



UniMAP

**REDUCING THE SEARCH SPACE AND TIME COMPLEXITY
OF NEEDLEMAN-WUNSCH ALGORITHM (GLOBAL
ALIGNMENT) AND SMITH-WATERMAN ALGORITHM
(LOCAL ALGORITHM) FOR DNA SEQUENCE ALIGNMENT**

by

FATIMAH NONI BINTI MUHAMAD

(0530210062)

A thesis submitted in fulfillment of the requirements for the degree of
Master of Science (Computer and Communication Engineering)

**SCHOOL OF COMPUTER AND COMMUNICATION
ENGINEERING**

UNIVERSITI MALAYSIA PERLIS

2011

UNIVERSITI MALAYSIA PERLIS

DECLARATION OF THESIS

Author's full name : Fatimah Noni Binti Muhamad
Date of birth : 15 April 1977
Title : Reducing The Search Space And Time Complexity Of Needleman-Wunsch Algorithm (Global Alignment) And Smith Waterman Algorithm (Local Alignment) For DNA Sequence Alignment
Academic Session : 2010 / 2011

I hereby declare that the thesis becomes the property of Universiti Malaysia Perlis (UniMAP) and to be placed at the library of UniMAP. This thesis is classified as :

- CONFIDENTIAL** (Contains confidential information under the Official Secret Act 1972)*
- RESTRICTED** (Contains restricted information as specified by the organization where research was done)*
- OPEN ACCESS** I agree that my thesis is to be made immediately available as hard copy or on-line open access (full text)

I, the author, give permission to the UniMAP to reproduce this thesis in whole or in part for the purpose of research or academic exchange only (except during a period of _____ years, if so requested above).

Certified by:

SIGNATURE

SIGNATURE OF SUPERVISOR

770415 – 06 – 5924

(NEW IC NO. / PASSPORT NO.)

Assoc. Prof. Dr. R. Badlishah bin Ahmad

NAME OF SUPERVISOR

Date : _____

Date : _____

NOTES : * If the thesis is CONFIDENTIAL or RESTRICTED, please attach with the letter from the organization with period and reasons for confidentiality or restriction.

Acknowledgements

First and foremost, I am greatly thankful to Allah s.w.t with the completed of this thesis, Alhamdulillah....

I would like to express my sincere gratitude to my thesis advisor, Assoc. Prof. Dr. R. Badlishah bin Ahmad, Assoc. Prof. Dr. Puteh binti Saad, and Puan Salina binti Mohd. Asi, this thesis could not have been completed without their help and guidance.

In regards to the writing of this thesis, a major thanks goes to my dearest husband, Muhamad Nasir bin Murad, and my beloved son, Muhammad Umar Mursyid bin Muhamad Nasir, without their support and happiness, I would not be where I am today.

I also would like to acknowledge the constant support and help from my family and friends at Embedded Cluster, without which, none of this would have been possible. Appreciation also goes to MOSTI (Ministry of Science, Technology & Innovation, Malaysia) for the 2 years scholarship.

Finally, I would like to gratefully acknowledge the panels for the viva voce, Dr. Ong Bi Lynn and Assoc. Prof. Dr. Norita Mohd. Norwawi, and also the staff at the Postgraduate Studies, UNIMAP.

Thank you very much.

TABLE OF CONTENTS

	1
	AGE
Declaration of Thesis	i
Acknowledgement	ii
Table of Contents	iii
List of Tables	viii
List of Figures	ix
List of Equations	xii
Abstrak	
Abstract	
CHAPTER 1 INTRODUCTION	
1.1 Introduction	1
1.2 Problem Statement	5
1.3 Objective	9
1.4 Summary	9
CHAPTER 2 LITERATURE REVIEW	
2.1 Biological Concept	11
2.1.1 Deoxyribonucleic Acid (DNA)	12
2.1.2 Chromosomes	17
2.1.3 Proteins	19

2.1.4	Ribonucleic Acid (RNA)	19
2.1.5	Data Representation	20
2.2	Bioinformatics	22
2.2.1	Scope and Application of Bioinformatics	25
2.3	DNA Sequence	26
2.3.1	Sequence Databases	27
2.3.2	DNA Sequence in FASTA Format	30
2.4	Sequence Comparison	31
2.4.1	Sequence Alignment	32
2.4.1.1	Pairwise Sequence Alignment	33
2.4.1.2	Multiple Sequence Alignment	34
2.4.2	Sequence Alignment and Similarity Search	35
2.4.2.1	Match and Mismatch	40
2.4.2.2	Gaps	41
2.4.2.3	Scoring Scheme	43
2.4.2.4	Optimal Alignment in Sequence Alignment	47
2.5	Sequence Alignment Technique	49
2.5.1	Heuristic	49
2.5.2	Linear Programming	51
2.5.3	Dynamic Programming	52
2.6	Sequence Alignment Algorithm	55
2.6.1	Global Alignment	55
2.6.1.1	Introduction of Needleman-Wunsch Algorithm	56
2.6.1.1.1	Time Complexity of Needleman-Wunsch	58
2.6.1.1.2	Space Complexity of Needleman-Wunsch	59
2.6.2	Local Alignment	60

2.6.2.1	Introduction of Smith-Waterman Algorithm	61
2.6.2.1.1	Time Complexity of Smith-Waterman	63
2.6.2.1.2	Space Complexity of Smith-Waterman	64
2.7	Parallel Computing	65
2.7.1	Parallel Execution Model	67
2.7.2	Types of Parallel Programming	70
2.7.3	OpenMP	72
2.7.3.1	OpenMP <i>Pragmas</i>	73
2.8	A Survey of Alignment Tools	76
2.8.1	BLAST (Basic Local Alignment Search Tools) – The Similarity (Searching and Homology) Tool	76
2.8.2	FASTA - The Similarity (Searching and Homology) Tool	77
2.8.3	CLUSTAL W	78
2.8.4	EMBOSS	79
2.9	Summary	83

CHAPTER 3 METHODOLOGY

3.1	Algorithms	84
3.1.1	Needleman_Wunsch Algorithm for <i>Needle</i> Program	88
3.1.1.1	Dot Plot Matrix of <i>Needle</i> Program	89
3.1.1.2	Scoring Matrix of <i>Needle</i> Program	93
3.1.1.3	Traceback of <i>Needle</i> Program	101
3.1.1.4	Optimal Alignment of <i>Needle</i> Program	108
3.1.1.5	Time Complexity of <i>Needle</i> Program	111
3.1.1.6	Space Complexity of <i>Needle</i> Program	111
3.1.2	Smith-Waterman Algorithm for <i>Smith</i> Program	111
3.1.2.1	Initialization of <i>Smith</i> Program	113
3.1.2.2	Scoring Matrix of <i>Smith</i> Program	116

3.1.2.3	Traceback of <i>Smith</i> Program	121
3.1.2.4	Optimal Alignment of <i>Smith</i> Program	128
3.1.2.5	Time Complexity of <i>Smith</i> Program	128
3.1.2.6	Space Complexity of <i>Smith</i> Program	129
3.1.3	Parallelization of <i>Needle</i> Program	129
3.1.3.1	Dot Plot Matrix of Parallel Program	131
3.1.3.2	Scoring Matrix of Parallel Program	132
3.1.3.3	Traceback of Parallel Program	133
3.2	Code Development	136
3.2.1	<i>Needle</i> Program	136
3.2.2	<i>Smith</i> Program	136
3.2.3	<i>Needle_Parallel</i> Program – Parallelization of <i>Needle</i> Program	137
3.2.4	Variables and Functions	137
3.2.5	Input Data Program	137
3.2.6	Initialization Program	138
3.2.6.1	Dot Plot for <i>NeedleArray[i][j]</i> Matrix	138
3.2.6.2	Initialization for <i>SmithArray[i][j]</i> Matrix	139
3.2.7	Scoring Program	140
3.2.7.1	Filling score in <i>Needle_scoring()</i>	141
3.2.7.2	Filling score in <i>Smith_scoring()</i>	143
3.2.8	Traceback Program	146
3.2.8.1	Traceback for <i>Needle</i> Program	146
3.2.8.2	Traceback for <i>Smith</i> Program	151
3.2.8.3	<i>Needle_Parallel</i> Program (Parallelization of Needleman-Wunsch Algorithm)	156
3.2.8.3.1	Defining OpenMP Library and Commands	156
3.2.8.3.2	Dot Plot for <i>Needle_Parallel</i> Program	157

	3.2.8.3.3	Scoring for <i>Needle_Parallel</i> Program	158
	3.2.8.3.4	Traceback for <i>Needle_Parallel</i> Program	160
3.3		Summary	164

CHAPTER 4 RESULTS AND DISCUSSION

4.1		Experiment Result	167
	4.1.1	Optimal Alignment Measurement	170
		4.1.1.1 Analysis on Similarity	173
		4.1.1.2 Analysis on Gap	174
		4.1.1.3 Analysis on Mismatch	177
		4.1.1.4 Analysis on Execution Time	179
		4.1.1.4.1 250 Length of DNA Sequences	179
		4.1.1.4.2 3000 Length of DNA Sequences	180
		4.1.1.4.3 Between 100 to 5000 Length of DNA Sequences	181
4.2		Summary	185

CHAPTER 5 CONCLUSION

5.1		Conclusion	186
5.2		Future Work	190

		References	191
		Appendix A	197
		Appendix B	201
		Appendix C	204
		Appendix D	207
		Appendix E	212

LIST OF TABLES

NO.	TITLE	PAGE
2.1	Simple scoring scheme.	44
2.2	Scoring scheme for chemical similarities.	45
2.3	Globally aligned sequences.	56
2.4	Locally aligned sequences.	61
3.1	A general rule to calculate the score for Needleman-Wunsch algorithm.	96
3.2	A general rule to calculate the score for Smith-Waterman algorithm.	117
4.1	Summary of optimal sequence alignment.	172
4.2	Summary of measurement data of similarity for the four methods.	173
4.3	Summary of measurement data of gap for the four methods.	175
4.4	Summary of measurement data of mismatch for the four methods.	177
4.5	Summary of measurement data of execution time (sec) on 250 length of sequences for the four methods.	179
4.6	Summary of measurement data of execution time (sec) on 3000 length of sequences for the four methods.	180

LIST OF FIGURES

NO.	TITLE	PAGE
1.1	Sub-sequences of DNA for Western Gorilla.	3
2.1	The basic nucleotide.	14
2.2	The base nucleotide pairings.	15
2.3	DNA has a structure like a spiral staircase.	16
2.4	The structure of Bioinformatics study.	24
2.5	Subsequences of DNA for the western gorilla genome retrieved from the GenBank database in FASTA format (NCBI, 2008a).	26
2.6	An example of DNA sequence in FASTA format for <i>Homo Sapiens</i> (human) (NCBI, 2008c).	31
2.7	BLOSUM62 substitution (scoring) matrix for proteins (S.F. Altschul, 1991).	46
2.8	A problem is broken into discrete parts.	67
2.9	Distributed Memory Model (cluster) (S. Ho, 2008).	68
2.10	Shared Memory Architecture (ARB, 2007).	69
2.11	The structure of Message Passing Interface, OpenMP – Shared Memory and SIMD (ARB, 2007).	72
2.12	OpenMP uses the <i>fork-join</i> parallelism model (S. Ho, 2008).	75
3.1	Techniques using Dynamic Programming method for Needleman-Wunsch and Smith-Waterman algorithm.	86
3.2	Flow chart of Dot Plot matrix.	91
3.3	Matches and mismatches of the sequences filled in the dot plot matrix.	93

3.4	Flow chart of returning the maximum score in the matrix, <i>NeedleArray[i][j]</i> .	98
3.5	The scoring matrix, <i>NeedleArray[i][j]</i> of <i>Needle</i> Program.	100
3.6	Flow chart of finding the highest scoring position.	103
3.7	Flow chart of trace matches in the matrix.	104
3.8	Flow chart of inserting gap in <i>AlignmentB</i>	106
3.9	Flow chart of inserting gap to <i>AlignmentA</i> .	107
3.10	The traceback of the alignment.	108
3.11	Flow chart of appending letters into <i>AlignmentA</i> and <i>AlignmentB</i> .	109
3.12	The optimal sequence alignment for global alignment.	110
3.13	The result for sequence A and sequence B.	110
3.14	Flow chart to create 0 value in <i>SmithArray[i][0]</i> and <i>SmithArray[0][i]</i> .	115
3.15	Flow chart to fill scores in <i>SmithArray[i][j]</i> matrix.	119
3.16	The score matrix for the alignment of two sequences.	120
3.17	Flow chart of finding maximum score in matrix.	123
3.18	Flow chart of finding maximum score in sub row and column of matrix.	124
3.19	Flow chart of trace match in the matrix and inserting gap in <i>MaxB</i> .	125
3.20	Flow chart of inserting gap in <i>MaxA</i> .	126
3.21	The traceback of the local alignment.	127
3.22	The optimal sequence alignment for local alignment.	128
3.23	Overall parallelization summary for <i>Parallel</i> program that run based on the Needleman-Wunsch algorithm.	130
3.24	The Dot Plot summary for <i>Parallel</i> program	132
3.25	The scoring summary for <i>Parallel</i> program.	133
3.26	The traceback summary for <i>Parallel</i> program.	134
4.1	Execution time for the length of DNA sequences from 100 to 500.	182
4.2	Execution time for the length of DNA sequences from 1000 to 5000.	183

LIST OF EQUATIONS

NO.	TITLE	PAGE
2.1	Time complexity of Needleman-Wunsch algorithm.	59
2.2	Time complexity of the Smith-Waterman algorithm.	64
3.1	To create a matrix for two sequences of DNA.	89
3.2	To create a matrix based on Needleman-Wunsch algorithm.	90
3.3	Filling scores in a global alignment matrix.	95
3.4	Scoring matrix for <i>Needle</i> program.	95
3.5	To initialize the <i>SmithArray</i> [i][j] matrix.	113
3.6	The cost of optimum alignment for <i>SmithArray</i> [i][j] matrix.	116
3.7	The basic score calculation for local alignments.	116

Mengurangkan Kompleksiti Ruang Carian Dan Masa Bagi Algoritma Needleman-Wunsch (Penjajaran Global) Dan Algoritma Smith-Waterman (Penjajaran Tempatan) Untuk Penjajaran Urutan DNA

Abstrak

Prosedur asas bagi menganalisis kandungan urutan ialah melalui perbandingan urutan. Perbandingan urutan boleh ditakrifkan sebagai permasalahan untuk mencari bahagian mana yang sama dan bahagian mana yang berbeza, iaitu dengan membandingkan dua urutan dan mengenalpasti persamaan serta perbezaan di antara mereka. Pendekatan khas untuk menyelesaikan masalah ini adalah dengan mencari kesejajaran yang baik dan wajar antara dua susunan. Penyelidikan utama dalam projek ini adalah untuk menyesuaikan susunan DNA dengan menggunakan algoritma Needleman-Wunsch bagi penjajaran global dan algoritma Smith-Waterman untuk penjajaran tempatan dengan berpandukan kepada algoritma Pengaturcaraan Dinamik. [Algoritma Pengaturcaraan Dinamik](#) adalah dijamin dalam mencari keselarasan yang optimum dengan mengeksplorasi semua kemungkinan penjajaran dan memilih yang terbaik melalui teknik *scoring* dan *traceback*. Algoritma yang dicadangkan dan dinilai adalah untuk mengurangkan *gap* dalam menyelaraskan susunan mahupun panjang urutan tanpa mengorbankan kualiti atau kesahihan hasilnya. Bagi projek ini, satu kajian tentang bagaimana mengaplikasikan kekuatan pengkomputeran bagi Pengkomputeran Selari untuk mempercepatkan proses perbandingan yang panjang tanpa harus berkompromi dengan hasil yang optimum. Keselarian tersebut hanya diterapkan pada algoritma Needleman-Wunsch. Untuk mengesahkan ketepatan dan konsistensi pengukuran adalah diperolehi dari data yang dibandingkan diantara algoritma Needleman-Wunsch dan Smith-Waterman dengan Emboss (*global*) dan Emboss (*local*) bagi 600 data uji. Hasil kajian juga menunjukkan bahawa program *Needle* dan *Smith* dapat mengurangkan *gap* dan *mismatch*, tetapi tidak menjejaskan ketepatannya. Manakala hasil bagi keselarian algoritma Needleman-Wunsch menunjukkan bahawa keselarian hanya berkesan untuk urutan DNA yang panjangnya 3000 dan keatas, tetapi tidak menunjukkan sebarang peningkatan pada urutan DNA yang panjangnya kurang dari 500 walaupun menggunakan platform teras berganda.

Reducing The Search Space And Time Complexity Of Needleman-Wunsch Algorithm (Global Alignment) And Smith-Waterman Algorithm (Local Algorithm) For DNA

Sequence Alignment

Abstract

The fundamental procedure of analyzing sequence content is sequence comparison. Sequence comparison can be defined as the problem of finding which parts of the sequences are similar and which parts are different, namely comparing two sequences to identify similarities and differences between them. A typical approach to solve this problem is to find a good and reasonable alignment between the two sequences. The main research in this project is to align the DNA sequences by using the Needleman-Wunsch algorithm for global alignment and Smith-Waterman algorithm for local alignment based on the Dynamic Programming algorithm. The Dynamic Programming Algorithm is guaranteed to find optimal alignment by exploring all possible alignments and choosing the best through the *scoring* and *traceback* techniques. The algorithms proposed and evaluated are to reduce the *gaps* in aligning sequences as well as the length of the sequences aligned without compromising the quality or correctness of results. In this project, a study on how to apply the computational power of Parallel Computing to speed up the lengthy process of comparing sequences without having to compromise on the optimal results. Parallelization is only applied to the Needleman-Wunsch algorithm. In order to verify the accuracy and consistency of measurements obtained in Needleman-Wunsch and Smith-Waterman algorithms the data is compared with *Emboss (global)* and *Emboss (local)* with 600 strands test data. Results show that the *Needle* and *Smith* programs are reduced *gaps* and *mismatch*, but do not affect the accuracy. Results on the parallelization of Needleman-Wunsch algorithm shows that the parallelization is only efficient for 3000 length of DNA sequences and above, but does not show any improvement for less than 500 lengths of DNA sequences although using multiple core platforms.

CHAPTER I

INTRODUCTION

1.1 Introduction

Nowadays, the current efforts in molecular biology are producing an abundance of biological data stored in numerous databases that spread all over the world. This wealth of data and the complexity of biological processes provided exciting opportunities for new knowledge discovery especially in Bioinformatics. Bioinformatics or Computational Biology refers to an emerging, interdisciplinary field in computer technology, including software, hardware and algorithm, are applied to solve problems arising in biology study. One research area of particular interest in the field is to develop tools for processing biomolecular data (C.Y. Chang, J.T.L Wang & R.K Chang, 1998).

Bioinformatics is the application of information technology to the field of molecular biology. Basically, it uses techniques and concepts from informatics, statistics, applied

mathematics, chemistry, biochemistry, physics, and linguistics to solve biological problems, typically through computer programs and mathematical models.

Research in bioinformatics includes method development for storage, retrieval, and analysis of the data. The major efforts in the field include sequence analysis, gene finding, genome assembly, protein structure alignment, protein structure prediction, prediction of gene expression and protein interactions and the modeling of evolution. Algorithmic development is an important part of bioinformatics, because, techniques and algorithms were specifically developed for the analysis of biological data, for example, the Dynamic Programming algorithm for sequence alignment, and the Machine Learning algorithm for data mining from micro arrays.

The biological or biomolecular data include *protein*, DNA (*deoxyribonucleic acid*), and RNA (*ribonucleic acid*). However, the data representation in this project is DNA sequences. DNA in a genome is not neatly arranged and stores an organism's genetic information in the form of a long sequence of molecules. Specifically, the information is encoded using four key chemicals, *adenine*, *thymine*, *guanine* and *cytosine* (abbreviated as A, T, G and C) (S.K. Moore, 2000).

An example of a genome sequence is shown in Figure 1.1. The DNA sequences of hundreds of organisms have been decoded and stored in databases. This biological sequence data can be obtained from variety of public and private databases.


```
CCTTCATCTAGGAGTTGAGAAGGGTAGATAAGATTCTTGGATACTAGGTATTTAAGAACTTTCTCAG
ATGAAAGGAAGCTGGGAACAAAAGTAAGAAAAGAATACCTTTTAGGATTCACAAAATTATGAGAAAGTCA
GCCACATACGGTAGGTCAGCTTTTTTAATGTATTTGCTCCTTTTCTTATTCATTGACCTGTGAAAGGA
AAATATCTTGGACCTCCAAAATCACTAAGAAAACCTCAAGCTGGATATTGCTTAGGGCAAACCTGCCT
ACCATTCATTTCAAAGTCACAGAGCATTACAAGAAGCGTATTATTATTATGCTTAAGTCAGAATATT
```

Figure 1.1 : Sub-sequences of DNA for Western Gorilla.

With the growing amount of data, it became impractical to analyze DNA sequences manually, so faster algorithms and tools are needed. Sequence analysis is the process used to find information about a *nucleotide* or *amino acid* sequence using computational methods (T.J. Vision & A. McLysaght, 2003). Common tasks in sequence analysis are comparing of sequences in order to find similar and dissimilar in compared sequences (sequence alignment), identifying gene-structures, determining the similarity of two genes, determining the protein coded by a gene, and determining the function of a gene by finding a similar gene in another organism with a known function. This work concentrates on efficient developing and evaluating algorithms for determining the similarity of two genes.

The fundamental procedure of analyzing sequence content is sequence comparison. Sequence comparison is the cornerstone of Bioinformatics. Sequence comparison is regarded as one of the most fundamental problems of computational biology, which is usually solved with a technique known as sequence alignment. Sequence alignment can be defined as the problem of finding which parts of the sequences are similar and which parts are different. Generally, it is the process of comparing two sequences to identify similarities and differences between them. So, a measure of how similar they are is also desirable, then,

a typical approach to solve this problem is to find a good and plausible alignment between the two sequences.

Dynamic Programming algorithm for sequence alignment is used (Needleman & Wunsch, 1970). Dynamic Programming algorithm is one of the major strategies for designing algorithms. There are many algorithms written that use the approach of Dynamic Programming. However, Needleman-Wunsch algorithm was the first to introduce Dynamic Programming to compare biological sequences for finding the global alignment between two sequences (Needleman & Wunsch, 1970). Later, the improvement from Needleman-Wunsch algorithm proposed Smith-Waterman algorithm to find the best local alignment between two sequences (Smith & Waterman, 1981). The Needleman-Wunsch and Smith-Waterman algorithms developed are to solve the sequence alignment problems that have computational complexities. In this project, the study on how to analyze large sequences and to reduce the search space and time complexity without compromising the accuracy and efficiency is presented. This is by evaluating the performance of Needleman-Wunsch and Smith-Waterman algorithms in finding the optimal alignment between a pair of DNA sequences.

Since the nature of Dynamic Programming requires to store and compute an $m \times n$ matrix, many researchers spending most of the effort to reduce calculation time. So, Parallel Computers are also used to reduce the computation time (S. Rajasekaran , V. Thapar, H. Dave & C. H. Huang, 2008). This project also applied the power of parallel computers to speed up the process of comparing sequences without having to compromise the optimal

result. However, this works only focused on the parallel process for Needleman-Wunsch algorithm.

As the database of proteins grows larger, faster algorithms become more important to be able to quickly compare a given sequence to the entire database. The concern on both computational speed and memory management had become a great concern in searching, matching and analyzing. High performance computing such as parallel computers are used, in order to speedup the processes. Researchers also worked on issues in memory management.

1.2 Problem Statement

There has been a tremendous thrust in Bioinformatics research over the last decade. Being a relatively new discipline, Bioinformatics offers a large number of challenging problems. Following are the problems arising from Bioinformatics :

- i. Many available algorithms and techniques in solving the problems of sequence alignment.

There are many algorithms that maximize speed and do not concern with the accuracy of the result alignment. And also, there are many algorithms that maximize accuracy and do not concern with the speed. Most current sequence comparison methods used in practice, such as, BLAST (S.F. Altschul, T.L. Madden, J. Zhang,

et. al, 1997) and FASTA (D.J Lipman. & W.R Pearson, 1985) are based on Heuristics (Michalewicz & Fogel, 2000) which are much faster, but do not provide optimal results. Although sequence comparison algorithms based on the Dynamic Programming method, such as Needleman-Wunsch and Smith-Waterman, provide optimal solutions, but tend to have very high computational complexities.

ii. The problem of finding optimal alignment.

Another problem in the field of Bioinformatics for sequence alignment is the comparison of two sequences of biological data, such as DNA sequences or proteins, and finding the optimal alignment for the sequences compared. Moreover, the comparison of a sequence against a large set of sequences, such as a large database and effort to search for the most similar sequences in the set is important. In order to apply Dynamic Programming to the optimal alignment problem, the problem structure must exhibit the properties, which the problem can be broken into subproblems and the optimal solution of the problem can be efficiently computed from the optimal solution of the subproblems. The goal is to produce the best alignment for a pair of DNA or protein sequences (represented as strings of characters). A good optimal alignment has zero or more gaps / mismatch inserted into the sequences to maximize the number of positions in the aligned strings that match. For example, consider aligning the sequences “ATTGGC” and “AGGAC”. By inserting gaps (“-”) in the appropriate place, the number of positions where the two sequences agree can be maximized: “ATTGG-C” and “A-GGAC” (K. Charter, J. Schaeffer & D. Szafron, 2000).

iii. The computational and space complexity of the sequence alignment algorithms.

The other main problem in sequence alignment is making sequences having the same size with the insertion of gaps in locations along the sequences and creating a correspondence between sequences. The sequence alignment problem can be illustrated as, given a scoring function that measures the score of aligning characters at the same position from each sequence, calculate the total score of the alignment by adding the scores of all positions and find the maximum total score of every possible alignments. The effort is still being studied by researchers in finding the best way to align sequences by reducing the gaps and sizes of the sequences.

iv. The effort to determine the degree of similarity for the optimal sequence alignment.

One of the problems in the comparison of sequences of biological data is an effort to determine their degree of similarity. Algorithms, such as the Smith-Waterman algorithm, are used to compare two sequences allowing for genetic mutations such as insertions, deletions, or substitutions. It is also useful to compare multiple sequences, however, the number of operations involved in non-heuristic algorithms is proportional to the product of the lengths of the sequences, and as a result, require a great deal of processing time. It is useful to determine the evolutionary history or phylogenic tree of a given a set of sequences. This problem is also computationally intensive. The non-heuristic algorithms also have an exponential time complexity based on the length of the protein, and as a result, it cannot be solved in a reasonable

amount of time for even medium length proteins. However, the accuracy in sequence alignment is compromised because many algorithms (including the Dynamic Programming) require human intervention while they are optimizing results. This intervention will have to be done by biologists who are very familiar with the data and thus the usage of such an algorithm is limited biologist availability. Often, it is not necessary to find the most accurate alignment among the sequences.

v. Improving the speed for the sequence alignment algorithms.

A search of a biological sequence database can take a great deal of time for two main reasons. First, sequence alignment algorithms have a significant time complexity because they must allow for genetic mutations, such as character insertions, deletions, or substitutions. Second, the size of a sequence data can be very large, some on the order of billions of characters. This biological sequence data can be obtained from a variety of public and private databases.

1.3 Objective

The objectives of this project are :

- a) To study and compare techniques in DNA sequence alignment.
- b) To develop the optimal sequence alignment algorithm for two sequences of DNA by finding the optimal alignment in much less time and space complexity.
- c) To design scoring functions without using the substitution matrices (e.g. PAM or BLOSSUM) in solving the computational complexities efficiently and for improving the performance of space and time complexity.
- d) To evaluate the accuracy and the complexity of optimal sequence alignment.
- e) To evaluate the effect on parallelization for Needleman-Wunsch algorithm.

1.4 Summary

Bioinformatics derives knowledge from computer analysis of biological data. Research in bioinformatics includes method development for storage, retrieval, and analysis of the data. The main kinds of information stored in biological databases are DNA, RNA

and proteins sequences. The data representation in this project is DNA sequences. As biological databases grow in size, faster algorithms and tools are needed. However, the field of bioinformatics consists of many computationally challenging problems, many of which involve very complex system. This work emphasizes the sequence alignment problems, where a variety of computational algorithms have been applied to the sequence alignment problems, including slow but formally optimizing methods like Dynamic Programming and efficient Heuristic or Probabilistic methods designed for large-scale database search.

In order to apply Dynamic Programming to the optimal alignment problem, the problem structure must exhibit the properties, which the problem can be broken into subproblems and the optimal solution of the problem can be efficiently computed from the optimal solution of the subproblems. It can be seen that the problem of finding optimal alignment obeys the properties and hence Dynamic Programming can be applied very fruitfully. Thus, the problem exhibits optimal substructure and Dynamic Programming can be applied by using Needleman-Wunsch algorithm for global alignment, Smith-Waterman algorithm for local alignment and parallelized the Needleman-Wunsch algorithm in finding the optimal alignment in much less time and space complexity.

In Chapter 2, the studies have been carried out in more detail with respect to sequence alignment algorithms and comparing the efficiencies of the various algorithms and see what sacrifices the algorithms make in exchange for speed and in reducing the search space.