

Ain Shams University
Faculty of Science
Department of Mathematics



On Solving the Gene Silencing Problem via Algorithmic Techniques

A thesis submitted in partial fulfillment of the requirements of the
M.Sc. degree in Computer Science

By

Soha Ibrahim Soliman Ibrahim

B.Sc. Statistics & Computer Science

Supervised by

***Prof. Dr. Fayed Fayek
Mohamed Ghaleb***

Ass. Prof. Emeritus of Mathematics,
Department of Mathematics,
Faculty of Science,
Ain Shams University

***Dr. Mohammad Hashim Aly
Abdel Rahman***

Lecturer of Computer Science,
Department of Mathematics,
Faculty of Science,
Ain Shams University

Submitted to:
Faculty of Science,
Ain Shams University,
Cairo - Egypt.
2018.

AIN SHAMS UNIVERSITY

Author: Soha Ibrahim Soliman Ibrahim
Title : On Solving the Gene Silencing problem
via Algorithmic Techniques
Division: Computer Science
Department: Department of Mathematics
Faculty: Faculty of Science
Degree: M.Sc.
Year: 2018

Contents

Acknowledgements	v
Abstract	vii
Summary	ix
Publication	xi
List Of Abbreviations	xiii
List of Figures	xv
List of Algorithms	xvii
List of Tables	xix
1 Basic Concepts	1
1.1 Bioinformatics	1
1.1.1 The definition of Bioinformatics	2
1.1.2 The Need for Bioinformatics:	3
1.1.3 Expected Benefits of Bioinformatics	4
1.1.4 ELSI: Ethical, Legal and Social Issues	5
1.2 Algorithmic techniques	5
1.3 Gene Silencing Problem	6
1.4 Objective of The Study	7
2 Algorithmic Techniques Solving Biological Problems	9
2.1 Biological problems	9
2.1.1 DNA restriction mapping problem	10
2.1.2 Motif Finding Problem	12
2.2 Hashing Technique	14
2.2.1 Introduction to Hashing	15
2.2.2 Methodology	15

2.2.3	Advantages of Hashing	17
2.2.4	Disadvantages of Hashing	17
2.2.5	Strategies of Collision resolving:	17
2.2.6	Load Factors	22
3	Gene Silencing	23
3.1	Some Biological Terms	23
3.1.1	Organism	23
3.1.2	Cell	25
3.1.3	Genome	25
3.1.4	DNA	27
3.1.5	Nucleotides	28
3.1.6	Genes	30
3.1.7	RNA	31
3.1.8	Proteins	32
3.2	Gene Silencing Definition	32
3.3	The need for Gene Silencing	33
3.4	Mechanisms Of Gene Silencing	34
3.5	Types of Gene Silencing	35
3.5.1	Transcriptional Gene Silencing (TGS)	35
3.5.2	Post-transcriptional gene silencing (PTGS)	36
3.6	RNA interference (RNAi)	36
4	Exogenous and Endogenous Controls of particular Gene	41
4.1	Some important notations	41
4.1.1	$l - mer$	41
4.1.2	Hamming distance	42
4.1.3	Mismatch Tolerance	42
4.1.4	Complement of DNA sequence	42
4.1.5	Reverse Complement(rc) of DNA sequence	42
4.2	Benefits of controlling a particular gene	43
4.2.1	Endogenous Control of Particular Gene	44
4.2.2	Exogenous Control of Particular Gene	45

4.3	Problem Statement	46
4.4	Related Work	48
4.4.1	Smith-Waterman Algorithm	48
4.4.2	Naive algorithm	49
4.4.3	BYP method	49
4.4.4	SOS-Hash Algorithm	50
4.4.5	SOS-A Solution Based on Radix Sorting . . .	51
4.4.6	Elhamy Algorithm	52
5	EGSH Algorithm	55
5.1	Preliminary Calculations	55
5.2	The Algorithm Description	57
5.2.1	Hash Table Description	57
5.2.2	First phase: Input Phase	59
5.2.3	Second Phase	60
5.2.4	Third Phase	62
5.2.5	Forth Phase	63
5.2.6	Final Phase	64
5.3	The Algorithm and its Analysis	64
5.3.1	EGSH Algorithm	65
5.3.2	Algorithm Analysis	66
5.4	Experimental Results and Discussion	67
5.4.1	Data Sets	67
5.4.2	Experimental Results	67
5.4.3	Discussion	68
6	Conclusion and Future Work	71
	References	73
A	Program Full Code	79
A.1	Class CompareGene	79
A.2	Class DNAUtility	80
A.3	Class GenerateHash	81

A.4 Class HashTable 84

A.5 Class HashTableIsFullException 92

A.6 Class Main 92

A.7 Class ReadFile1 96

Acknowledgements

First of all, cordial thanks to **ALLAH**, from start to end, for enabling me to do and finish this study successfully.

I would like to express my deepest thanks and gratitude to **Prof. Dr. Fayed Fayek M. Ghaleb**, *Prof. of Mathematics, Department of Mathematics, Faculty of Science, Ain Shams University*, for his kind supervision, valuable advice, continuous encouragement, support, and for his great efforts in reviewing the present manuscript.

Great appreciations and thanks are due to **Dr. Mohammad H. Abdel Rahman**, *Lecturer of Computer Science, Department of Mathematics, Faculty of Science, Ain Shams University*, for suggesting and planning the present point of research, and for his valuable suggestions and fruitful advice throughout the present study, for his kind supervision, guidance and participation in reviewing the current manuscript.

Sincere thanks are due to the staff members and my colleagues in the *Department of Mathematics, Faculty of Science, Ain Shams University*, for their cooperation and encouragement during the present study.

All my love and respect are due to my dear parents for their kind support, encouragement, and always praying for me; and to my beloved husband and daughter for the patience, encouragement and missing me most of the time.

Abstract

This thesis introduces a new algorithm for solving one of the hottest problems in biology and medicine which is the gene silencing problem.

The new algorithm is called “Exogenous Gene Silencing using Hashing ” (EGSH). This algorithm is specially designed to solve the exogenous silencing of a specific target gene, taking into consideration the possibility of both exact and partial matching between the target gene and small interfering RNA that produced by the new algorithm.

The theoretical analysis of the running time and memory complexity of the EGSH algorithm confirms that the EGSH algorithm achieves a remarkable speeding up for the running time and reducing of the memory space required comparable with other previously introduced algorithms.

Also, the experimental result obtained from implementing the EGSH algorithm (using Java Programming language) emphasizes the theoretical analysis results. It shows that for any specific Human gene, the program takes, in average, about three minutes running time consuming less one Gigabyte memory.

Summary

Some genetic disorder diseases are caused by mutant genes that synthesize harmful proteins causing the initiation and progression of such disorders.

One of the ways to avoid these diseases is to turn off the mutant genes to prevent them from producing the harmful proteins which is known as the Gene Silencing.

The Exogenous gene silencing means designing short interfering RNA sequences in laboratories called siRNA then injecting them into the cell to target a particular messenger RNA(mRNA) of the mutant gene and cause its degradation. The big challenge in this approach is that matching between the designed siRNA and target mRNA do not need to be a perfect matching.

In this thesis, a new algorithm Exogenous Gene Silencing using Hashing (**EGSH**) is introduced to solve the exogenous gene silencing problem. Although few previous algorithms were introduced to solve this problem but they focused on the totally matching between siRNA and target mRNA.

The new feature of the proposed algorithm is it solves both the total and the partial matching problem. The proposed algorithm uses the hashing technique that takes linear execution time and it also requires relatively small memory space compared to the previous algorithms.

This thesis is organized into six chapters:

Chapter 1: Explains the basic concepts of bioinformatics and why scientists use it, and gives a brief summary of algorithmic techniques, gene silencing, and the objective of the study.

Chapter 2: Is divided into two parts. The first part shows the importance of the computational tools like algorithmic techniques in solving some biological serious problems. The second part talks about the hash technique, it's definition, methodology, strategies, hash functions, it's characteristics and the advantage, and disadvantage of using the hash technique.

Chapter 3: Composed of two main parts. The first part illustrates some important biological terms that help to understand the problem of the thesis like Organism, Cell, Genome, Deoxyribonucleic Acid (DNA), Nucleotides, Gene, Ribonucleic Acid (RNA) and Proteins. The second part illustrates all the details of the gene silencing problem such as it's definition, the need for the gene silencing, mechanism of gene silencing, different types of gene silencing, finally the definition of RNA interference (RNAi) and how it works.

Chapter 4: Explains firstly the benefits of controlling a particular gene, this control can be done by one of two ways endogenous or exogenous gene silencing control. Secondly present some of the previous work that tries to solve the problem with it is advantages and disadvantages.

Chapter 5: Introduces in detail a new algorithm called EGSB for the exogenous gene silencing problem of a particular gene, also illustrates it is analysis and the experimental results of the new EGSB algorithm.

Chapter 6: The conclusion and future work; summarizes the results achieved in the thesis including the new algorithm advantages and limitation, and provide conclusion and future work in this field.

Publication

Soliman, S.I.; Ghaleb, F.F.M. and Abdel-Rahman, M.H. (2017): A Hashing Algorithm Improving the Exogenous Gene Silencing Problem. First International Conference of Faculty of Science, Ain Shams University, Hurghada, EGYPT.