis submitted in fulfilment of the requirements for the degree of Doctor of Philosophy

School of Information Technology
Faculty of Engineering and Information Technology

Monther Tarawneh
March 2008

*Science is nothing but perception*
*Plato (Ancient Greek*
*Philosopher*
*428-348 BC)*

*To my mother*

*To my father*


أهداء الى امي وابي

# Acknowledgment

First and foremost, I would like to thank my supervisor Dr. Bing Bing Zhou. Without him this thesis would not have been possible. I am grateful to his kind supports, excellent advice, intensive collaboration throughout the whole time we were working together, for intensively reading this thesis, and the financial support he had gave me during my PhD research. Furthermore, I thank my associate supervisor Prof. Ablert Zomaya for his support and opportunity to work in the environment of advanced network research. Also, for his encouragement and help during my PhD.

Thanks to Chen, Daniel, Penghao, and every one attended our very intensive enjoyable collaboration during the last 3 years, every meeting result in new idea, some key definition and deep sight. I am especially grateful to Daniel Chu who contributed by implementing the distributed version of QB algorithm.

I appreciate very much the unfinished work with Abdur Sikder and Gowrie, it dose open my eyes into many future direction, hope we can finish.

I would like to thank every one in the ANRG group at the school of IT especially, Michael Charleston. Besides, Friday meeting is a good experience for every research student.

Thanks to David London, Greg and the workshop people for their technical support during my time in the school of information technology/University of Sydney.

I wish to thank several colleagues for their continuous support in any way, and creating such friendly environment especially Abdur, Gowrie, Shwen, Penghao, Khalid, Mohammad, Tanveer, and all members in the school of IT at Sydney University.

Finally, I would like to thanks my parent for their ever lasting support of my work and study.

# Declaration

I, Monther Tarawneh, do hereby declare that this submission is my own work and that, to the best of my knowledge and belief, it contains no material previously published or written by another person, nor material which to a substantial extent has been excepted for the award of any degree or diploma of a university or other institute of higher learning, except where due acknowledgments were made in the text.

…………………..

Monther A Tarawneh

# Abstract

Molecular Evolution is the key to explain the divergence of species and the origin of life on earth. The main task in the study of molecular evolution is the reconstruction of evolutionary trees from sequences data of the current species.

This thesis introduces a novel algorithm for inferring evolutionary trees from genetic data using quartet-based approach. The new method recursively merges sub-trees based on a global statistical provided by the global quartet weight matrix. The quarte weights can be computed using several methods. Since the quartet weights computation is the most expensive procedure in this approach, the new method enables the parallel inference of large evolutionary trees.

Several techniques developed to deal with quartets inaccuracies. In addition, the new method we developed is flexible in such a way that can combine morphological and molecular phylogenetic analyses to yield more accurate trees. Also, we introduce the concept of critical point where more than one possible merges are possible for the same sub-tree. The critical point concept can provide information about the relationships between species in more details and show how close they are. This enables us to detect other reasonable trees.

We evaluated the algorithm on both synthetic and real data sets. Experimental results showed that the new method achieved significantly better accuracy in comparison with existing methods.

# List of Figures

# List of Tables

# Table of Contents