

Palestine Polytechnic University



**College of Engineering & Technology
Electrical and Computer Engineering Department**

Software Project

"Twin Gene"

Project Team

**Ala Abu Al_Reesh
Gharam Al-Skafi
Nida Al-Amlh**

**Project Supervisor
Dr.Mahmood Al-saheb**

Hebron – Palestine Polytechnic University

June, 2004

Palestine Polytechnic University



**College of Engineering & Technology
Electrical and Computer Engineering Department**

Software Project

"Twin Gene"

Project Team

**Ala Abu Al_Reesh
Gharam Al-Skafi
Nida Al-Amlh**

**Project Supervisor
Dr.Mahmood Al-saheb**

Hebron – Palestine Polytechnic University

June, 2004

Palestine Polytechnic University

Hebron – Palestine

"TwinGene"

By

Ala Abu-Reesh

Gharam al-Skafi

Nida Al-Amlh

According to the directions of the project supervisor and by agreement all of the committee members, this project was submitted to Department of Electrical and Computer Engineering in College of Engineering and Technology to partially fulfill to the Bachelor requirements for the department.

Supervisor Signature

Name:.....

Dept. Head Signature

Name:.....

Name:.....

Committee Members

1. Name:.....

2. Name:.....

ABSTRACT

Computer science has the characteristic that it is related to all fields of our life, one of these fields is the medical field; computer and medical fields were collected in the last few five years into bioinformatics science, so bioinformatics is the science that studies the relation of computer and medical fields. From here we decided to invest our knowledge in the bioinformatics science, so we decided to work in the field of genes. The main important problem that faces experts in medical field is that they sometimes discover that two genes are the same gene or that one gene has another similar gene that scientists think it is the same gene, so we worked in this project on the alignment of genes that is bringing similar genes for the query gene then writing each gene in one line based on the alignment rules and then dividing the genes into groups based on the correlation in differences with the query gene.

Dedication

To any one appreciates knowledge

To every one helped us

To our families and friends

Ala'

Gharam

Nida'

Acknowledgement

Our sincere thanks and appreciation to the person who gave us his guidance and support to complete this project, specially our instructor Dr. Mahmood Al-saheb and Dr Yaqub Al_Ashab.

Thanks to everyone shared in the success of the project.

Table of Contents

Dedication.....	I
Table of Contents.....	II
Table of Figures.....	VI
List of Tables.....	VII
Acknowledgement.....	VIII
Abstract.....	IX

CHAPTER ONE INTRODUCTIONError! Bookmark not defined.

1.1 Overview..... **Error! Bookmark not defined.**
1.2 Project Importance **Error! Bookmark not defined.**
1.3 What is Twin Gene **Error! Bookmark not defined.**
1.4 Project Objectives **Error! Bookmark not defined.**
1.5 The Report Outline **Error! Bookmark not defined.**

CHAPTER TWO SYSTEM PALNNING ...Error! Bookmark not defined.

2.1 Overview..... **Error! Bookmark not defined.**
2.2 System Risks..... **Error! Bookmark not defined.**
2.3 Development Requirements..... **Error! Bookmark not defined.**
 2.3.1 Hardware Development Resources..... **Error! Bookmark not defined.**
 2.3.2 Software Development Resources **Error! Bookmark not defined.**
2.4 Cost Estimation..... **Error! Bookmark not defined.**
 2.4.1 Development Cost..... **Error! Bookmark not defined.**
2.5 Operational Requirements: 11
2.6 Time Scheduling 11

CHAPTER THREE SYSTEM REQUIREMENTSError! Bookmark not defined.2

3.1 Overview..... 13
3.2 Requirement Definition 13

3.2.1 Functional Requirements Definition.....	Error! Bookmark not defined.	3
3.2.2 Non-Functional Requirements	Error! Bookmark not defined.	4
3.3 Functional Requirements Analysis		14
3.3.1 Import Query Sequence	Error! Bookmark not defined.	4
3.3.2 Displaying Genes Headers.....	Error! Bookmark not defined.	5
3.3.3 Making Multiple Alignment	Error! Bookmark not defined.	5
3.3.3 Making Global Alignment	Error! Bookmark not defined.	8
3.3.4 Making single alignment	Error! Bookmark not defined.	
3.3.5 Dividing the aligned EST's into groups	Error! Bookmark not defined.	
3.3.6 Counting the number groups and the number of EST's in each group.....	Error! Bookmark not defined.	
3.3.7 Calculating the % of intersected nucleotide	Error! Bookmark not defined.	
3.3.8 Draw some statistical graphs	Error! Bookmark not defined.	0
3.4 Functional Requirements Specification	Error! Bookmark not defined.	1
3.4.1 Import the Query.....	Error! Bookmark not defined.	1
3.4.2 Displaying genes headers.....	Error! Bookmark not defined.	2
3.4.3 Single Alignment	Error! Bookmark not defined.	3
3.4.4 Multiple Alignments	Error! Bookmark not defined.	4
3.4.4 Global Alignments	Error! Bookmark not defined.	5
3.4.5 Grouping of Genes.....	Error! Bookmark not defined.	

CHAPTER FOUR SYSTEM DESIGN**Error! Bookmark not defined.**

4.1 Overview.....	Error! Bookmark not defined.	
-------------------	-------------------------------------	--

4.2 System Flow Charts	Error! Bookmark not defined.
4.3 Design Models	Error! Bookmark not defined.
4.3.1 Data Flow Diagram(DFD)	Error! Bookmark not defined.
4.3.2 Flowcharts.....	Error! Bookmark not defined.1
4.3.2.1 Multiple Alignment Flowchart	Error! Bookmark not defined.1
4.2.2.2 Grouping Flowchart.....	Error! Bookmark not defined.4
4.3.3 Structure Chart	Error! Bookmark not defined.5
4.4 Data Design.....	36
4.4.1 Arrays.....	36
4.4.2 Records	37
4.4.3 Files.....	37
4.5 User Interface Design	Error! Bookmark not defined.1

CHAPTER FIVE IMPLEMENTATION ..Error! Bookmark not defined.6

5.1 Overview.....	47
5.2 Building the System.....	47
5.3 Prepare the Platform of the System	47
5.3.1 Install and Configure Windows XP professional	48
5.3.2 Install Visual Basic.NET	48
5.3.3 Install BLAST Program	48
5.3.3.1 Blast System recommendations:.....	Error! Bookmark not defined.2
5.4 Multiple Alignment Implementations.....	Error! Bookmark not defined.2

5.5 Single Alignment Implementation.....	58
5.6 Grouping Implementation.....	69

CHAPTER SIX TESTING.....Error! Bookmark not defined.2

6.1 Introduction.....	Error! Bookmark not defined.3
6.2 Testing Procedures.....	Error! Bookmark not defined.
6.3 Testing Strategies.....	Error! Bookmark not defined.5
6.3.1 Black Box Testing	Error! Bookmark not defined.5
6.3.2 White Box Testing	Error! Bookmark not defined.6

CHAPTER SEVEN CONCLUSIONS AND FUTURE WORKError! Bookmark not d

7.1 Conclusions.....	Error! Bookmark not defined.1
7.2 Future Work.....	Error! Bookmark not defined.2

CHAPTER ONE INTRODUCTIONError! Bookmark not defined.

1.1 Overview.....	Error! Bookmark not defined.
1.2 Project Importance	Error! Bookmark not defined.
1.3 What is Twin Gene	Error! Bookmark not defined.
1.4 Project Objectives	Error! Bookmark not defined.
1.5 The Report Outline	Error! Bookmark not defined.

CHAPTER TWO SYSTEM PALNNING ...Error! Bookmark not defined.

2.1 Overview.....	Error! Bookmark not defined.
2.2 System Risks.....	Error! Bookmark not defined.
2.3 Development Requirements.....	Error! Bookmark not defined.
2.3.1 Hardware Development Resources.....	Error! Bookmark not defined.
2.3.2 Software Development Resources	Error! Bookmark not defined.
2.4 Cost Estimation.....	Error! Bookmark not defined.
2.4.1 Development Cost.....	Error! Bookmark not defined.
2.5 Operational Requirements:	11
2.6 Time Scheduling	11

CHAPTER THREE SYSTEM REQUIREMENTSError! Bookmark not defined.2

3.1 Overview.....	13
3.2 Requirement Definition	13
3.2.1 Functional Requirements Definition.....	Error! Bookmark not defined.3
3.2.2 Non-Functional Requirements	Error! Bookmark not defined.4
3.3 Functional Requirements Analysis	14
3.3.1 Import Query Sequence	Error! Bookmark not defined.4
3.3.2 Displaying Genes Headers.....	Error! Bookmark not defined.5
3.3.3 Making Multiple Alignment	Error! Bookmark not defined.5
3.3.3 Making Global Alignment	Error! Bookmark not defined.8

3.3.4 Making single alignment	Error! Bookmark not defined.
3.3.5 Dividing the aligned EST's into groups	Error! Bookmark not defined.
3.3.6 Counting the number groups and the number of EST's in each group.	Error! Bookmark not de
3.3.7 Calculating the % of intersected nucleotide	Error! Bookmark not defined.
3.3.8 Draw some statistical graphs	Error! Bookmark not defined.0
3.4 Functional Requirements Specification	Error! Bookmark not defined.1
3.4.1 Import the Query.....	Error! Bookmark not defined.1
3.4.2 Displaying genes headers.....	Error! Bookmark not defined.2
3.4.3 Single Alignment	Error! Bookmark not defined.3
3.4.4 Multiple Alignments	Error! Bookmark not defined.4
3.4.4 Global Alignments	Error! Bookmark not defined.5
3.4.5 Grouping of Genes.....	Error! Bookmark not defined.

CHAPTER FOUR SYSTEM DESIGN.....Error! Bookmark not defined.

4.1 Overview.....	Error! Bookmark not defined.
4.2 System Flow Charts	Error! Bookmark not defined.
4.3 Design Models	Error! Bookmark not defined.
4.3.1 Data Flow Diagram(DFD)	Error! Bookmark not defined.
4.3.2 Flowcharts.....	Error! Bookmark not defined.1
4.3.2.1 Multiple Alignment Flowchart	Error! Bookmark not defined.1
4.2.2.2 Grouping Flowchart.....	Error! Bookmark not defined.4
4.3.3 Structure Chart.....	Error! Bookmark not defined.5

4.4 Data Design.....	36
4.4.1 Arrays.....	36
4.4.2 Records	37
4.4.3 Files.....	37
4.5 User Interface Design	Error! Bookmark not defined.1

CHAPTER FIVE IMPLEMENTATION ..Error! Bookmark not defined.6

5.1 Overview.....	47
5.2 Building the System.....	47
5.3 Prepare the Platform of the System	47
5.3.1 Install and Configure Windows XP professional	48
5.3.2 Install Visual Basic.NET	48
5.3.3 Install BLAST Program.....	48
5.3.3.1 Blast System recommendations:.....	Error! Bookmark not defined.2
5.4 Multiple Alignment Implementations.....	Error! Bookmark not defined.2
5.5 Single Alignment Implementation.....	58
5.6 Grouping Implementation.....	69

CHAPTER SIX TESTING.....Error! Bookmark not defined.2

6.1 Introduction.....	Error! Bookmark not defined.3
6.2 Testing Procedures.....	Error! Bookmark not defined.

6.3 Testing Strategies.....	Error! Bookmark not defined.	5
6.3.1 Black Box Testing	Error! Bookmark not defined.	5
6.3.2 White Box Testing	Error! Bookmark not defined.	6
References	Error! Bookmark not defined.	4

CHAPTER ONE

INTRODUCTION

1.1 Overview

1.2 Project Importance

1.3 What is Twin Gene

1.4 System Objectives

1.5 The Report Outline

1.1 Overview

Computer science is related to all sciences and to all aspects of our life; its relation to biology science is strengthened by Bioinformatics which is defined as conceptualizing biology in terms of molecules and then applying informatics techniques (derived from disciplines such as applied math, computer science and statistics) to understand and organize information associated with these molecules [1].

1.2 Project Importance

As the number of the projects that deal with the organization of biological information have grown explosively in the recent years the Bioinformatics research appeared which is a newly emerging interdisciplinary research area and can be defined as the interface between biological and computational sciences. Thus, the people working in this field in most cases either have a training in biology or computer science, and they learned about the other field by dealing with problems or using the tools of the other one.

Twin Gene System allows the user to make some processing on a query gene such as alignment with similar genes, alignment with one gene and determining the groups of the genes contained in one database of genes.

1.3 What is Twin Gene

Twin Gene system can be defined as the system that divides the genes into groups. It accepts a huge database taken from Gene Bank and with complex comparisons, it finds all the similar EST's (Expressed Sequence Tags) for the query Gene (using public program called BLAST). The BLAST output should be converted into a multiple alignment (MA) file. A multiple alignment arranges a set of sequences in a scheme where positions believed to be homologous (similar) are written in a common column, then it process the MA file to determine the Twin Gene (divided into groups based on correlation of differences with the query gene).

1.4 Project Objectives

1. Learn how to design and implement twin gene in a system software project.
2. To connect the biological field to the computer science.
3. To explore and evaluate the phenomenon of pseudogene expression in Human Genome by checking the prevalence of paralogous EST's that are misannotated to one gene.
4. To provide solutions for designing molecular techniques that fit genes with paralogous copy with high sequence similarity.
5. The identification of new unknown copy or copies of a given gene will lead us to more accurate mapping of the different genes in the Human Genome.

1.5 The Report Outline

This report shows all important points related to the project, this documentation helps the reader to get comprehensive view about the project.

We tried to present a documentation that can be understood easily by dividing the documentation to a set of subject related parts.

The report consists of a number of chapters. Each chapter discusses a subject related to the project, and these chapters are sorted according to the logical relations of subjects.

Chapter One: “Introduction”; gives overview about the project, project importance, and Twin Gene system objectives.

Chapter Two: “System planning”; presents system risks, development requirements, cost estimation, operational Requirements and time scheduling

Chapter Three: “System Requirements”; it includes definition, analyses and specification.

Chapter Four: “System Design”; consists of the system design model which includes the data flow diagram, structural chart and the data and control analysis.

Chapter Five: “System Implementation”; presents general view; prepares the platform of the system.

Chapter six:” System Testing”; shows how the system is going to be tested.

Chapter seven: “Future Work and Conclusions”; which is in the last chapter, it shows the suggestions for the future development, and points out the conclusions.

CHAPTER TWO

SYSTEM PALNNING

2.1 Overview

2.2 System Risks

2.3 Development Requirements

2.4 Cost Estimation

2.5 Operational Requirements

2.6 Time Scheduling

2.1 Overview

This chapter talks about system risks, hardware and software requirements and time scheduling.

2.2 System Risks

1. This is the first project in this field so we have not previous projects to help us.
2. Because this project contains biomedical concepts, so we need specialists in this field to help and teach us many concepts.
3. The size of the output file is so large because the database is so large and may reach to nine compact CD's or higher, as a result we may display we display samples of the output.

2.3 Development Requirements

Development Requirements includes hardware resources, software resources and the human resources.

2.3.1 Hardware Development Resources

As system needs complex searching and comparisons to find a Twin Gene, we need a professional language to help our for the complex comparisons that we need for our project and in this project we use Visual Basic.NET , so we need a computer that can implement this system; this computer should have the following:

Requirement	Required	Recommended
Processor	PC with a Pentium II-class processor, 450 MHz	Pentium III-class, 600MHz
Ram	Windows XP Professional: 160 MB	192 MB for XP Professional
Hard disk	2 Giga	10Giga
Video card	800 x 600, 256 colors	High Color 16-bit

Table 2.1 Hardware development resources [2].

2.3.2 Software Development Resources

The PC computer should contain the following operating system and software:

Operating system	Windows XP professional
Programs	Visual studio.net (2002)

Table 2.2 Software development resources [2].

2.4 Cost Estimation

This part contains the estimation of the development costs and the implementation costs.

2.4.1 Development Cost

Hardware Costs:

Hardware Component	Cost
Petium3_600Mhz	4000\$

Table (2.3): Development Hardware Costs

Software Costs:

Software Component	Cost
Windows XP	299 \$
Visual Basic .NET	1,794\$
Total	2093

Table (2.4): Development Software Costs

Human Resources Costs:

Number of Developers	Work Hours/day (120 days)	Cost per hour for one Developer	Total
3	5	20\$	36000\$

Table (2.5): Development Human Costs

Hardware Dev Costs	Software Costs	Human Dev. costs	Total
400\$	2093\$	36000\$	38493\$

Table (2.6): Total development Costs

2.5 Operational Requirements:

Required Processor	Recommended Processor	Required RAM	Recommended RAM	Programs	OS
Pentium 90 MHz	Pentium 90 MHz or faster	32 MB	96 MB or higher	Microsoft Internet Explorer 5.01 .Net framework	Win 98, win me or higher

Table 2.7 Operational Requirements [2]

2.6 Time Scheduling

Our project has to be developed during 17 weeks . figure 1.1 shows the time plan of the project .

Task	week	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
System definition & planning.		■	■	■	■	■											
Requirements , specification,& analysis				■	■	■	■										
System Design						■	■	■	■	■	■	■					
Implementation									■	■	■	■	■	■	■	■	■
Testing											■	■	■	■	■	■	■
Documentation		■	■	■	■	■	■	■	■	■	■	■	■	■	■	■	■

Table 2.8: Time Scheduling

INDEIECES

INDEX A

Theoretical Background

List of Tables

Table 2.1 Hardware development resources	8
Table 2.2 Software development resources	8
Table 2.3 Development Hardware Costs	9
Table 2.4 Development Software Costs	9
Table 2.5 Development Human Costs.....	10
Table 2.6 Total development Costs	10
Table 2.7 Operational Requirements.....	10
Table 2.8 Time Scheduling.....	11
Table 6.1 Testing Level.....	77
Table 6.2 Black Box Testing	77
Table 6.3 White box testing.....	80

List of Figures

Figure 3.1 Alignment cases 1.a	16
Figure 3.2 Alignment cases 1.b	16
Figure 3.3 Alignment cases 1.c.....	16
Figure 3.4 Alignment cases 2.a.a	17
Figure 3.5 Alignment cases 2.a.b	17
Figure 3.6 Alignment cases 2.a.c.....	17
Figure 3.7 Alignment cases 2.a.d	18
Figure 3.8 Alignment cases 2.a.e.....	18
Figure 4.1 System Flow Charts.....	29
Figure 4.2 Data Flow Diagram of the system.....	31
Figure 4.3 Flowchart for MA.....	32
Figure 4.4 Flowchart of the Grouping system	34
Figure 4.5 Structure Chart for the system	35
Figure 4.6 Query Gene	37
Figure 4.7 Example of BLAT Output	39
Figure 4.8 Global Alignment	40
Figure 4.9 Multiple Alignment File.....	41
Figure 4.10 Interface for welcoming the user	42
Figure 4.11 The main interface.....	43
Figure 4.12 Interface for MA.....	44
Figure 4.12 Interface for SA	45
Figure 5.1 Example of BLAT output	51
Figure 6.1 White Box Testing	77

DNA (deoxyribonucleic acid) and proteins are biological macromolecules built as long linear chains of chemical components. In the case of DNA these components are the so-called nucleotides, of which there are four different ones, each denoted by one of the letters A, C, G and T.

DNA	adenine	guanine	cytosine	thymine
	A	G	C	T/U
RNA	adenine	guanine	cytosine	uracil

So the nucleotides can be made up of 20 different amino acids (or "residues") which are denoted by 20 different letters of the alphabet.

	One-letter code	Three-letter-code	Name
1	A	Ala	Alanine
2	C	Cys	Cysteine
3	D	Asp	Aspartic Acid
4	E	Glu	Glutamic Acid
5	F	Phe	Phenylalanine
6	G	Gly	Glycine
7	H	His	Histidine
8	I	Ile	Isoleucine
9	K	Lys	Lysine
10	L	Leu	Leucine
11	M	Met	Methionine
12	N	Asn	Asparagine
13	P	Pro	Proline
14	Q	Gln	Glutamine
15	R	Arg	Arginine
16	S	Ser	Serine
17	T	Thr	Threonine
18	V	Val	Valine
19	W	Trp	Tryptophan
20	Y	Tyr	Tyrosine

To compare between different nucleotides they insertions into the chain and deletions from the chain. The elementary operations allowed in the definition of sequence similarity are chosen to correspond to these events. To visualize the relationship between two similar sequences they are represented in the form of an alignment:

V-LSPADKTNVKAAWGKVGAGHAGEYGAEALERMFLSFP'TTKTYFPHF-DL	HAHU
VHLTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDL	HBHU
SH-----GSAQVKGHGKKVADALTNVAHVDDMPNALSALSDLHAHKLRV	HAHU
STPDAVMGNPKVKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHV	HBHU
DPVNFKLLSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTISKYR	HAHU
DPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHKYH	HBHU

The two amino acid sequences compared here are the alpha chain of human hemoglobin (abbreviated HAHU) and its beta chain (HBHU). With the sequences being approximately 150 amino acids long, each block of lines contains part of the first sequence in the upper and of the second sequence in the lower line. Residues on top of each other in one block are equivalenced. Some residues are conserved (the amino acids in the column are identical), some have been exchanged and part of the chain has been deleted from the one sequence or (equivalently) inserted in the other. Insertions or deletions are indicated by a letter paired with a dash, the gap-character. An alignment can also be interpreted as representing the operations necessary to transform a sequence into another one using the same operations as evolution does.