# SOLVING MULTIPLE SEQUENCE ALIGNMENT PROBLEM

# UTILIZING INTEGER LINEAR PROGRAMMING

**PUDIAHWAI ANTON WIBOWO**

**0303010281**

**UNIVERSITY OF INDONESIA**

**FACULTY OF MATHEMATICS AND SCIENCE**

**DEPARTMENT OF MATHEMATICS**

**DEPOK**

**2008**

# SOLVING MULTIPLE SEQUENCE ALIGNMENT PROBLEM

# UTILIZING INTEGER LINEAR PROGRAMMING

This *skripsi* is submitted in partial fulfillment
of the requirements for the degree of
*Sarjana Sains*

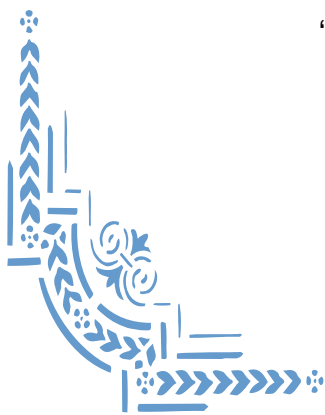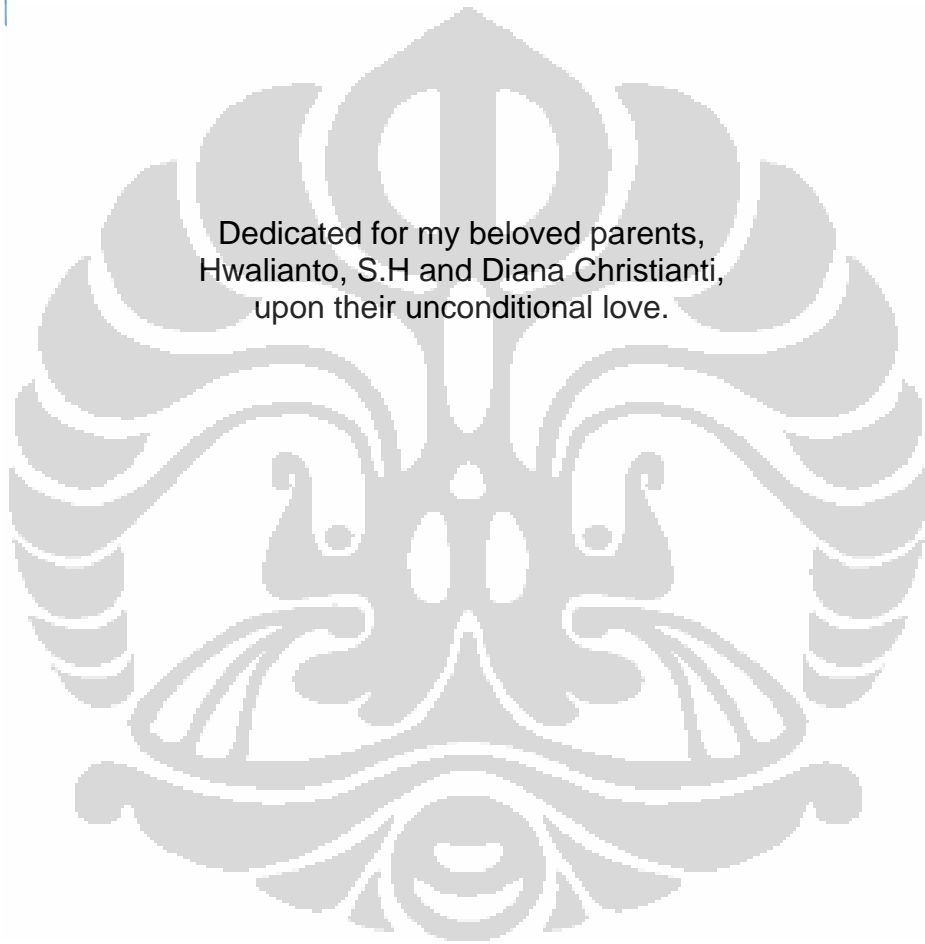By:
**PUDIAHWAI ANTON WIBOWO**
**0303010281**

**DEPOK**
**2008**

Dedicated for my beloved parents,
Hwalianto, S.H and Diana Christianti,
upon their unconditional love.

"*Sikap sama dengan tujuan.*"
-My Father-

*SKRIPSI* : SOLVING MULTIPLE SEQUENCE ALIGNMENT PROBLEM

UTILIZING INTEGER LINEAR PROGRAMMING.

NAME : PUDIAHWAI ANTON WIBOWO

NPM : 0303010281

THIS *SKRIPSI* HAS BEEN APPROVED

DEPOK, JULY 17$^{th}$, 2008

| | |
|---|---|
| <u>Dra. Denny R. Silaban, M.Kom</u> | <u>Dr. Kiki A. Sugeng, M.Si</u> |
| ADVISOR I | ADVISOR II |

Date of completion of *sidang sarjana* examination: July 17$^{th}$, 2008

Examiner I : Dra. Denny R. Silaban, M.Kom

Examiner II : Bevina D. Handari, Ph.D

Examiner III : Dra. Yahma Wisnani, M.Kom

# ABSTRACT

One of the dominant problems in computational molecular biology is multiple sequence alignment (MSA) of DNA. Many methods have been proposed to solve MSA problem such as dynamic programming and heuristic. A method has been proposed by Althaus *et al.* to solve MSA problem which is based on integer linear programming (ILP). The general ILP formulation of the MSA is derived from the graph representation of the MSA problem. Although we have the general ILP formulation of the MSA problem, constructing the ILP model of an MSA that can be solved directly using an ILP solver is not straightforward. We develop a MATLAB program that can generate and solve the ILP model of an MSA problem. The method that is used to solve the ILP model is branch-and-bound. The constructed program can generate the ILP model of any given MSA problem but can only solve an MSA problem of a small number of short DNA sequences. The result of the program is the aligned sequences of the MSA problem.

**Keywords**: MSA problem of DNA sequences, gapped extended alignment graph, gapped trace, ILP, branch-and-bound method.

# ABSTRAK

Salah satu dari masalah-masalah dominan pada komputasi biologi molekuler adalah penyejajaran barisan berganda (*Multiple Sequence Alignment* - MSA) dari DNA. Banyak metode yang telah diajukan untuk menyelesaikan masalah MSA seperti pemrograman dinamik dan heuristik. Satu metode telah diajukan oleh Althaus *et al.* untuk menyelesaikan masalah MSA yang didasarkan pada pemrograman linear bilangan bulat (*Integer Linear Programming* - ILP). Formulasi ILP umum dari masalah MSA diturunkan dari representasi graf dari masalah MSA. Walaupun formulasi ILP umum dari masalah MSA diketahui, membentuk model ILP dari suatu masalah MSA yang dapat diselesaikan langsung menggunakan suatu *solver* ILP tidaklah mudah. Sebuah program yang dapat membangun dan menyelesaikan model ILP dari sebuah masalah MSA menggunakan MATLAB telah dibuat. Metode yang digunakan untuk menyelesaikan model ILP tersebut adalah branch-and-bound. Program yang telah dibuat dapat menghasilkan model ILP dari sembarang masalah MSA yang diberikan tetapi hanya dapat menyelesaikan masalah MSA dari sejumlah kecil barisan-barisan DNA yang pendek. Hasil dari program tersebut adalah penejajaran barisan-barisan DNA dari masalah MSA yang diberikan.

**Kata kunci**: masalah MSA barisan-barisan DNA, *gapped extended alignment graph*, *gapped trace*, ILP, *branch-and-bound method*.

xiv + 115 p.

Bibliografi: 8 (1981 - 2006)

# CONTENTS

# LIST OF FIGURE

x

# LIST OF TABLE

# LIST OF APPENDIX