

**UNIVERSITI TEKNOLOGI MARA**

**DEOXYRIBONUCLEIC ACID (DNA)  
FRAGMENTATION ACCELERATION  
TECHNIQUES USING LINEAR  
DIMENSIONAL ARRAY (1D) ON  
SMITH WATERMAN AFFINE GAP  
PENALTY (SWAGP) AND VITERBI  
PAIR HIDDEN MARKOV MODELS  
(VPAIRHMMS)**

**NUR FARAH AIN BINTI SALIMAN**

Thesis submitted in fulfillment  
of the requirements for the degree of  
**Master of Science**

**Faculty of Electrical Engineering**

June 2018

## ABSTRACT

Computational biology is rapidly pushing the advancement bioinformatics field which involves a wide range of areas, including the Deoxyribonucleic Acid (DNA) sequence alignment similarity region. The technique that involve in similarity region is sequence alignment for finding the similarity region between sequences. This technique the of the DNA sequence alignment gaining close attention due great impact through the various area such as in biological data assembly, comparative biological data, drug design, finding homology and specific sequences and evolutionary biology trees. With the increased number of DNA database, it causes the finding for a similar sequence over large biological database is unable to produce results within a realistic time. Hence, for acceleration improvement over acceleration technique that take account the accuracy, sensitivity, speed and size of architecture. These studies proposed an acceleration technique over the Smith-Waterman Affine Gap Penalty (SWAGP) and Viterbi pair Hidden Markov Models (VpairHMMs) based on Field Programmable Gate Array (FPGA). In order to facilitate clearer of the designs, SWAGP and VpairHMMs development is divided into four stages. First, the SWAGP and VpairHMMs algorithm are started with the theoretical verification. Secondly, the development based on Software Version (SV) on FPGA for both algorithms. Next, optimization has been implemented for accelerate both algorithms by upgrade the SV into Custom Instruction (CI). Finally, the algorithms are designed and implemented on hardware accelerator into Linear Dimensional Array (1D) acceleration technique. Twelve (12) tests for each design (SV, CI and 1D) with ranged from 1 until 2048 length base-pair was conducted. The cores has been designed using Verilog HDL and implemented on DE2-115 FPGA (EP4CE115F29C7). The results show that the proposed structural design archived runtime performances of between 8 until 51 percent (SWAGP) and about 7 percent (VpairHMMS) for SV against CI. Next, SV against 1D showed up that the average speed up results range is between 11 until 76 percent (SWAGP) and 10 until 11 percent (VpairHMMS). The CI case against 1D showed that the speed up results range is between 4 until 52 percent (SWAGP) and 3 until 11 percent (VpairHMMS). In conclusion, the runtime depends on the computational method that being used such as software, hardware and system design task implementation.

## ACKNOWLEDGEMENT

All praises due to ALLAH, the most gracious, the most merciful. I would like to express my sincere thanks to those who have involved and supported me towards the successful completion of my Master study.

First and foremost, my deepest gratitude goes to my supervisor, Associate Professor Zulkifli Abd. Majid and co-supervisor Dr. Syed Abdul Mutalib Al-Junid of the Faculty Electrical Engineering at Universiti Teknologi MARA (UiTM). Thank you very much indeed, the door to Associate Professor Zulkifli Abd. Majid and Dr. Syed Abdul Mutalib Al-Junid office that was always open whenever I ran into a trouble spot, had a question about my research or writing, for suggestion, criticism and guidance. This invaluable efforts, patient and kindness to increasing my knowledge especially in hardware design, technical writing and more.

Secondly, this appreciation goes to my labmate or teammate Mrs. Nur Dalilah Ahmad Sabri who has studied together about the fundamental parts of bioinformatics field, sequence alignment, hardware description language, and hardware debugging skills, C programming and software design task debugging skills. Without her passionate participation, the research project could not have been successfully done. In addition, this appreciation also goes Mr. Christopher Lee as the reader of this thesis which I am really gratefully indebted to his very valuable comments on this thesis.

Last but not least, I must express my very profound gratitude to my husband Mr. Luqman Hakim Mat Sout, to my lovely son Mr. Muhammad Ar Rafif Luqman Hakim, to my parents (Mr. Saliman Abu Denan and [REDACTED]) and sibling (Miss Nurul Natasha Saliman and Mr. Muhammad Afieq Saliman) for providing me with unfailing support and continuous encouragement throughout my years of study and through the process of researching and writing this thesis. This accomplishment would not have been possible without them. Thank you.

# TABLE OF THE CONTENTS

	<b>Page</b>
<b>CONFIRMATION BY PANEL OF EXAMINERS</b>	<b>ii</b>
<b>AUTHOR`S DECLARATION</b>	<b>iii</b>
<b>ABSTRACT</b>	<b>iv</b>
<b>ACKNOWLEDGEMENT</b>	<b>v</b>
<b>TABLE OF CONTENTS</b>	<b>vi</b>
<b>LIST OF TABLES</b>	<b>x</b>
<b>LIST OF FIGURES</b>	<b>xi</b>
<b>LIST OF ABBREVIATIONS</b>	<b>xiv</b>
<b>CHAPTER ONE: INTRODUCTION</b>	<b>1</b>
1.1 Background	1
1.2 Problem Statement	3
1.3 Research Objectives and Aims	4
1.4 Research Activities	5
1.5 Scope of Work	8
1.6 Thesis Overview	9
<b>CHAPTER TWO: SEQUENCE ALIGNMENT METHODS</b>	<b>11</b>
2.1 Bioinformatics Field	11
2.2 DNA Sequence Alignment	13
2.2.1 DNA Sequence Alignment Methods	16
2.2.1.1 Global Alignment Methods	16
2.2.1.2 Local Alignment Methods	20
2.2.1.3 Viterbi Algorithm	25
2.3 Smith-Waterman Algorithm Work Background	29
2.3.1 Smith-Waterman Affine Gap Penalty Review	30
2.4 Viterbi Algorithm Work Background	32
2.4.1 Viterbi Algorithm Review	33

2.4.1.1	Hidden Markov Models (HMMs)	33
2.4.1.2	Viterbi Pair Hidden Markov Models (VpairHMMs)	35
2.5	DNA Sequence Alignment Acceleration Types	37
2.5.1	Software Acceleration	37
2.5.1.1	Fast Alignment Search Tool (FASTA)	37
2.5.1.2	Basic Local Alignment Search Tool (BLAST)	38
2.5.1.3	HMMER Software	38
2.5.2	Hardware Acceleration	38
2.5.2.1	Field Programmable Gate Array (FPGA)	39
2.5.2.2	Single Instruction Stream Multiple Data Stream (SIMD)	42
2.6	Field Programmable Gate Array Tools	42
2.7	Summary	44
<b>CHAPTER THREE: SMITH-WATERMAN OPTIMIZATION</b>		
<b>TECHNIQUES</b>		
3.1	Research Methodology	45
3.2	The Procedure Verification of SWAGP	46
3.3	The Specification Verification of SWAGP	48
3.4	Software Version (SV) Configuration	48
3.4.1	Properties for Software Version (SV)	49
3.4.2	Theoretical and Result Verification for Software Version (SV)	50
3.4.2.1	Theoretical Verification for Software Version (SV)	51
3.4.2.2	Result Verification for Software Version (SV)	53
3.5	Custom Instruction (CI) Configuration	54
3.5.1	Properties for Custom Instruction (CI)	54
3.5.2	Theoretical and Result Verification for Custom Instruction (CI)	58
3.5.2.1	Theoretical Verification for Custom Instruction (CI)	58
3.5.2.2	Result Verification for Custom Instruction (CI)	59
3.6	Linear Dimensional Array (1D) Configuration	62
3.6.1	A Smith-Waterman Processing Element Array (SWPEA)	64
3.6.1.1	The single PE for SWAGP	64
3.6.1.2	Cascading the Single Processing of the SWAGP	67
3.6.2	Properties for Linear Dimensional Array (1D)	71
3.6.3	Theoretical and Result Verification for Smith-Waterman Processing	