

**SSR MARKER DEVELOPMENT IN SORGHUM
&
PHYLOGENETIC STUDIES IN CEREALS.**

*Dissertation Submitted In Partial Fulfillment Of
Requirement For The Award Of Degree Of*

**MASTER OF TECHNOLOGY
in
BIOTECHNOLOGY**

By

HEMABINDU CHINTALA



CENTRE FOR BIOTECHNOLOGY
Institute of postgraduate studies & research
Jawaharlal Nehru Technology University
Hyderabad-500028

20002

**TO
LORD
SHRI SHIRIDHI SAI BABA**



ICRISAT

International Crops Research Institute for the Semi-Arid Tropics

Patancheru 502 324
Andhra Pradesh
India



CGIAR

Tel +91 40 3296161 (19 lines)
Fax +91 40 241239
+91 40 3296182
Email ICRISAT@CGIAR.ORG

CERTIFICATE

This is to certify that the work reported in the dissertation entitled “**SSR MARKER DEVELOPMENT IN SORGHUM & PHYLOGENETIC STUDIES IN CEREALS**” Submitted by **CH.Hemabindu** have been carried out under my supervision. This work is towards the partial fulfillment of her **M.Tech** Degree from **Jawaharlal Nehru Technological University**, Hyderabad. This work is original and has not been submitted in part or full for any other degree or diploma of any university

[Dr. V. MAHALAKSHMI]
Principal Scientist,
GT-1,
ICRISAT.

Visit our worldwide web site at <http://www.icrisat.org>

ICRISAT is part of the global agricultural research network called the Consultative Group on International Agricultural Research (CGIAR)



JAWAHARLAL NEHRU TECHNOLOGICAL UNIVERSITY
CENTRE FOR BIOTECHNOLOGY

Institute of Post Graduate Studies and Research

Mahaveer Marg, Hyderabad - 500 028, A.P., India

Phone No. : 040-3373020

Dr. M. LAKSHMI NARASU

Associate Professor & Head

CERTIFICATE

This is certified that the work reported in the dissertation entitled "**SSR MARKER**

DEVELOPMENT IN SORGHUM & PHYLOGENETIC STUDIES IN CEREALS"

Submitted in partial fulfillment for the award of **M.Tech** in biotechnology from

Jawaharlal Nehru Technological University, Hyderabad, is a bonafied work carried out

by **Ms CH. HEMEBINDU** under the guidance of **Dr.V.Mahalakshmi**,

Senior scientist, International Crops Research Institute For Semi-Arid Tropic,**(ICRISAT)**

[Dr.M.Lakshmi Narasu]

DECLARATION

I Hemabindu.Ch, a bonafied student of IPGSR, JNTU, Hyderabad here by declare that the dissertation entitled “**SSR Marker Development In Sorghum And Phylogenetic Studies In Cereals**” is solely done by me under the expertise guidance of Dr.V.Mahalakshmi at International Crops Research Institute For Semi-Arid Tropics (ICRISAT), Hyderabad.

The facts and figures enumerated in this project work are in accordance with the results of the modeling done in computer. This project work has not been submitted to any university or institution for the award of any degrec or diploma.

ACKNOWLEDGEMENTS

This project work is carried out with valuable suggestions and guidance under the supervision of Dr. V. Mahalakshmi, senior scientist, GREP, ICRISAT, Patancheru, Hyderabad. Her unstilted encouragement and deep concern helped me to complete my project work in time. I am highly grateful and indebted to her.

It gives me immense pleasure in expressing my deep sense of gratitude to Dr. Lakshmi Narasu, Head, Centre for biotechnology, Jawaharlal Nehru University, Hyderabad. I take this opportunity to thank her for suggesting me to join in the ICRISAT as an Apprentice and helped me in finishing my project work.

My sincere thanks to Dr. Prameela Devi, Dr. Archana Giri, associate professors and Mr. Kiran, Ms. Anuradha and Ms. Radhika, academic associates, Centre for biotechnology, JNTU, for their valuable guidance.

I would like to thank my friends Ms. Manjula, Ms. Leela, Mr. Kumar, Mr. PVNS Prasad, Ms. Rekha and Mr. A. Balakrishna and other colleagues in ICRISAT for their encouragement and cooperation during my course of project work.

I show my gratitude to my beloved mother, father, brother's and sister & my Husband for their constant encouragement and good support throughout this course of work.

HEMABINDU.CH

ABSTRACT

Bioinformatics, the application of computational techniques to analyze the information associated with bio-molecules on a large scale and encompasses a wide range of subject areas from structural biology, genomics to gene expression studies.

Chapter 1 gives an introduction and overview of the SSR marker development in sorghum . For this purpose bioinformatics tools like Tandem Repeats Finder, Primer3, Windows software and MS Access were used. More than 50,000 records were collected and placed in the INTRANET of ICRISAT, which could be accessed by the scholars and scientists for their requirements.

A tandem repeat in DNA is two or more contiguous approximate copies of a pattern of nucleotides. Extensive knowledge about pattern size, copy number, mutational history, etc. for tandem repeats has been limited by the inability to easily detect them in genomic sequence data. These sequences range in size from 3kb up to 700kb. A World Wide Web server interface at http://c3.biomath.mssm.edu/trf/advanced_submit.html has been established for automated use of the program.

Primer design is crucial for the success of PCR. Inappropriate primers cause low yield, equivocal results and misinterpretation. An ideal primer should only bind with its 3' end to a unique sequence.

Chapter 2 is dealt with an introduction about Phylogenetic studies and describes completely about the phylogeny of selected cereals for conserved enzymes. . For this purpose we utilized the software tools like CLUSTAL W for multiple alignment, JALVIEW for alignment analysis and phylogenetic tree construction, AND PRIMER 3 for primer designing which is crucial for PCR. Phylogeny is about evolution and is used to reconstruct evolutionary events. It is now possible to construct phylogenetic evolution at a molecular level through analysis of molecular sequences, namely proteins & nucleic acids. To construct phylogenetic tree among grass family, the sequences of conserved enzymes from mitochondria, chloroplast and nucleus are probed using bio-informatics tools.

CONTENTS

1. SSR Marker Development For Sorghum Data Base1
1.1 Introduction1
1.2 Different Tools For Repeat Finder.....	...2
1.2.1 Censor_xxx_Humrep [Censor].....	...3
1.2.2 Repeatmasker_Xxx Primate [Repmask].....	...3
1.2.3 SST_Xxx_HumRep [SST].....	...3
1.2.4 Xnun_Repeat_Default [XNUN].....	... 3
1.2.5 TANDEM_Xxx_DEFAULT [TANDEM].....	...3
1.2.6 Inverted_Xxx_Default [INVERTED].....	...3
1.3 Tandem Repeats Finder4
1.4.1 Basic.....	...4
1.4.2 Intermediate.....	...4
1.4.3 Advanced.....	...4
1.5 Advanced Tandem Repeat Finder Program Parameters6
1.5.1 Alignment Parameters.....	...6
1.5.2 Minimum Alignment Score.....	...6
1.5.3 Maximum Period Size.....	...6
1.5.4 Detection parameters.....	...6
1.5.5 Options6
1.5.5.1 Flanking sequence.....	...6
1.5.4.2 Masked sequence File6
1.5.4.3 Data File.....	...6
1.6 Procedure For Finding Tandem Repeats.....	...8
1.7 Entering the sequence for Finding Tandem Repeats...8
1.7.1 Sequence.....	...8
1.7.2 Fasta Format.....	...8

1.7.3 Submit sequence.....	9
1.8 Table Explanation.....	10
1.9 Alignment Explanation.....	11
1.10 Primer Design.....	12
1.10.1 Introduction.....	12
1.11 Primer Design Programs.....	12
1.12 Primer Design Considerations.....	13
1.13 Features Of Primer Design.....	14
1.14 Limits Of Primer Design.....	14
1.15 Primer Design Parameters.....	15
1.15.1 Primer Length.....	15
1.15.2 Primer Sequence.....	15
1.15.3 GC content.....	16
1.15.4 Melting temperature.....	16
1.15.5 Secondary structure formation.....	17
1.15.6 Specificity.....	17
1.15.7 Primer ends.....	17
1.16 Primer3.....	18
1.17 Procedure For Primer Design By Using Primer3.....	18
1.18 Primer3 Input Parameters.....	19
1.18.1 Source Sequence.....	19
1.18.2 Sequence Id.....	19
1.18.3 Targets.....	19
1.18.4 Excluded Regions.....	19
1.18.5 Product Size.....	19
1.18.6 Number To Return.....	20
1.18.7 Max 3' Stability.....	20
1.18.8 Max Mispriming.....	20
1.18.9 Pair Max Mispriming.....	20
1.18.10 Primer Size.....	20

1.18.11 Primer T_m	21
1.18.12 Maximum T_m Difference	21
1.18.13 Product T_m	21
1.18.14 Primer GC%	21
1.18.15 Max Complementarity	21
1.18.16 Max 3' Complementarity	22
1.18.17 Max Poly-X	23
1.18.18 Included Region	23
1.18.19 Start Codon Position	23
1.18.20 Mispriiming Library	23
1.18.21 CG Clamp	23
1.18.22 Salt Concentration	23
1.18.23 Annealing Oligo Concentration	24
1.18.24 Max Ns Accepted	24
1.18.25 Liberal Base	24
1.18.26 First Base Index	24
1.18.27 Inside Target Penalty	24

1.18.28 Outside Target Penalty.....	24
1.18.29 Sequence Quality.....	25
1.18.30 Min Sequence Quality.....	25
1.18.31 Min 3' Sequence Quality.....	25
1.18.32 Sequence Quality Range Min.....	25
1.18.33 Sequence Quality Range Max.....	25
1.18.34 Penalty Weights.....	25-26
1.18.35 Hyb Oligos (Internal Oligos)	27
1.19 Methods For Finding Tandem Repeats.....	30
1.19.1 Steps Involved After Downloading The Sequence	
Of Sorghum From NCBI.....	30
1.19.1.1 Fasta Format.....	30
1.19.1.2 Enter The Sequence For Finding Tandem Repeats.....	30-31
1.19.1.3 Submit Sequence.....	32
1.19.1.4 Tandem Repeats Report.....	33
1.19.1.5 Summary Table.....	33
1.19.1.6 Alignment Explanation Table.....	34
1.19.1.7 Primer3.....	35-37
1.19.1.7 Out Put Of Primer3.....	38-40
1.20 Results.....	41-48
1.21 Discussion.....	49
2. Phylogenetic studies in cereals.....	50
2.1 Introduction	50
2.2 Phylogenetic terms.....	51-56
2.3 Phylogenetic classifications.....	56
2.4 Methods of phylogenetic analysis.....	57

2.4.1 Cladistic Method.....	57
2.4.2 Phenetic Method.....	57
2.4.3 Multiple Alignment Method.....	58
2.5 Clustalw	58
2.5.1 upload a file	59
2.5.2 Sequences	60
2.5.3 search title	61
2.5.4 CPU mode.....	61
2.5.5 alignment.....	61
2.5.6 Output.....	61
2.5.7 JalView.....	62
2.5.8 out order.....	62
2.5.9 color.....	62
2.5.10 Fast pairwise alignment options.....	62
2.5.11 multiple sequence alignment options.....	62
2.5.12 Gap open.....	63
2.5.13 End gap.....	63
2.5.14 Gapext.....	63
2.5.15 Gapdist.....	63
2.5.16 Phylogenetic Tree.....	63
2.5.17 kimura correction of distances.....	63
2.5.18 Ignore Gaps In Alignment	64
2.6 Phylogenetic Tree.....	64-65
2.7 Methods For The Phylogenetic Studies Of Cereals.....	66
2.8 Reasons For Taking Enzymes In Phylogenetic Studies.....	66
2.8.1 Nuclear Enzyme.....	66
2.8.2 Mitochondrial Enzyme.....	66
2.8.3 Chloroplast Enzyme.....	67
2.9 Multiple Alignment Method.....	67
2.10 Steps Involved In Multiple Alignment Method.....	67
2.10.1Select The Most Conserved Enzymes.....	67

2.10.2 Search The Sequences From NCBI.....	68-69
2.10.3 Accession Number.....	72-80
2.10.4 Exon Regions.....	81
2.10.5 Multiple Alignment.....	81
2.10.6 Run Clustalw.....	82
2.10.7 JalView.....	83-87
2.10.8 Alignment Graph.....	88-93
2.10.9 Phylogenetic Tree.....	93
2.10.10 Tree Analysis.....	94
2.10.11 Design Primers For The Sequences Of Maize.....	94
2.10.12 Selection Of First Set Of Primers.....	95
2.10.12.1 Primer3 For Left Primer.....	95-97
2.10.12.2 Primer3 Output For Left Primer.....	98
2.10.12.3 Primer3 For Right Primer.....	99-101
2.10.12.4 Primer3 Output For Right Primer.....	102
2.10.12.5 Calculation Of Product Size.....	103
2.10.13 Selection Of Second Set Of Primers.....	103
2.10.13 .1 Primer3 For Left Primer.....	103
2.10.13 .2 Primer3 Output For Left Primer.....	106
2.10.13 .3 Primer3 For Right Primer.....	107-109
2.10.13 .4 Primer3 Output For Right Primer.....	110
2.10.13 .5 Calculation Of The Product Size.....	111
2.10.14 Mitochondrial And Nuclear Enzyme.....	111-112
2.11 Results.....	113-114
2.12 Discussion.....	115

WEBSITES USED IN THE PROJECT WORK

- 1) http://www-genome.wi.mit.edu/cgi-bin/primer/primer3_www.cgi
- 2) <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi>
- 3) <http://www.ebi.ac.uk/clustalw/>
- 4) <http://www.ncbi.nlm.nih.gov/>
- 5) <http://c3.biomath.mssm.edu/trf.html>
- 6) <http://c3.biomath.mssm.edu/example.html>
- 7) <http://c3.biomath.mssm.edu/trf.definitions.html#fasta>
- 8) <http://c3.biomath.mssm.edu/trf.submit.options.html>
- 9) <http://c3.biomath.mssm.edu/trf.advanced.submit.html>
- 10) <http://c3.biomath.mssm.edu/trf.upload.form.html>
- 11) <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Nucleotide>
- 12) <http://www2.ebi.ac.uk/~michele/jalview/contents.html>
- 13) <http://www.expasy.ch/enzyme/>
- 14) <http://www.genome.ad.jp/>
- 15) <http://www.mssm.edu/school.html>

LIST OF FIGURES

S.No	Subject	page No.
1	Tandem Repeat Finder Option Page.	5
2	Tandem Repeat Finder Advanced Submit Page.	7
3	Submission Of Sequence To Advance Tandem Repeat Finder.	9
4	Primer3.	28- 29
5	Alignment Explanation Table.	34
6	NCBI Home Page.	70
7	NCBI Nucleotide Page.	71
8	ClustalW Submission Form.	82

BIOINFORMATICS

Introduction:

Bioinformatics is conceptualizing biology in terms of molecules (in the sense of physical-chemistry) and then applying “informatics” techniques (derived from disciplines such as applied math, Computer Science, and statistics) to understand and organize the associated with these molecules, on a large scale.

In the last few decades, advances in molecular biology and the equipment available for research in this field have allowed the increasingly rapid sequencing of large portions of the genomes of several species. In fact, to date, several bacterial genomes, as well as those of some simple eukaryotes (e.g., *Saccharomyces cerevisiae*, or baker's yeast) have been sequenced in full. The Human Genome Project, designed to sequence all 24 of the human chromosomes, is also progressing. Popular sequence databases, such as Genbank and EMBL, have been growing at exponential rates. This deluge of information has necessitated the careful storage, organization and indexing of sequence information. Information science has been applied to biology to produce the field called **Bioinformatics**.

The most pressing tasks in bioinformatics involve the analysis of sequence information. **Computational Biology** is the name given to this process, and it involves the following:

1. Finding the genes in the DNA sequences of various organisms
2. Developing methods to predict the structure and/or function of newly discovered proteins and structural RNA sequences.
3. Clustering protein sequences into families of related sequences and the development of protein models.
4. Aligning similar proteins and generating phylogenetic trees to examine evolutionary relationships.

The Need For Bioinformatics

1. Whole Genome Analyses and Sequences
2. Experimental Analyses involving Thousands of Genes simultaneously
3. DNA Chips and Array Analyses
4. Expression Arrays
5. Comparative Analyses between Species and Strains
6. Proteomics: 'Proteome' of an Organism ... 2D gels, Mass Spec
7. Medical applications: Genetic Disease ... SNPs
8. Pharmaceutical and Biotech Industry
9. Forensic applications
10. Agricultural applications

Evolution Of Bioinformatics

After years of research in structure-function relationships of genes and proteins, the last decade proved to be extremely important and immensely satisfying due to its technical advances in genome sequences of several species and protein structure and identification. To handle this ever-increasing voluminous data, computer processing power and disk storage has been instrumental. Besides gathering all these data, it is necessary to compare these nucleotide and amino acid sequences to find similarities and differences. Since it is not convenient to compare the sequences, that are several hundreds of nucleotides by hand, several computational techniques were developed to approach this problem. In addition these are less error-prone than the manual approach. So bioinformatics has taken its place to cater the needs of biological community.

Divisions Of Bioinformatics

Bioinformatics is a multi disciplinary subject. Though only about a decade old, it has become very important for the growth of biosciences, biotechnology, and the economic prosperity of nations. Three well-identified subdivisions of Bioinformatics are:

- a) Molecular Bioinformatics
- b) Cellular and sub-cellular Bioinformatics and
- c) Organismic and community Bioinformatics.

Out of these three, most Bioinformatics scientists and workers practice molecular Bioinformatics. The other two areas are more recent and are at different stages of development. In the next 5-10 years, cellular and sub-cellular Bioinformatics that will include metabolic pathways, epigenetic, and neuro Bioinformatics on one hand and Bioinformatics of Species Diversity, behavior, evolution and the effect of pollutants on higher as well lower species, on the other will occupy the main stage.

Global Importance Of Bioinformatics

Bioinformatics has acquired great importance due to its recent application in vast amount of data generated in the Genome sequence projects. The target of decoding the three billion base pairs of the human DNA has become achievable only through the use of various innovative techniques and methods evolved by the Bioinformatics scientists. Bioinformatics has become an essential component of biotechnology based product and process development. The process of drug design and development is expensive and time consuming. The application of the tools and techniques of Bioinformatics has resulted in the reduction in cost and the development cycle of the drugs. This aspect has a tremendous impact on the society. If a newly discovered drug is a life-saving one, then the resulting gains are not only in terms of financial savings but also in saving the lives of several million people. Major pharmaceutical and biotechnology companies have set-up large R&D groups in Bioinformatics.

Current research has identified biotechnology the fastest growing sector of production technology. Further advances in this sector will depend quite a lot upon the progress of Bioinformatics and hence there is a great emphasis on Bioinformatics world over. The following are a few important Websites for international networks/ institutions/ groups on Bioinformatics.

The Role Of Bioinformatics In Research

As traditional way of research is very time consuming and laborious, in these days there is introduction of computers and even biological sciences are no exception to this. The long sequences and enormous flow data in molecular biology it becomes difficult to manage data by persons or by organizations. Further with the help of computers through Internet and Intranet, sharing of data with others in the organization or organizations through out the world is becoming routine. This aids in the research and makes it easier for scientists to download, analyze, compare and exchange information. This also saves a lot of time and finances.

Uses Of Bioinformatics

1. The main role of bioinformatics for today's world is to increase productivity and quality of all plants and animals.
2. For understanding evolution, comparative genomics will add hither to unimaginable comprehensiveness to the study of the relationships between species.
3. In cell biology, the components of cellular activity, how these components interact and how they are influenced by environmental states can be identified comprehensively.
4. In medicine, sequences will provide a basis for the study of susceptibility to disease, its pathogenesis and the development of new preventive and therapeutic approaches.
5. Bioinformatics is extremely useful to mankind as it can help in increasing crop yields and animal produces, to cater the needs of ever increasing human

population. It also helps to develop more effective drugs to protect public health as well as to control pests and diseases on the crop.

1. SSR Marker Development In Sorghum

1.1 Introduction

DNA molecules are subject to a variety of mutational events. One of the less well understood is **TANDEM DUPLICATION** in which a stretch of DNA, which we call the pattern, is converted into two or more copies, each following the preceding one in a contiguous fashion. for example we could have

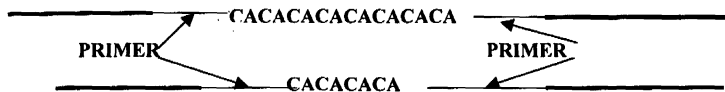
...TCGGA... □ ...TCGGCGGCGGA...

In which the single occurrence of triplet **CGG** has been transformed into three identical, adjacent copies. The result of a tandem duplication event is termed as **TANDEM REPEAT**. Over time, individual copies within a tandem repeat may undergo additional, uncoordinated mutations so that so typically, only approximate tandem copies are present. Tandem repeats are ubiquitous sequence features in both prokaryotic and eukaryotic genomes. These are highly useful as genetic markers. They are codominant, occur in high frequency, and appear to be distributed through out the genomes of most, not all the higher plants and animals. They also display a high level of polymorphism, even among closely related accessions, and are amenable to simple and inexpensive Polymerase Chain Reaction (**PCR**) assays (Brown et. al. 1996).

Tandem repeats are usually classified among satellites (spanning megabases of DNA, associated with heterochromatin), minisatellites (repeat units in the range 6-100 bp, spanning hundreds of base pairs) and microsatellites (repeat units in the range 1-5 bp, spanning a few tens of nucleotides). The minisatellites are also called "various number tandem repeats" or **VNTRs**. The microsatellites are also called "short tandem repeats" or **STRs**. (Short tandem repeats, **STRs**) contain 2-5 bp repeats **VNTR 'S** are scattered at various locations in the Genome are regions that are highly variable. These regions contain a type of DNA sequence called variable number Tandem Repeat. Tandem repeats are multiple copies of sequence of base pairs arranged in head to tail fashion. For example, a frequently found Tandem Repeat is **CA**, and one strand containing this type of repeat reads **CACACA...** notated as **(CA)_n** the other strand would read

GTGTGT.... In this example, the number of repeating base pairs is two, but it can be more. When the repeating unit is less than four, the VNTR is called a **MICROSATELLITE** and when the repeating unit is longer it is a **MINISATELLITE**.

MICROSATELLITE are DNA regions with variable numbers of short tandem repeats flanked by a unique sequence. Microsatellites make good genetic Markers because they each have many different 'alleles'-i. e. There can be many different lengths of the repeat region. The number of repeats there are at the same location defines an allele. With many alleles, most individuals are heterozygous, giving power to note association between marker allele and performance in progeny inheriting a favorable linked QTL allele.



Tandemly repeated sequences are especially liable to undergo misalignments during chromosome pairing, and the size of tandem clusters tends to be highly polymorphic. The smaller clusters of this simple sequence can be used to characterize individual genomes in the technique of "DNA finger printing". Comparisons of corresponding regions of simple sequence DNA with in and between species are informative about the mechanisms involved in manipulating sequences

1.2 Different Tools For Repeat Finder

1.2.1 `cenor_xxx_humrep` [cenor]

This tool searches for genome wide repeats and classifies its findings. Each incarnation of this tool produces two tables containing ALU (xxx=ALU) and non ALU (xxx=NONALU) repeats. In this incarnation the "Humrep" repeat database [humrep] is used and the search sensitivity is set to "moderate". If the input sequence was derived from other mammals than human the "MamRep" repeat database - which comes along with CENSOR too - is used.

1.2.2 repeatmasker_XXX primate [repmask]

This tool searches for both genome wide and simple repeat and classifies its findings. Each incarnation of this tool produces three tables containing ALU (xxx=ALU), non ALU (xxx=NONALU) and simple (xxx=REPEAT) repeats. In this incarnation the "Primate" repeat database is used. The search sensitivity is set to "high" and tagging of simple repeats is turned on. If the input sequence was derived from other mammals than human the "Rodent" or "Vertebrate" repeat database - which comes along with RepeatMasker too - is used.

1.2.3 sst_XXX_humrep [sst]

This tool searches for both genome wide and simple repeats and classifies its findings. Each incarnation of this tool produces three tables containing ALU (xxx=ALU), non ALU (xxx=NONALU) and simple (xxx=REPEAT) repeats. The "HumRep" repeat database [[humrep](#)] is used.

1.2.4 xnun_repeat_default [xnun]

This tool searches for short-period repeats (micro-satellites) and masks its findings. In order to reduce the rate of false positives the default probability cut is lowered from 1 % to 0.1 %.

1.2.5 tandem_XXX_default [tandem]

This tool searches for short-period repeats (micro-satellites) and reports the consensus repeat unit. Each incarnation of this tool produces two tables. The first table (xxx=STRAND) contains all repeats found together with the consensus repeat unit and the strand information. The second table (xxx=REGION) contains just the repeat regions and is cross-linked with the first table.

1.2.6 inVerted_XXX_default [inverted]

This tool searches for inverted repeats within a window of 2 *kb*. The window size is an internally fixed parameter of the analysis tool itself. Each incarnation of this tool produces two tables. The first table (xxx=STRAND) contains two entries for each

inverted repeats. One for the forward and one for the reverse strand repeat unit. The second table (xxx=REGION) contains just the repeat regions and is cross-linked with the first table.

1.3 Tandem Repeats Finder

A Tandem repeat in DNA is two or more adjacent, approximate copies of a pattern of nucleotides. Tandem Repeats Finder is a program to locate and display Tandem Repeats in DNA sequences. In order to use this program, we have to submit the sequence in **FASTA** format. There is no need to specify the pattern, the size of the pattern or any other parameter. The program's analysis is sent back as two files, a summary table file and an alignment file. The summary table contains information about each repeat, including its location, size, number of copies and nucleotide content. Clicking on the location indices for one of the table entries opens a second web browser that shows an alignment of the copies against a consensus pattern. The program is very fast, analyzing sequences on the order of .5Mb in just a few seconds. Submitted sequences may be up to 5Mb in length. Repeats with pattern size in the range from 1 to 500 bases are detected.

1.4 Levels Of Tandem Repeat Finder:

There are 3 LEVELS of tandem repeat finders.

1.4.1 Basic

It uses default parameters (recommended for beginners.)

1.4.2 Intermediate

It provides the parameter Maximum period size, and options Flanking sequence and Masked Sequence File.

1.4.3 Advanced

It provides the parameters like Alignment parameters (match, mismatch and indels), Minimum Alignment Score To Report Repeat, Maximum Period Size and options like Flanking sequence, Masked sequence file, Data file.

TANDEM REPEATS FINDER

Submitting Your Sequence

Please Select one of the following options:

[Beginner](#) Use default parameters (recommended for beginners)

[Intermediate](#) Set some parameters and select options

[Advanced](#) Set all parameters and options



Department of
Mathematical
Sciences

Last revised July 22, 1999
Send any questions or comments to
inf@math.utoronto.ca

TANDEM REPEAT FINDER OPTION PAGE

1.5 Advanced Tandem Repeat Finder Program Parameters:

Input to the program consists of a sequence file and the following parameters:

1.5.1 Alignment Parameters

Weights for match, mismatch and indels. Lower weights allow alignments with more mismatches and indels. Match weight is +2 in all options here. Mismatch and indels weights [interpreted as negative numbers] are 3, 5, or 7. A 3 is more permissive and a 7 is less permissive of these types of alignment choices.

1.5.2 Minimum Alignment Score

The alignment score must meet or exceed this value for the repeat to be reported.

1.5.3 Maximum Period Size

The period size must be no larger than this value for the repeat to be reported. The program will find all repeats with period size between 1 and 500, but the output table can be limited to some other range.

1.5.4. Detection parameters

Matching probability P_m and indel probability P_i . $P_m = .80$ and $P_i = .10$ by default and it cannot be modified in this version of the program.

1.5.5 Options

1.5.5.1 Flanking sequence

Flanking sequence consists of the 200 nucleotides on each side of a repeat. Flanking sequence is recorded in the alignment file. This may be useful for PCR primer determination.

1.5.4.2 Masked sequence File

The masked sequence file is a FASTA format file containing a copy of the sequence with every character that occurred in a tandem repeat changed to the letter 'N'. The word "masked" is added to the sequence description line just after the '>' character.

1.5.4.3 Data File: The data file is a text file, which contains the same information, in the same order, as the summary table file, plus consensus sequences. This file contains no labeling and is suitable for additional processing, for example with a perl script, outside of the program.

TANDEM REPEATS FINDER

Parameters:

Alignment Parameters (match, mismatch, indel)

Minimum Alignment Score To Report Repeat

Maximum Period Size

Options:

Flanking Sequence Masked Sequence File Data File

Sequence:

Your data must be a DNA sequence in FASTA format.
Choose one of the following ways to send your data

- 1 Upload a file from your directory

- 2 Cut and paste sequence

TANDEM REPEAT FINDER ADVANCED SUBMIT PAGE

1.6 Procedure For Finding Tandem Repeats:

Web sites like <http://c3.biomath.mssm.edu/trf/advanced.submit.html> are used to find out the tandem repeats of particular sequence (source sequence). This web site is provided by The Department of Biomathematical sciences, Mount Sinai School of Medicine.

- Download the source sequence from the www.ncbi.nlm.nih.gov/entrez
- Open the page <http://c3.biomath.mssm.edu/trf/advanced.submit.html>
- Enter the source sequence in **cut and paste sequence** blank.
- Click the **submit sequence** button.

1.7 Entering The Sequence For Finding Tandem Repeats:

1.7.1 Sequence:

Data must be a DNA sequence in FASTA format.

1.7.2 FASTA format:

The FASTA format looks like this,

```
gij18063613|gb|BM324719.1|BM324719 PIC1_34_G05.b1_A002 Pathogen-infected
compatible 1 (PIC1) Sorghum bicolor cDNA, mRNA sequence
GCACGAGGCACACACACACTGACGTACGATATCGAACACACTGACACAA
CACAACCGCCGGGGCGCCGGCTTTGTCTGTC AACGTCACCATTCCTTAACA
CACTCTCTACATACACGTGCGAGCCGGCCGGCAGCTCATCTTCTGTATTTAA
ACGCGCTCCATCCATTCCGGTCAGGGCAGGAGTCGTCACCTCCTAGCTCCGC
GCGCACCGAGCTCCCACACTCTCGAGCAATGGCACGCGGGCGGCTCTACC
GCGGCTCCTCCACCCCTCCCGCCGCTGCTGCTGCTGCTGGTGGTAGTAGTGG
TGTCCTCGGCCTTTTCGCTCGCTGCTGCCACGCCGCGCCGGGGTGC ACTG
GTGACGAGCGTCCCGGGGTACACGGCGGCAGCGGCGCTGCCGTCGAAG
CACTACGCCGGGTACGTGACGGTGGACGAGGCGCACGGCAGGAGGCTCTTC
TACTACCTGGTGGAGCCGAGCGTGACCCCGCAAGACCCCGTCGTGCTGTG
GCTCAACGGCGGGCCTGGCTGCTCCAGCTTCGACGGC
```


1.8 Table Explanation:

The summary table includes the following information:

1. Indices of the repeat relative to the start of the sequence.
2. Period size of the repeat.
3. Number of copies aligned with the consensus pattern.
4. Size of consensus pattern (may differ slightly from the period size).
5. Percent of matches between adjacent copies overall.
6. Percent of indels between adjacent copies overall.
7. Alignment score.
8. Percent composition for each of the four nucleotides.
9. Entropy measure based on percent composition.
10. If the output contains more than 140 repeats, multiple linked tables are produced.

47zEZYDucydeZ.2.3.5.80.10.30.10.1.html - Microsoft Internet Explorer

Address: http://zbc.uh.tmc.edu/47zEZYDucydeZ.2.3.5.80.10.30.10.1.html

Gary Benzon
Department of Biomedical Sciences
Baylor College of Medicine
Version 2.00

Please cite:
G. Benzon,
"Tandem repeat finder: a program to analyze DNA sequences"
Nucleic Acid Research(1999)
Vol. 27, No. 5, pp. 573-580.

Sequence: 1000513
Parameters: 2 2 5 00 10 30 10
Length: 556

Tables: 1
This is table 1 of 1
Click on indices to view alignment

Table description

Indices	Period (bp)	Copy Number	Consensus Size	Exact Matches	Percent Indels	Score	A	C	G	T	Entropy (0.2)
9-15	7	230	7	50	1.2	50	42	6	10	10	1.74
100-101	6	37	6	27	0	34	0	31	40	27	1.56
337-345	9	24	9	14	0	34	0	31	40	27	1.56

The End

SUMMARY TABLE

1.9 Alignment Explanation:

The alignment is presented as follows:

1. In each pair of lines, the actual sequence on top and a consensus sequence for all the copies is on the bottom.
2. Each pair of lines is one period except very small patterns.
3. The 10 sequence characters before and after a repeat is shown.
4. Symbol * indicates a mismatch.
5. Symbol - indicates an insertion or deletion.
6. Statistics refers to the matches, mismatches and indels overall between adjacent copies in the sequence, not between the sequence and consensus pattern.

7. Distances between matching characters at corresponding positions are listed as distance, number at that distance, percentage of all matches.
8. A, C, T, G count is percentage of each nucleotide in the repeat sequence.

1.10 Primer Design

1.10.1 Introduction:

Designing PCR and sequencing primers are essential activities for molecular biologists around the world. Primer design was developed to find suitable primers for PCR or oligo nucleotides for probes and DNA sequencing. Primer design is crucial for the success of PCR. Inappropriate primers cause low yield and misinterpretation. Primers that bind to multiple DNA loci can synthesize side products and render sequencing illegible, especially with high amplification of small amounts of DNA and with impure DNA. They are generally the result of short DNA sequence repeats. An ideal primer should only bind to a unique sequence. To ensure this the given sequence must be compared with itself to identify repeats.

Primer Design is a DOS-program to choose primer for PCR or oligonucleotide probes. Napiwotzki, J. and Becker, A. wrote this program in 1995. It is tailored to check known sequences for repeats and unique sequences and subsequently to create proper primers according to this data.

1.11 Primer Design Programs:

- **PrimerGen** searches strings of amino acid residues in order to reverse-translate oligonucleotide primers of a desired range of lengths and maximum number of degeneracies.
- **Primer** (Stanford) Sun Sparcstations only
- **Primer** (Whitehead) Unix, Vms (and DOS and Mac if you can compile it)

- **Amplify**, Bill Engels (Macintosh only), is for use in designing, analyzing, and simulating experiments involving the polymerase chainreaction (PCR). You can obtain your copy of Amplify [here](#)
- **PrimerDesign 1.04**, a DOS-program to choose primer for PCR or oligonucleotide probes. See also the [PrimerDesign Welcome Page](#).
- **PC-Rare**, a new software by R. Griffais, which uses a rare octamer at the 3' termini of the primer. This powerful (but user friendly) software is available for Macintosh and Windows environment.
- **Primer Design**, a Java applet by Luca Ida Giovanni TOLDO.
- **CODEHOP**, PCR primers designed from protein multiple sequence alignments (Local mirror site at WIS).
- **Primer3**, an online service to pick PCR primers from nucleotide sequence. (Local mirror site at WIS).
- **NetPrimer** (PREMIER Biosoft International). NetPrimer combines the latest primer design algorithms with a web-based interface allowing the user to analyze primers over the internet

1.12 Primer Design Considerations:

One of the single most important factors in successful automated DNA sequencing is proper primer design. It is important that a primer has the following characteristics:

1. Primers should be at least **18-20 nucleotides** in length to minimize the chance of encountering problems with a secondary hybridization site on vector or insert.
2. Primers with long runs of a single base should generally be avoided. It is especially to avoid **3** or more G's or C's in a row.
3. For cycle sequencing, primers with melting temperatures above **55°C** are generally produce better results than primers with lower melting temperatures.

4. Primers should have a **G/C** content between **40** and **60** percent. For primers **GC** content of less than **50%**, it may be necessary to extend the primer sequence beyond 18 bases to keep the melting temperature above the recommended lower limit of **55°C**.
5. Primers should be "stickier" on their **5'** ends than on their **3'** ends. A "sticky" **3'** end as indicated by a high **G/C** content could potentially anneal at multiple sites on the template.
6. "**G**" or "**C**" is desirable at the **3'** end.
7. Primers should not contain complementary (palindromes) within themselves, which is they should not form hairpins. If this state exists a primer will fold back on itself and result in an unproductive priming event, which decreases the overall signal obtained.
8. Primers should not contain sequences of nucleotides that would allow one primer molecule to anneal to it self or to the other primer used in a **PCR** reaction (primer dimer formation).
9. If possible, run a computer search against the vector and insert **DNA** sequences to verify that the primers, and especially the 8-10 bases of its **3'** end, are unique.
10. Do not design degenerate primers. Do not request inosine in sequencing primers. They either do not work or give poor cycle sequencing results.

1.13 Features Of Primer Design:

- Creating of new **primer pairs**.
- Creating one suitable primer to a given primer.
- Finding of **repeats** within a sequence.
- Finding of **unique** sequences within a sequence.
- Handling of sequences up to **32,000bp**.

1.14 Limits Of Primer Design:

- The sequence length, which can be used for primer design, repeat and unique search is limited to **32,000bp**.

- Maximal **16000** repeats can be found and sorted.
- Primer combinations can be explored up to **6000 pairs**.

1.15 Primer Design Parameters:

The most critical parameter for successful PCR is the design of primer. The primer sequence determines several things such as the length of the product, its melting temperature and ultimately the yield. A poorly designed primer can result in little or no product due to non-specific amplification or primer dimer formation, which can be competitive enough to suppress product formation. The main parameters for primer design are:

1.15.1 primer length

Specificity, the temperature and the time of annealing are at least partly dependent on primer length. Length of primer is critical parameter in primer designing. It is also proportional to annealing efficiency. For general studies primers of typically 17-34 nucleotides in length are the best primers, <16 nucleotides generally anneal non-target DNA sequences. Longer primers are required when a perfect complementary sequence to the entire template is not found.

1.15.2 primer sequence

- 3' self/cross complementarity** – 3 or more bases of complementarity at the 3' end can give a strong primer dimer and reduce PCR efficiency.
- Runs of single bases** – avoid long runs of single bases (i.e.5 or more).
- Internal self-complementarity** - We should check and avoid sequences of internal self-complementarity otherwise these regions can form hairpin loops.
- Internal sequence structure** – Primer stability is affected by nearest neighbour interactions between bases. Consequently it is better to have Gs and Cs together in pairs rather than dispersed within As and Ts.

1.15.3 GC content:

Melting temperature and annealing temperature are strictly dependent on each other. GC is an important characteristic of DNA and provides information about the strength of annealing (GC have hydrogen bonds between them). The base composition of primer should be between 45% and 55% of GC. The primer sequence must be chosen such that there is no poly G or poly C stretches that can promote non-specific annealing. Poly A, poly T stretches also to be avoided as these will breath and open up stretches of primer template complex.

1.15.4 Melting temperature:

- a) **Matching the primer for melting temperature** – the stability of a primer-template DNA duplex can be measured by its melting **temperature**(T_m). The T_m of any duplex of DNA is defined as the temperature at which **50%** is dissociated as **ssDNA**. T_m is affected by sequence, length, pH , salt concentration etc. The T_m of any primer pair (or multiplexed set) should be reasonably well matched otherwise amplification efficiency may be compromised.

There are many formulas used for calculating T_m the most accurate account of nearest neighbour interactions and are best left to computer programs.

$$T_m(^{\circ}C) = 2(A+T) + 4(G+C)$$

This formula becomes inaccurate for longer primers.

Alternatively another more versatile formula is as follows: -

$$T_m(^{\circ}C) = 69.3 + 0.41(\%G+C) - (650/L)$$

$\%G + C =$ %age GC content of the primer

$L =$ length of primer (bp)

After calculating the T_m s of a primer pair they differ by too much it is best to add further bases to the less stable of the pair from the 5' end.

- b) Choosing an annealing temperature for a primer set** – the first time a pair of primers is used it is best to attempt an annealing temperature 5^o C below the average T_m for the primer pair. Annealing temperature is then adjusted depending on the preliminary results or alternatively a Mg²⁺ titration can be useful especially where the product size is large or the GC content is high.

1.15.5 Secondary structure formation:

An important factor to consider when designing a primer is the presence of secondary structures. Presence of secondary structures greatly reduces the number of primer molecules available for the reaction. The presence of hairpin loops reduces the efficiency by limiting the ability to bind to the target site. The single stranded nucleic acid sequences may have secondary structure due to the presence of complementary sequences in the primer sequence. Secondary structure formation occurs mainly because of repetitive sequence.

1.15.6 Specificity:

Primer specificity is partly dependent on primer length. Primers must be chosen so that they have a unique sequence with in the template DNA that is to be amplified. A primer designed with a high repetitive sequence will result a smear when amplifying genomic DNA. How ever the same primer may give a single band if clone from a genomic library is amplified. Because Taq DNA Polymerase is active over a broad range of temperature, primer extension will occur at the lower temperatures of annealing. If the temperature is too low, non-specific priming may occur which can be extended by the Polymerase if there is a short homology at the 3' end. In general, a melting temperature of 55^o C to 72^o C gives the best results.

1.15.7 Primer ends:

The 3' terminal position in PCR primers is essential for the control of mis-priming. Another variable to look at is the inclusion of a G or C residue at the 3' end of primer. This

will help to improve the efficiency of the reaction by minimizing local unwind of DNA double helix. This occurs because of AT rich sequence.

1.16 PRIMER3

To design primers for a region of interest, Genotator i.e. *Primer3* is used. The development of *Primer3* and the *Primer3* WWW interface were funded by Howard Hughes Medical Institute and by the National Institutes of Health, National Human Genome Research Institute, under grants ROI – HG00257 and P50-HG0098.

Primer 3 started as a reimplement of *Primer .5* as software component; the design of *Primer 3* draws heavily on the design of *Primer .5* and *Primer v2* and WWW interface designed by Richard Resnick for *Primer v2*.

Primer 3 is a computer program that suggests PCR primers for a variety of applications.

- a) To create STS (sequence tagged sites).
- b) To amplify sequences for single nucleotide polymorphism discovery.
- c) To select single primers for sequencing reactions.
- d) Do design oligo nucleotide hybridization probes.

1.17 Procedure For Primer Design By Using Primer3:

- Select the query sequence and Id, which has the Tandem repeats in it.
- Paste source sequence in **FASTA** format.
- Paste the sequence Id number in the sequence Id blank.
- Then put the tandem repeats in brackets [], which are present in the source sequence.
- Then adjust parameters according to our requirement.
- Then click the **Pick Primers** option.
- Then it shows the results as output.

1.18 Primer3 Input Parameters:

1.18.1 Source Sequence

The sequence from which to select primers or hybridization oligos.

1.18.2 Sequence Id

An identifier that is reproduced in the output to enable the user to identify the chosen primers.

1.18.3 Targets

If one or more Targets are specified then a legal primer pair must flank at least one of them. A Target might be a simple sequence repeat site (for example a CA repeat) or a single-base-pair polymorphism.

1.18.4 Excluded Regions

Primer oligos may not overlap any region specified in this tag. The associated value must be a space-separated list of

Start, length

Pairs where *start* is the index of the first base of the excluded region, and *length* is its length.

E.g. 401,7 68,3 forbids selection of primers in the 7 bases starting at 401 and the 3 bases at 68. e.g. ...ATCT<CCCC>TCAT. Forbids primers in the central CCCC.

1.18.5 Product Size

Minimum, Optimum, and Maximum lengths (in bases) of the PCR product. Primer3 will not generate primers with products shorter than Min or longer than Max, and with default arguments Primer3 will attempt to pick primers producing products close to the Optimum length,

1.18.6 Number To Return

The maximum number of primer pairs to return. Primer pairs returned are sorted by their "quality", in other words by the value of the objective function (where a lower number indicates a better primer pair). Setting this parameter to a large value will increase running time.

1.18.7 Max 3' Stability

The maximum stability for the five 3' bases of a left or right primer. Bigger numbers mean more stable 3' ends.

1.18.8 Max Mispriming

The maximum allowed weighted similarity with any sequence in Mispriming Library. Default is 12.

1.18.9 Pair Max Mispriming

The maximum allowed sum of similarities of a primer pair (one similarity for each primer) with any single sequence in Mispriming Library. Default is 24

1.18.10 Primer Size

Minimum, Optimum, and Maximum lengths (in bases) of a primer oligo. Primer3 will not pick primers shorter than Min or longer than Max, and with default arguments will attempt to pick primers close with size close to Opt. Min cannot be smaller than 1. Max cannot be larger than 36. (This limit is governed by maximum oligo size for which melting-temperature calculations are valid.) Min cannot be greater than Max.

1.18.11 Primer T_m

Minimum, Optimum, and Maximum melting temperatures (Celsius) for a primer oligo. Primer3 will not pick oligos with temperatures smaller than Min or larger than Max, and with default conditions will try to pick primers with melting temperatures close to Opt.

1.18.12 Maximum T_m Difference

Maximum acceptable (unsigned) difference between the melting temperatures of the left and right primers.

1.18.13 Product T_m

The minimum, optimum, and maximum melting temperature of the amplicon. Primer3 will not pick a product with melting temperature less than min or greater than max.

$$T_m = 81.5 + 16.6(\log_{10}([Na^+])) + .41*(\%GC) - 600/\text{length},$$

where $[Na^+]$ is the molar sodium concentration, (%GC) is the percent of Gs and Cs in the sequence, and length is the length of the sequence.

1.18.14 Primer GC%

Minimum, Optimum, and Maximum percentage of Gs and Cs in any primer.

1.18.15 Max Complementarity:

The maximum allowable local alignment score when testing a single primer for (local) self-complementarity and the maximum allowable local alignment score when testing for complementarity between left and right primers. For example, the alignment

```
5' ATCGNA 3'  
  |||  
3' TA-CGT 5'
```

is allowed (and yields a score of 1.75), but the alignment

5' ATCCGNA 3'

|| ||

3' TA--CGT 5'

is not considered. Scores are non-negative, and a score of 0.00 indicates that there is no reasonable local alignment between two oligos.

1.18.16 Max 3' Complementarity:

The maximum allowable 3'-anchored global alignment score when testing a single primer for self-complementarity, and the maximum allowable 3'-anchored global alignment score when testing for complementarity between left and right primers. The 3'-anchored global alignment score is taken to predict the likelihood of PCR-priming primer-dimers, for example

5' ATGCCCTAGCTTCCGGATG 3'

||| ||||

3' AAGTCCTACATTTAGCCTAGT 5'

or

5' AGGCTATGGGCCTCGCGA 3'

|||||

3' AGCGCTCCGGGTATCGGA 5'

The scoring system is as for the Max Complementarity argument. In the examples above the scores are 7.00 and 6.00 respectively. Scores are non-negative, and a score of 0.00 indicates that there is no reasonable 3'-anchored global alignment between two oligos. In order to estimate 3'-anchored global alignments for candidate primers and primer pairs, Primer assumes that the sequence from which to choose primers is presented 5'→3'. It is nonsensical to provide a larger value for this parameter than for the Maximum (local) Complementarity parameter because the score of a local alignment will always be at least as great as the score of a global alignment.

1.18.17 Max Poly-X:

The maximum allowable length of a mononucleotide repeat, for example AAAAAA.

1.18.18 Included Region:

A sub-region of the given sequence in which to pick primers. For example, often the first dozen or so bases of a sequence are vector, and should be excluded from consideration. The value for this parameter has the form

start,length

where *start* is the index of the first base to consider, and *length* is the number of subsequent bases in the primer-picking region.

1.18.19 Start Codon Position:

This parameter should be considered EXPERIMENTAL at this point. Some erroneous inputs might cause an error in Primer3. Index of the first base of a start codon. This parameter allows Primer3 to select primer pairs to create in-frame amplicons.

1.18.20 Mispriming Library:

This selection indicates what mispriming library (if any) Primer3 should use to screen for interspersed repeats or for other sequence to avoid as a location for primers.

1.18.21 CG Clamp:

Require the specified number of consecutive Gs and Cs at the 3' end of both the left and right primer. (This parameter has no effect on the hybridization oligo if one is requested.)

1.18.22 Salt Concentration:

The millimolar concentration of salt (usually KCl) in the PCR. Primer3 uses this argument to calculate oligo melting temperatures.

1.18.23 Annealing Oligo Concentration:

The nanomolar concentration of annealing oligos in the PCR. Primer3 uses this argument to calculate oligo melting temperatures.

1.18.24 Max Ns Accepted:

Maximum number of unknown bases (N) allowable in any primer.

1.18.25 Liberal Base:

This parameter provides a quick way to get Primer3 to accept IUB / IUPAC codes for ambiguous bases (i.e. by changing all unrecognized bases to N).

1.18.26 First Base Index:

The index of the first base in the input sequence. For input and output using 1-based indexing (such as that used in GenBank and to which many users are accustomed) set this parameter to 1. For input and output using 0-based indexing set this parameter to 0. (This parameter also affects the indexes in the contents of the files produced when the primer file flag is set.) In the WWW interface this parameter defaults to 1.

1.18.27 Inside Target Penalty:

Non-default values valid only for sequences with 0 or 1 target regions. If the primer is part of a pair that spans a target and overlaps the target, then multiply this value times the number of nucleotide positions by which the primer overlaps the (unique) target to get the 'position penalty'. The effect of this parameter is to allow Primer3 to include overlap with the target as a term in the objective function.

1.18.28 Outside Target Penalty:

Non-default values valid only for sequences with 0 or 1 target regions. If the primer is part of a pair that spans a target and does not overlap the target, then multiply this value times the number of nucleotide positions from the 3' end to the (unique) target to get the

'position penalty'. The effect of this parameter is to allow Primer3 to include nearness to the target as a term in the objective function.

1.18.29 Sequence Quality

A list of space separated integers. There must be exactly one integer for each base in the Source Sequence if this argument is non-empty. High numbers indicate high confidence in the base call at that position and low numbers indicate low confidence in the base call at that position.

1.18.30 Min Sequence Quality:

The minimum sequence quality (as specified by Sequence Quality) allowed within a primer.

1.18.31 Min 3' Sequence Quality:

The minimum sequence quality (as specified by Sequence Quality) allowed within the 3' pentamer of a primer.

1.18.32 Sequence Quality Range Min:

The minimum legal sequence quality (used for interpreting Min Sequence Quality and Min 3' Sequence Quality).

1.18.33 Sequence Quality Range Max:

The maximum legal sequence quality (used for interpreting Min Sequence Quality and Min 3' Sequence Quality).

1.18.34 Penalty Weights:

This section describes "penalty weights", which allow the user to modify the criteria that Primer3 uses to select the "best" primers. There are two classes of weights: for some parameters there is a 'Lt' (less than) and a 'Gt' (greater than) weight. These are the weights

that Primer3 uses when the value is less or greater than (respectively) the specified optimum. The following parameters have both 'Ll' and 'Gt' weights:

- Product Size
- Primer Size
- Primer T_m
- Product T_m
- Primer GC%
- Hyb Oligo Size
- Hyb Oligo T_m
- Hyb Oligo GC%

For the remaining parameters the optimum is understood and the actual value can only vary in one direction from the optimum:

- Primer Self Complementarity
- Primer 3' Self Complementarity
- Primer #N's
- Primer Mispriming Similarity
- Primer Sequence Quality
- Primer 3' Sequence Quality
- Primer 3' Stability
- Hyb Oligo Self Complementarity
- Hyb Oligo 3' Self Complementarity
- Hyb Oligo Mispriming Similarity
- Hyb Oligo Sequence Quality
- Hyb Oligo 3' Sequence Quality

The following are weights are treated specially:

- Position Penalty Weight

- Determines the overall weight of the position penalty in calculating the penalty for a primer.
- Primer Weight
- Determines the weight of the 2 primer penalties in calculating the primer pair penalty.
- Hyb Oligo Weight

Determines the weight of the hyb oligo penalty in calculating the penalty of a primer pair plus hyb oligo.

The following govern the weight given to various parameters of primer pairs (or primer pairs plus hyb oligo).

- T_m difference
- Primer-Primer Complementarity
- Primer-Primer 3' Complementarity
- Primer Pair Mispriming Similarity

1.18.35 Hyb Oligos (Internal Oligos):

Parameters governing choice of internal oligos are analogous to the parameters governing choice of primer pairs.

Primer3

(New interface) pick primers from a DNA sequence

Clipboard Paste
 Paste source sequence below (5' to 3', string of ACGTNaactm -- other letters treated as N -- numbers and blanks ignored). FASTA format ok. Please N-out undesirable sequence (vector, ALUS, LINES, etc.) or use a MASK option

LIBRARY SEQUENCE LIBRARY [NONE]

<input type="checkbox"/> Pick left primer or use left primer below.	<input type="checkbox"/> Pick hybridization probe (internal oligo) or use oligo below.	<input type="checkbox"/> Pick right primer or use right primer below (5' to 3' on opposite strand).
---	--	---

SEQUENCE _____ A string to identify your output.
 E.g. 50.2 requires primer to surround the 2 bases at positions 50 and 51. Or mark the sequence with | and |. e.g. ...ATCC|TGGCC|CAT... means that primers must flank the central CCCC. E.g. 401.7 60.3 forbids selection of primers in the 7 bases starting at 401 and the 3 bases at 60. Or mark the sequence with < and >. e.g. ...ATCC<CCCC>CAT... forbids primers in the central CCCC.

MINIMUM SIZE Min: Opt: Max:
MINIMUM LENGTH **MAXIMUM LENGTH**
MAX. MAXIMUM **MAX. MIN. DIFFERENCE**

General Primer Picking Conditions

PRIMER SIZE Min: Opt: Max: **MIN. MIN. DIFFERENCE**
MINIMUM IN Min: Opt: Max:
MINIMUM OUT Min: Opt: Max:
MIN. SIZE Min: Opt: Max:
MIN. SIZE **MIN. SIZE**
MIN. SIZE **MIN. SIZE**
MIN. SIZE **MIN. SIZE**
MIN. SIZE **MIN. SIZE**
MIN. SIZE **MIN. SIZE**
 General Best Show Primer Picking

Other Per-Sequence Inputs

INCLUDED REGION _____ E.g. 20,400: only pick primers in the 400 base region starting at position 20. Or use (end) in the sequence string to mark the beginning and end of the included region: e.g. in ATC (TTC...TCT)AT the included region is TTC...TCT.
Start Codon Position _____
Sequence Quality _____

Sequence Quality

Min. Sequence Reads: Min. Read Coverage: Sequence Quality Filter: Sequence Quality Filter:

Objective Function Penalty Weights for Primers

Len L: G:
 GC% L: G:
 GC% L: G:
 Self-Complementarity:
 3' Self-Complementarity:
 dN%:
 Mismatches:
 Sequence Quality:
 Read Sequence Quality:
 Position Penalty:
 End Penalty:

Objective Function Penalty Weights for Primer Pairs

Product Size L: G:
 Product Len L: G:
 Len Difference:
 Avg. Self-Complementarity:
 3' Self-Complementarity:
 Over-Mismatches:
 Primer-Primer Mismatch:
 Hyb. Class Penalty Weight:

Hyb. Oligo (Internal Oligo) Per-Sequence Inputs

Hyb. Class:

Hyb. Oligo (Internal Oligo) General Conditions

Hyb. Oligo Size: Min Opt Max
 Hyb. Oligo Len: Min Opt Max
 Hyb. Oligo GC%: Min Opt Max
 Hyb. Oligo Self-Complementarity: Hyb. Oligo Max. 3' Self-Complementarity:
 Max. dN%: Hyb. Oligo Max. Poly-X:
 Hyb. Oligo Mismatch Label: Hyb. Oligo Max. Mismatch:
 Hyb. Oligo Min. Sequence Quality:
 Hyb. Oligo Self-Contribution: Hyb. Oligo DNA Concentration:

Objective Function Penalty Weights for Hyb. Oligos (Internal Oligos)

Hyb. Oligo Len L: G:
 Hyb. Oligo Len L: G:
 Hyb. Oligo GC% L: G:
 Hyb. Oligo Self-Complementarity:
 Hyb. Oligo dN%:
 Hyb. Oligo Mismatches:
 Hyb. Oligo Sequence Quality:

1.19 Methods For Finding Tandem Repeats

All available nucleotide sequences of sorghum from the NCBI (<http://www.ncbi.nlm.nih.gov/entrez/>) were downloaded in FASTA format on to a local database.

1.19.1 Steps Involved After Downloading The Sequence Of Sorghum From NCBI.

1.19.1.1 Fasta Format

The FASTA format looks like this,

```
gi|18063613|gb|BM324719.1|BM324719 PIC1_34_G05.b1_A002 Pathogen-infected
compatible 1 (PIC1) Sorghum bicolor cDNA, mRNA sequence
GCACGAGGCACACACACACTGACGTACGATA'TCGAACACACTGACACAA
CACAACCGCCGGGCGCCGGCTTTG'TCGT'CAACGTCACCATTCTTAACAC
ACTCTCTACATACAGTTCGACGCCGGCCGAGCTCATCTTCTGTATTTAAAC
GCGTCCATCCATTCCGGTCAGGGCACGGAGTCGTC'ACTCCTAGCTCCGCGC
GCACCGAGCTCCACACACTCTCGCAGCAATGGCACGCGGGCGGCTCTACCGC
GGCTCCTCCACCCCT'CCUGCCGCTGCTGCT'GCTGCTGG'FGGTAGTAGTGGTGT
CCTCGGCCTTTTTCGCTCGCTGCTGCCACGCCGCGCCGCCGGGTGCACTGGTG
ACGAGCGTCCCGGGGTACACGGCGGGCAGCGGGCGCGCTGCCGTGAAGCAC
TACGCCGGGTACGTGACGGTGGACGAGGGCGCACGGCAGGAGGCTCTTCTACT
ACCTGGTGGAGCCGAGCGTGACCCCGCCAAGGACCCCGTCGTGCTGTGGCTC
AACGGCGGGCCTGGCTGCTCCAGCTTCGACGGC
```

1.19.1.2 Enter The Sequence For Finding Tandem Repeats.

Each of the sequence was searched for tandem repeats region and motif using the Tandem Repeat finder at <http://c3.biomath.mssm.edu/trf/advanced.submit.html>. This program was developed by the department of biomathematical sciences, Mount Sinai School of Medicine (<http://www.mssm.edu/school.html>). Advanced options were used to be find the tandem repeat region.

Parameters are set according the requirement. Alignment parameters (match, mismatch, indels) are set to (2,3,5). Minimum alignment score to report repeat is set to 30 and option period size to 10.

TANDEM REPEATS FINDER

Parameters:

Alignment Parameters (match, mismatch, indels)

Minimum Alignment Score To Report Repeat

Maximum Period Size

Options:

Flanking Sequence Masked Sequence File Data File

Sequence:

Your data must be a DNA sequence in FASTA format.

Choose one of the following ways to send your data:

- Upload a file from your directory.

- Cut and paste sequence.

```

>111
GGATTGCGAGTTCCGCTATCCGAGACCGCCAGCCGCGGTCGCGCCAGCTTTTCT
TAGGGGCGAGGAGAGGCGGCGGCGTATAGCGTATGCGCTTGGCAGCTTCGGAT
CGCGGCGCTATGAGCGCGCGGCGGCGGCGGCGGCGGCGGCGGCGGCGGCGG
CGCATATTAATTAATTAATTAATTAATTAATTAATTAATTAATTAATTAATTA
TTTGAATTAATTAATTAATTAATTAATTAATTAATTAATTAATTAATTAATTA
    
```

Submit sequence

Copy the alignment

Last revised July 11, 1998

Send any questions or comments to submit@ncbi.nlm.nih.gov

SUBMISSION OF SEQUENCE TO ADVANCE TANDEM REPEAT FINDER

1.19.1.3 Submit Sequence

By clicking the submit sequence button which is on the advanced tandem repeat finder we will get program analysis of DNA sequence on screen

file Edit View Favorites Open Help

address:

TANDEM REPEATS FINDER

The Tandem Repeats Finder Server.

Please cite:
G. Benson,
"Tandem repeats finder: a program to analyze DNA sequences,"
Nucleic Acids Research (1999)
Vol. 27, No. 2, pp. 573-580.

[Click here for...](#)

If you save the summary table and alignment files, use the default name supplied by the browser to preserve the automatic cross-linking

Copyright 1998

1.19.1.4 Tandem Repeats Report

Click on the tandem repeats reports button on the program to analyze DNA sequence. We will get the summary table.

Tandem Repeats Finder Program version 2.32

Gary Benson
Department of Biomathematical Sciences
Mount Sinai School of Medicine
Version 2.32

Please cite:
G. Benson,
"Tandem repeats finder: a program to analyze DNA sequences"
Nucleic Acid Research 1998;
Vol. 26, No. 4, pp. 573-580.

Sequence: 112
Parameter: 2 3 5 80 10 10 10
Length: 1124

Tables: 1

This is table 1 of 1

Click on indices to view alignment

TABLE SUMMARY

Indices	Period Size	Copy Number	Consensus Size	Percent Matches	Percent Gaps	Score	A	C	G	T	Entropy (bits)
1-112	2	13.7	2	97	0	70	56	24	3	6	1.0

The End!

SUMMARY TABLE

1.19.1.5 Summary Table

By clicking the column of indices of table explanation, which is present in summary table, we will get the alignment explanation.



Address: [redacted]

Tandem Repeat Finder Program written by:

Gary Benson
 Department of Biomedical Sciences
 Mount Sinai School of Medicine

Version 1.02

Sequence: 123

Parameters: 2 2 5 00 10 10 13

Repeat: CAA, Period: 0.10
 Tuple sizes: 0,1,2,3
 Tuple Distances: 0, 25, 150, MAXIMUM

Length: 1224
 Alignment: A10.11, C10.12, G10.13, T10.14

Found 01 11114 original direct final sites:

11114 CAACTCAGCCG

Indices: 1110--1159 Score: 61
 Period size: 9 CopyNumber: 11.7 Consensus size: 9

1110 CAACTCAGCCG

 * * * * *
 1119 CAG CAG CAG CAG CAA CAA CAA CAA CAA CAA CAA CAA CAA
 1 CAA CAA CAA CAA CAA CAA CAA CAA CAA CAA CAA CAA CAA

1159 CAACTCAGCCG

Statistics
 Matches: 37, Mismatches: 1, Indels: 0
 0.87 0.03 0.10

Matches are distributed among these distances:
 1 37 1.00

Alignment: A10.12, C10.14, G10.15, T10.16

Consensus pattern (5 bp):
 CAA

Done.

ALIGNMENT EXPLANATION TABLE

1.19.1.6 Alignment Explanation Table

Alignment explanation table explains the number of times the repeat is repeated. for 1x -
 In the above case the repeat CAA is repeated 9 times

1.19.1.7 Primer3

PRIMER3 program to used design the primers in sequences where repeat region were found. The sequence and its Id. Were pasted in the sequence text box in FASTA format. Tandem repeats regions were masked by placing them in brackets []. The parameters were set according to the requirement. The selected options for product size Minimum 100, Optimum 200, maximum 400. the primer annealing temperature Minimum 59.0, optimum 60.0, Maximum 61.0 degree celcius. Primer GC% should be Minimum 20.0 and Maximum 80.0, primer size should be minimum 18bp, Optimum 20bp maximum 27bp. Then click the **pick primer** option.

Sequence Quality

Max Sequences
 Max. Seq. Sequences Quality
 Sequences Coverage Range
 Sequences Quality Range

Objective Function Penalty Weights for Primers

Tm Lt: Gt:
 Self Lt: Gt:
 GC Lt: Gt:
 Self Complementarity
 N- Self Complementarity
 Hair
 Maximize
 Sequence Quality
 Seq. Coverage Quality
 Sequence Length
 Seq. Stability

Objective Function Penalty Weights for Primer Pairs

Forward-Reverse Lt: Gt:
 Right-Left Lt: Gt:
 Tm Difference
 Self Complementarity
 N- Complementarity
 Self Complementarity
 Hair
 Maximize
 Sequence Quality
 Seq. Coverage Quality

Hyb Oligo (Internal Oligo) Per-Sequence Inputs

Hyb Oligo Excluded Sequences

Hyb Oligo (Internal Oligo) General Conditions

Hyb Oligo Size: Min Opt Max
 Hyb Oligo Tm: Min Opt Max
 Hyb Oligo GC: Min: Opt: Max:
 Hyb Oligo Self Complementarity: Hyb Oligo Max. Self Comp. Complementarity:
 Max. Hair: Hyb Oligo Max. Hair:
 Hyb Oligo Maximize Library: Hyb Oligo Max. Maximize:
 Hyb Oligo Max. Sequence Quality:
 Hyb Oligo Seq. Coverage Quality: Hyb Oligo Seq. Concentration:

Objective Function Penalty Weights for Hyb Oligos (Internal Oligos)

Hyb Oligo Tm Lt: Gt:
 Hyb Oligo Size Lt: Gt:
 Hyb Oligo GC Lt: Gt:
 Hyb Oligo Self Complementarity
 Hyb Oligo Hair
 Hyb Oligo Maximize
 Hyb Oligo Sequence Quality

1.19.1.7 Out Put Of Primer3

Primer3 Output

No mispriming library specified
Using 1-based sequence positions
OLIGO
LEFT PRIMER 1088 20 59.96 55.00 4.00 3.00 GGCCCATCTAACCGTACAGA
RIGHT PRIMER 1285 20 59.50 45.00 1.00 1.00 TGCAACTGTTCTGTGTGCTG
SEQUENCE SIZE: 2124
INCLUDED REGION SIZE: 2124

PRODUCT SIZE: 198, PAIR ANY COMPL: 4.00, PAIR 3' COMPL: 2.00
TARGETS (start, len): 1131,27

1 GAAATTCGAGCTCGGTACCCAGACCTCCCAAACCCATGCTCGCCACGTTTGTAGGCCAAGG

61 AGGCGCCGGCCACACCTACCAGGCCCTACCAATCCGCCATGTCTAGCTATCAAGCCCTCC

121 TTGCCCTCAGTGAAATGTCAAGATTGTTTCACCATATTTATTTATGACAGAGGACTTGATAA

181 TTTTTTTCTTGTAATCAAAGTTAATTAAGTTTGTCAAATTTACACAATACTAAAGCA

241 ATTGTAATCAAGAAACACAGGGAGTGCCTGTGATAATAGCTATAAAGCATTCTAATTTGTAC

301 ATTCTATTTGTGTGCATACATCTGTACATACTGGGATTTCAATTTGTCTTGATCTTGTAG

361 CATTTTTCATCATTTGATGAACAACCTTCATCTAACTACGTTGCAAGACAAAATAGTACAGT

421 AGTACAACAAAGTCCCTTTGATAAAGGCTTTGATATACATGAGCAAGTCATAAATTTACTT

481 GCACATCATGTCTGTAAAGTGAACATTTCTGATGTGGCTAAGGCTATAACATGTGTAAG

541 GGTGAAGTGATGCTCACTCCTCATTTATTCGAAAAGTTCCAATAGAAAATGACAACCTTTTCC

601 TTGTAAGTAGTGGAAATTTGCTTTCTTCTACACAGACCATATAATCCAATAAAAATTGATAAC

661 TAAATGTCAAAAATCGACTAGGTGCCATGTCATCTATAGTTTATCTGTTGTTCGCAAAAAGC

721 CAAAATCTAAACAGATATCTATGAGCTCTCACTCATATAAATAGGCCCCAAATCAGTAGT

781 TAAACCATCGCCATAACATTGAGAGGAATTAGAAAAATACCAAGTGAACGAACTAGCA

2041 CAAAATAACGGCAAAGCAAACATAGTCAAATCAGATGCATGTATGGGGATCCTCTAGA

2101 GTCGACCTGCAGGCATGCAAGCTT

KEYS (in order of precedence):

***** target.

>>>>> left primer

<<<<< right primer

ADDITIONAL OLIGOS

1	LEFT PRIMER	1081	20	60.15	50.00	5.00	0.00	AATCCTTGCCCATCTAACC
	RIGHT PRIMER	1285	20	59.50	45.00	7.00	1.00	TGCAACTGTTGTTGCTG
	PRODUCT SIZE:	205,	PAIR ANY	COMPL: 3.00,	PAIR 3'	COMPL: 2.00		
2	LEFT PRIMER	1080	20	59.39	50.00	5.00	0.00	CAATCCTTGCCCATCTAACC
	RIGHT PRIMER	1285	20	59.50	45.00	7.00	1.00	TGCAACTGTTGTTGCTG
	PRODUCT SIZE:	206,	PAIR ANY	COMPL: 3.00,	PAIR 3'	COMPL: 3.00		
3	LEFT PRIMER	1001	20	60.00	60.00	6.00	2.00	CCTAGCTGTGCAGGCCTATC
	RIGHT PRIMER	1188	20	59.06	55.00	8.00	3.00	CGGCTAGTGCACCTGAGTGAT
	PRODUCT SIZE:	188,	PAIR ANY	COMPL: 5.00,	PAIR 3'	COMPL: 3.00		
4	LEFT PRIMER	1002	20	60.00	60.00	6.00	2.00	CTAGCTGTGCAGGCCTATCC
	RIGHT PRIMER	1188	20	59.06	55.00	8.00	3.00	CGGCTAGTGCACCTGAGTGAT
	PRODUCT SIZE:	187,	PAIR ANY	COMPL: 5.00,	PAIR 3'	COMPL: 3.00		

Statistics

	con sid ered	too many Ns	in tar get	in excl reg	bad GC	no GC clamp	tm too low	tm too high	high any compl	high 3' compl	poly X	high end stab
ok												
Left 771	10645	0	0	0	448	0	5809	3518	5	22	14	58
Right 577	8773	0	0	0	142	0	3498	4419	11	24	24	78

Pair Stats:

considered 11174, unacceptable product size 10518, high end compl 117, ok 539
primer3 release 0.9

(primer3_www results.cgi v 0.2)

1.20 Results

All available (approximately 110055) nucleotide sequences were analyzed for the presence of tandem repeats up to 80bp (maximum) length of repeat motif and no penalty gaps or indels were allowed. Although such methods may not detect some repeats containing insertions or deletions, it was considered that such a conservative approach might not be affecting the results very dramatically. The summary of the total number of potential SSRs in each of the nucleotide length group is presented in table 1. As described in the above procedure the nucleotide sequences are allowed to find tandem repeats by submitting them to Tandem repeat finder. Possible primers were found to these sequences which had repeat regions using Primer3. All the sequences with tandem repeats do not yield primers either due the sequence being short and less than minimum product size that is prescribed in the parameter or the tandem repeat regions could be skewed to the sides resulting in the loss as a flanking region. The results of Primer3 output were collated and entered into a database table, with right primer sequence, left primer temperature, right primer temperature and total sequence as shown in Table 2.

Of the total 9232 a little over (6232) were from tri-nucleotide groups, the tetra-nucleotide groups that potentially form the polypeptide repeats were the next largest groups followed by di-nucleotide and penta-nucleotide. The SSRs summary by various nucleotide groups for a minimum repeat size length of 10bp and the units that are part of the group, are presented in Table3. The largest group is the tri-nucleotide group 'CCG' consisting of 3312 accessions. Of these, the largest sub-group were from the 'GCG' and 'GCC' group which both code for alanine. In the di-nucleotide repeats the most common repeat was the 'AG' group which translated to di-oligopeptide of Arginine and Glutamine.

Table 1 . Summary of the repeat motif by nucleotide unit length of more than 10bp repeat length

Nucleotide unit length	Number of SSRs
Di	962
Tri	6232
Tetra	1468
Penta	569
Total	9232

Table 2

Id	Accessio	Details	Sequence	Repeat	Leftprimer	Rightprimer	Lefttem	Rightte	Totalsequence
1	6674175	AW284211	AAGAGGAGCCAAGAA AAAGAAGAAGAAGGAG AAGAAGGCCAAAAAGG CCAAGAAGAAGAAGAA GTCCAAAAAGGATTGA GGTGGTCAGGCAAAG GCGGGGCAATAAAATC C	AAG(5)	GAGGAGCCA AAGAAAAAG AAGA	ATTTTATTGC CCCGCCTTT	59.13	60.63	112
2	7554456	AW680658	GCACGAGGATCGATCA GCGTGCATCCCTGACA AAAAAAGCACACACA CACATACATACACCTA CTAGTCTTCGTGGCTG CCAAGCTTAAGCTCGA GAAGGATTAACAATGG CGGTAGTCTCGCGG	CA(6)	GGATCGATC AGCGTGCAT	CCGCCATTG TTAATCCTTC T	60.78	59.04	126
3	8087874	AW922049	TCAATGGGTGGACCAA ACACAATTGCAAACCC TCAAGCTGCTGCTGCT GCTCAACTGACCCTCT TACAGCAGCAACGGAT CCTGCCTATGCAACAA CAGCATCAG	TGC(4)	ATGGGTGGA CCAAACACA AT	CTGATGCTGT TGTTGCATAG G	59.95	59.36	105
4	8089381	AW923556	CGCTTTAATGTCGCGA GATTGATAGCAAGGT CCTGCTGCTGCTGCTG CTGCTCTGGAACAGTG GGAACACTGCTGTCT TGCTCCGTTGTGTCTA GGCCGGGGTGTGCA TCGCTGGATT	CTG(6)	CGCGAGATT GATAGCAAG G	AATCCAGCG ATGCACACC	59.53	60.66	122

5	9847911	BE592838	GCACGAGGCTCGGCC AAAAATCTGCGGCGAG CTCCGCTCCGATCCGA TCCGATCCGACTGCTC CCCACGCCGTTGAGG CCGAGCGGGGGTCCA GGAGAGGAGAGGAGA GGCGCGGC	TCCGA(4)	AGGCTCGGC CAAAAATCT	GCCTCTCT CTCCTCTCCT G	59.77	60.63	116
6	11680518	BF588194	GCACGAGGGGATTGCT TTCCCCACCCACCTC CGTCTCCGTCTCTTGC CGACCGAAGCGACTCC ACCTGAGGAGAGGAG AGGAGAGGACAGAGG AGAGGACAGCAAGCTA CCGCCCCCGTCCCG	GAGGA(4)	CACGAGGGG ATTGCTTTC	CGGTAGCTT GCTGTCCTCT C	59.17	60.16	125
7	11920692	BF655560	GCACGAGGCTCTCTCT CTCTCTCTCTCTCTCT TCTTCTCTCGAATG GTTGTGTAGGCGGCG GGCCAGGAGGGAGG GAGGGAGGAGAGAGA GAGGGGAGCTAGGGT TTT	GGAG(4)	GCACGAGGC TCTCTCTCTC T	CCCTAGCTC CCCTCTCTCT C	59.03	59.53	112
8	12500452	BG049068	TACTACTAGTAGGACA GGAGTGCTCCATTATT AATTAATTATTCTATTA TAAAGTGCTTGCTTGC TTGCTGCTAATCCATG CACCATTAGGAAGAGA ACCCTCTCTCCGTCGT GCATGCAGCTTGATC	TGCT(4)	AGGACAGGA GTGCTCCAT TATT	CACGACGGA GAGAGGGTT C	59.13	60.81	128

9	13588582	BC559584	GAAC TTTGGTTTCGTTT CCCTGAATGATCAACC ATCATCATCATCATCAT CATGTGTCGAACTCTG CTGTCCCAAGGATTT CCGTCTCCTCCTTTCC TCTCCTGCTCTTCT	ATC(6)	CTTTGGTTTC GTTTCCTG A	AGCAGGAGA GGAAAGGAG GA	60.08	60.47	112
10	14570584	BI099002	GCACGAGGCTTAGC GCCCAGCACGCTCTCG CAGTAGCAGCGAGCG AGCGTCCCGAGGGAG TGAGAGGTCCGGAGC AGTCGCCCGCGTCCA CCTTCGTTGTCGGAGA TGAAGGAGAGAGC	AGCG(3)	CAGCACGCT CTCGCAGTA	TCTCTCCTTC ATCTCCGACA A	60.03	59.93	120

Table 3. Summary by nucleotide group and members of the group

Click on the group or the individual unit to retrieve data

Group	SSR Count	SSR Units in Group
AC	255	AC CA GI IG
AG	357	AG CI GA IC
AT	260	AT IA
CG	91	CG GC
AAC	39	AAC ACA CAA GTI IGI ITG
AAG	294	AAG AGA CII GAA ICI IIC
AAT	29	AAT ATA AII TAA IAI IIA
ACC	281	ACC CAC CCA GGI GTG IGG
ACG	1378	ACG AGC CAG CGA CGI CIG GAC GCA GCI GIC ICG IGC
ACT	278	ACT AGT ATC AIG CAT CIA GAT GIA IAC IAG ICA IGA
AGG	621	AGG CCI CIC GAG GGA ICC
CCG	3312	CCG CGC CGG GCC GCG GGC
AAAC	31	AAAC AACA ACAA CAAA GTII IGI IITG
AAAG	29	AAAG AAGA AGAA CIII GAAA ICTI IICI ITIC
AAAT	20	AAAT AATA ATAA AIII IAAA IATII ITAI ITIA
AAAC	27	AAAC ACCA CAAC CCAA GGTII GTTG IGGI ITGG
AAACG	46	AAACG AAGC ACGA AGCA CAAG CGTI CTIG GAAC GCAA GCTI GTTC IGGI IGGI IICG ITGG

<u>AACT</u>	25	<u>AACT</u> <u>AATC</u> <u>ACTA</u> <u>AGTT</u> <u>ATCA</u> <u>ATTG</u> <u>CAAT</u> <u>CTAA</u> <u>GATT</u> <u>GTTA</u> <u>TAAC</u> <u>TAGT</u> <u>TCAA</u> <u>TGAT</u> <u>TTAG</u> <u>TTGA</u>
<u>AAGG</u>	29	<u>AAGG</u> <u>AGGA</u> <u>CCTT</u> <u>CTTC</u> <u>GAAG</u> <u>GGAA</u> <u>TCCT</u> <u>TTCC</u>
<u>AAGT</u>	15	<u>AAGT</u> <u>AATG</u> <u>ACTT</u> <u>AGTA</u> <u>ATGA</u> <u>ATTC</u> <u>CAAT</u> <u>CTTA</u> <u>GAAT</u> <u>GTAA</u> <u>TAAG</u> <u>TACT</u> <u>TCAT</u> <u>TGAA</u> <u>TTAC</u> <u>TTCA</u>
<u>AATT</u>	7	<u>AATT</u> <u>ATTA</u> <u>TAAT</u> <u>TTAA</u>
<u>ACAG</u>	7	<u>ACAG</u> <u>AGAC</u> <u>CAGA</u> <u>CTGT</u> <u>GACA</u> <u>GTCT</u> <u>TCTG</u> <u>IGTC</u>
<u>ACAT</u>	77	<u>ACAT</u> <u>ATAC</u> <u>ATGT</u> <u>CATA</u> <u>GTAAT</u> <u>TACA</u> <u>TATG</u> <u>TGTA</u>
<u>ACCC</u>	12	<u>ACCC</u> <u>CACC</u> <u>CCAC</u> <u>CCCA</u> <u>GGGT</u> <u>GGTG</u> <u>GTGG</u> <u>TGGG</u>
<u>ACCT</u>	108	<u>ACCT</u> <u>AGGT</u> <u>ATCC</u> <u>ATGG</u> <u>CATC</u> <u>CCAT</u> <u>CCTA</u> <u>CTAC</u> <u>GATG</u> <u>GGAT</u> <u>GGTA</u> <u>GTAG</u> <u>TACC</u> <u>TAGG</u> <u>TCCA</u> <u>TGGA</u>
<u>ACGC</u>	97	<u>ACGC</u> <u>CACG</u> <u>CGCA</u> <u>CGTG</u> <u>GCAC</u> <u>GCGT</u> <u>GTGC</u> <u>TGCG</u> <u>GCCA</u> <u>GCTG</u> <u>GGTC</u> <u>TGGC</u> <u>ACCG</u> <u>AGCC</u> <u>CAGC</u> <u>CCAG</u> <u>CCGA</u> <u>CGGT</u> <u>GACC</u>
<u>ACGG</u>	84	<u>ACGG</u> <u>AGGC</u> <u>CAGG</u> <u>CCGT</u> <u>CCTG</u> <u>CGBA</u> <u>CGTC</u> <u>CTGC</u> <u>GACG</u> <u>GCAG</u> <u>GCCT</u> <u>GGAC</u> <u>GGCA</u> <u>GTCC</u> <u>TCCG</u> <u>TGCC</u>
<u>ACGT</u>	181	<u>ACGT</u> <u>ATGC</u> <u>CATG</u> <u>CGTA</u> <u>GCAAT</u> <u>GTAC</u> <u>TACG</u> <u>TGCA</u> <u>AGTC</u> <u>CAGT</u> <u>GTCA</u> <u>TCAG</u> <u>TGAC</u>
<u>ACTC</u>	71	<u>ACTC</u> <u>AGTG</u> <u>CACT</u> <u>CTCA</u> <u>GAGT</u> <u>GTGA</u> <u>TCAC</u> <u>TGAG</u>
<u>AGAT</u>	194	<u>AGAT</u> <u>ATAG</u> <u>ATCT</u> <u>CTAT</u> <u>GATA</u> <u>TAGA</u> <u>TATC</u> <u>TCTA</u>
<u>AGCG</u>	54	<u>AGCG</u> <u>CGAG</u> <u>CGCT</u> <u>CTCG</u> <u>GAGC</u> <u>GCGA</u> <u>GCTC</u> <u>TCCG</u>
<u>AGCT</u>	228	<u>AGCT</u> <u>ATCG</u> <u>CGAT</u> <u>CTAG</u> <u>GATC</u> <u>GCTA</u> <u>TAGC</u> <u>TCGA</u>
<u>AGGG</u>	52	<u>AGGG</u> <u>CCCT</u> <u>CCTC</u> <u>CTCC</u> <u>GAGG</u> <u>GGAG</u> <u>GGGA</u> <u>TCCC</u>
<u>CCCG</u>	74	<u>CCCG</u> <u>CCGC</u> <u>CGCC</u> <u>CGGG</u> <u>GCCC</u> <u>GCCG</u> <u>GGCG</u> <u>GGGC</u> <u>CCGG</u> <u>CGGC</u> <u>GCCG</u> <u>GGCC</u>
<u>AAAAG</u>	69	<u>AAAAG</u> <u>AAACA</u> <u>AAATG</u> <u>AAGAA</u> <u>ATACA</u> <u>AITAC</u> <u>CATAT</u> <u>GTTTT</u> <u>TAATG</u> <u>TCATT</u> <u>TGATA</u> <u>TGTTA</u> <u>TGTTT</u> <u>TTCIT</u> <u>TITCI</u> <u>TTTGT</u> <u>TTTTG</u> <u>TTTTC</u> <u>ICTIT</u>

<u>AAAAT</u>	12	<u>AAAAT</u>	<u>ATAAA</u>	<u>TAAAA</u>	<u>TTATT</u>														
<u>AAAGC</u>	155	<u>AAAGC</u>	<u>AAGCA</u>	<u>AAGCT</u>	<u>ACATC</u>	<u>ACTGT</u>	<u>AGAGA</u>	<u>AGCTA</u>	<u>AGGAA</u>	<u>AGTAC</u>	<u>AGTGT</u>	<u>ATGGA</u>	<u>ATTCCG</u>	<u>CAGAT</u>					
		<u>CATCT</u>	<u>CATTG</u>	<u>CCAAT</u>	<u>CTAGA</u>	<u>CTCAT</u>	<u>CTGAT</u>	<u>CITCT</u>	<u>GAATC</u>	<u>GAGAA</u>	<u>GAGTA</u>	<u>GATCA</u>	<u>GAICT</u>	<u>GATIG</u>					
		<u>GCAAA</u>	<u>GGT</u>	<u>GTGTA</u>	<u>GTGT</u>	<u>GTTG</u>	<u>GTTTC</u>	<u>TACTG</u>	<u>TAGCA</u>	<u>TAGGA</u>	<u>TAGTC</u>	<u>TCAGA</u>	<u>TCCAT</u>	<u>TCCTA</u>					
		<u>TCGAA</u>	<u>TCTTC</u>	<u>TCTCT</u>	<u>TGATC</u>	<u>TGCAA</u>	<u>TGCTT</u>	<u>TGTAC</u>	<u>TGTTG</u>	<u>TTCAG</u>	<u>TTCCT</u>	<u>TGCA</u>	<u>TGCT</u>	<u>TGTG</u>					
		<u>TTCG</u>	<u>TGATT</u>																
<u>AACGG</u>	227	<u>AACGG</u>	<u>ACAGG</u>	<u>ACCGA</u>	<u>AGAGC</u>	<u>AGAGG</u>	<u>AGCAG</u>	<u>AGCGA</u>	<u>AGCTG</u>	<u>AGGGA</u>	<u>AGGTG</u>	<u>AGTCG</u>	<u>AGTGG</u>	<u>ATCCC</u>					
		<u>ATCGG</u>	<u>ATGCG</u>	<u>CACCT</u>	<u>CACGC</u>	<u>CAGCT</u>	<u>CAGTG</u>	<u>CCACA</u>	<u>CCCTT</u>	<u>CCTAC</u>	<u>CCTGA</u>	<u>CCTTC</u>	<u>CGATC</u>	<u>CGATG</u>					
		<u>CGGTT</u>	<u>CGTCT</u>	<u>CGTGT</u>	<u>CTAGC</u>	<u>CTCCT</u>	<u>CTCGT</u>	<u>CTCTC</u>	<u>CTCTG</u>	<u>CTGAG</u>	<u>CTGCA</u>	<u>CTGCT</u>	<u>GAGAG</u>	<u>GAGCA</u>					
		<u>GAGCT</u>	<u>GAGGA</u>	<u>GAGGT</u>	<u>GATCG</u>	<u>GCACA</u>	<u>GCACT</u>	<u>GCGTT</u>	<u>GCTCA</u>	<u>GCTCT</u>	<u>GCTTG</u>	<u>GGAGA</u>	<u>GGCTA</u>	<u>GGGAA</u>					
		<u>GGGAT</u>	<u>GGTGA</u>	<u>GTGCA</u>	<u>GTGTG</u>	<u>TCCGA</u>	<u>TCCGT</u>	<u>TCCGA</u>	<u>TCCGT</u>	<u>TCCCG</u>	<u>TCCGC</u>	<u>TCCGC</u>	<u>TCCGC</u>	<u>TCCGC</u>					
		<u>TGGTG</u>	<u>TTC</u>																
<u>ACCGC</u>	91	<u>ACCGC</u>	<u>AGGCC</u>	<u>CACCC</u>	<u>CAGCG</u>	<u>CAGGC</u>	<u>CCCCT</u>	<u>CCCCT</u>	<u>CCCTG</u>	<u>CCGTC</u>	<u>CCTCC</u>	<u>CCTGC</u>	<u>CGCCT</u>	<u>CGCGA</u>					
		<u>CGGAG</u>	<u>CGGCA</u>	<u>CGGCT</u>	<u>CTCCC</u>	<u>CTCCG</u>	<u>CTCGC</u>	<u>CTGCC</u>	<u>CTGGC</u>	<u>GAGGG</u>	<u>GCCGT</u>	<u>GCGGA</u>	<u>GCTCC</u>	<u>GCTGC</u>					
		<u>GGGAC</u>	<u>GGGGA</u>	<u>GTGCG</u>	<u>TCCCC</u>	<u>TCCGC</u>	<u>TCCCG</u>	<u>TCCGC</u>	<u>TGCC</u>	<u>TGCC</u>	<u>TGCC</u>	<u>TGCC</u>	<u>TGCC</u>	<u>TGCC</u>					
<u>CCCCG</u>	15	<u>CCCCG</u>	<u>CCCCG</u>	<u>CCGCC</u>	<u>CGCCG</u>	<u>CGCGC</u>	<u>CGGCG</u>	<u>GCCCC</u>	<u>GCGGC</u>	<u>GCGCG</u>									

1.21 Discussion:

De novo generation of micro-satellite markers through laboratory-based screening of SSR- enriched genomic libraries is highly time consuming and expensive. An alternative is to screen the public database of related model species where abundant sequence data is already available. Beyond cost savings, this approach also offers the possibility of identifying rare micro-satellite motifs, which would be uneconomical to identify through laboratory protocols. The availability of massive amounts of nucleotide sequence data has led to the development of innovative ways to examine these data as reflected in their functions. Various types of DNA markers have been used in plant breeding and of these the most extensively used are the micro-satellite markers. The reasons for their extensive use are due to their mode of transmission, which is bi-parental-nuclear with few loci and many alleles per locus. Mode of gene action being co- dominance with exception of null allele at some loci, show large variation within populations and are generally found in non-coding regions, which may contribute to the genome stability. Genome sequence and protein sequence information is publicly available for large -scale analysis from Gene bank at (NCBI) and European Molecular biology Laboratory (EMBL). Today, the search for a gene of interest starts with sequence

Information, including expressed sequence tags (EST). Genome related public database are an invaluable part of the scientific community and most notably the model organism database, have two major consumers: the focused scientific community actively studying that system, and the large scientific community interested in relating this specialized information to and from other systems. The thrust of any high-throughput facility is the creation of large, well-organized, rigorous datasets. The model system database can be mined by other related crop specialist to design markers for marker-assisted selection and -aided introgression methods. Such an approach can save valuable resources both in terms of time and funds.

2. Phylogenetic studies in cereals

2.1 Introduction

Phylogeny is really about Evolution and is used to reconstruct evolutionary events at a molecular level through analysis of molecular sequences, namely proteins & nucleic acids. Phylogeny is a diagram (a phylogenetic tree or cladogram) that depicts the evolutionary relationships among organisms. Comparative morphological, anatomical, embryological, molecular, behavioral, physiological, chemical, geographical, and fossil data can all be used, together or separately, to construct the phylogeny. Phylogeny provides the historical perspective from which to interpret the evolution of characters, patterns and processes of diversification, rates of evolution, historical biogeography, and co-evolutionary phenomena, such as the relationships between hosts and parasites or plants and herbivores. Phylogeny is used to classify organisms on the basis of their inferred evolutionary relationships (the phylogenetic approach to classification). Phylogeny is a hypothesis based on the best interpretation of the data at hand and subject to further evaluation (and possibly change) as new data become available. Phylogenetic focuses on the construction of descent relationships of species or groups of species and on how to incorporate these relationships into classification systems. Probably the first steps on establishing relationships and classifications were by Carolus Linnaeus in the 18th century. He is the founder of taxonomy, the biological branch of study dedicated to naming and classifying life forms. Linnaeus first established the binomial naming system, which continues today. In fact, many of the names he developed are still being used as scientific names. He also established a filing hierarchy, but avoided making his classifications based on evolutionary relationships. He was a natural theologian and believed that species were permanent creations that existed according to God's plan.

2.2 Phylogenetic terms

Adaptation	Change in an organism resulting from natural selection; a structure which is the result of such selection.
Anagenesis	Evolutionary change along an unbranching lineage; change without speciation.
Ancestor	Any organism, population, or species from which some other organism, population, or species is descended by reproduction.
Basal group	The earliest diverging group within a clade; for instance, to hypothesize that sponges are basal animals is to suggest that the lineage(s) leading to sponges diverged from the lineage that gave rise to all other animals.
Character	Heritable trait possessed by an organism; characters are usually described in terms of their states, for example: "hair present" vs. "hair absent," where "hair" is the character, and "present" and "absent" are its states.
Clade	A monophyletic taxon; a group of organisms which includes the most recent common ancestor of all of its members and all of the descendants of that most recent common ancestor. From the Greek word "klados", meaning branch or twig.
Cladogenesis	The development of a new clade; the splitting of a single lineage into two distinct lineages; speciation.
Cladogram	A diagram, resulting from a cladistic analysis, which depicts a hypothetical branching sequence of lineages leading to the taxa under consideration. The points of branching within a cladogram are called nodes. All taxa occur at the endpoints of the cladogram.

Convergence	Similarities, which have arisen independently in two or more organisms that are not closely related. Contrast with homology.
Crown group	All the taxa descended from a major cladogenesis event, recognized by possessing the clade's synapomorphy. See: stem group.
Derived	Describes a character state that is present in one or more subclasses, but not all, of a clade under consideration. A derived character state is inferred to be a modified version of the primitive condition of that character, and to have arisen later in the evolution of the clade. For example, "presence of hair" is a primitive character state for all mammals, whereas the "hairlessness" of whales is a derived state for one subclade within the Mammalia.
Diversity	Term used to describe numbers of taxa, or variation in morphology.
Endosymbiosis	When one organism takes up permanent residence within another, such that the two become a single functional organism. Mitochondria and plastids are believed to have resulted from endosymbiosis.
Evolution	Darwin's definition: descent with modification. The term has been variously used and abused since Darwin to include everything from the origin of man to the origin of life.
Evolutionary tree	A diagram, which depicts the hypothetical phylogeny of the taxa under consideration. The points at which lineages split represent ancestor taxa to the descendant taxa appearing at the terminal points of the cladogram.
Extinction	When all the members of a clade or taxon die, the group is said to be extinct.

Gradualism	A model of evolution that assumes slow, steady rates of change. Charles Darwin's original concept of evolution by natural selection assumed gradualism. Contrast with punctuated equilibrium.
Homology	Two structures are considered homologous when they are inherited from a common ancestor, which possessed the structure. This may be difficult to determine when the structure has been modified through descent.
Hypothesis	A concept or idea that can be falsified by various scientific methods.
In-group	In a cladistic analysis, the set of taxa which are hypothesized to be more closely related to each other than any are to the outgroup.
Lineage	Any continuous line of descent; any series of organisms connected by reproduction by parent of offspring.
Monophyletic	Term applied to a group of organisms, which includes the most recent common ancestor of all of its members and all of the descendants of that most recent common ancestor. A monophyletic group is called a clade.
Outgroup	In a cladistic analysis, any taxon used to help resolve the polarity of characters, and which is hypothesized to be less closely related to each of the taxa under consideration than any are to each other.
Paraphyletic	Term applied to a group of organisms, which includes the most recent common ancestor of all of its members, but not all of the descendants of that most recent common ancestor.
Parsimony	Refers to a rule used to choose among possible cladogram, which states that the cladogram implying the least number of changes in character states is the best.

Phylogenetic	Field of biology that deals with the relationships between organisms. It includes the discovery of these relationships, and the study of the causes behind this pattern.
Phylogeny	The evolutionary relationships among organisms; the patterns of lineage branching produced by the true evolutionary history of the organisms being considered.
Plesiomorphy	A primitive character state for the taxa under consideration.
Polarity of characters	The states of characters used in a cladistic analysis, either original or derived. Original characters are those acquired by an ancestor deeper in the phylogeny than the most recent common ancestor of the taxa under consideration. Derived characters are those acquired by the most recent common ancestor of the taxa under consideration.
Polyphyletic	Term applied to a group of organisms, which does not include the most recent common ancestor of those organisms; the ancestor does not possess the character shared by members of the group.
Primitive	Describes a character state that is present in the common ancestor of a clade. A primitive character state is inferred to be the original condition of that character within the clade under consideration. For example, "presence of hair" is a primitive character state for all mammals, whereas the "hairlessness" of whales is a derived state for one subclade within the Mammalia.
Pseudoextinction	The apparent disappearance of a taxon. In cases of pseudoextinction, this disappearance is not due to the death of all members, but the evolution of novel features in one or more lineages, so that the new clades are not recognized as belonging to the paraphyletic ancestral

	group, whose members have ceased to exist. The Dinosauria, if defined so as to exclude the birds, is an example of a group that has undergone pseudoextinction.
Punctuated equilibrium	A model of evolution in which change occurs in relatively rapid bursts, followed by longer periods of stasis.
Radiation	Event of rapid cladogenesis, believed to occur under conditions where a new feature permits a lineage to move into a new niche or new habitat, and is then called an adaptive radiation.
Rank	In traditional taxonomy, taxa are ranked according to their level of inclusiveness. Thus a genus contains one or more species; a family includes one or more genera, and so on.
Relatedness	Two clades are more closely related when they share a more recent common ancestor between them than they do with any other clade.
Reticulation	Joining of separate lineages on a phylogenetic tree, generally through hybridization or through lateral gene transfer. Fairly common in certain land plant clades; reticulation is thought to be rare among metazoans.
Selection	Process which favors one feature of organisms in a population over another feature found in the population. This occurs through differential reproduction -- those with the favored feature produce more offspring than those with the other feature, such that they become a greater percentage of the population in the next generation.
Sister group	The two clades resulting from the splitting of a single lineage.
Stasis	A period of little or no discernible change in a lineage.

Stem group	All the taxa in a clade preceding a major cladogenesis event. They are often difficult to recognize because they may not possess synapomorphies found in the crown group.
Synapomorphy	A character which is derived, and because it is shared by the taxa under consideration, is used to infer common ancestry.
Systematics	Field of biology that deals with the diversity of life. Systematics is usually divided into the two areas of phylogenetics and taxonomy.
Taxon	Any named group of organisms, not necessarily a clade.
Taxonomy	The science of naming and classifying organisms.
Vicariance	Speciation, which occurs as a result of the separation and subsequent isolation of portions of an original population.

2.3 Phylogenetic classifications

Integrating the results of a phylogenetic analysis into a classification is an extremely contentious issue in systematics. Basically, taxa are recognized on the basis of monophyly.

1. Linnacian system
 - a) Named groups are monophyletic
 - b) Not all groups are necessarily named
 - c) Ranks are arbitrary

2. Phylogenetic or "rankless" systems
 - a) Groups given unranked names
 - b) Groups defined by ancestry - i.e. phylogenetic tree
 - c) Groups described/diagnosed by a character(s) on the branch of the monophyletic

2.4 Methods of phylogenetic analysis

There are three methods

- 1) Cladistic
- 2) Phenetic
- 3) Multiple alignment

2.4.1 Cladistic Method

Cladistics is a particular method of hypothesizing relationships among organisms. Like other methods, it has its own set of assumptions, procedures, and limitations. Cladistics is now accepted as the best method available for phylogenetic analysis, for it provides an explicit and testable hypothesis of organismal relationships. The basic idea behind cladistics is that members of a group share a common evolutionary history, and are "closely related", more so to members of the same group than to other organisms. These groups are recognized by sharing unique features, which were not present in distant ancestors. These *shared derived* characteristics are called synapomorphies. The Cladistic approach takes into account both ancestral relationships between organisms and current data stored in their sequences. Here, morphological similarities/differences are also considered as part of the evolutionary history of organisms being compared. The software employed in this approach includes parsimony and maximum likelihood methods.

2.4.2 phenetic Method

Phenetics, also known as numerical taxonomy, involves the use of various measures of overall similarity for the ranking of species. There is no restriction on the number or type of characters (data) that can be used, although all data must be first converted to a numerical value, without any character "weighting." Each organism is then compared with every other for all characters measured, and the number of similarities (or differences) is calculated. The organisms are then clustered in such a way that the most similar are grouped close together and the more different ones are linked more distantly. The taxonomic clusters, called phenograms, that result from such an analysis do not necessarily reflect genetic similarity or evolutionary relatedness. The lack of evolutionary

significance in Phenetics has meant that this system has had little impact on animal classification, and as a consequence, interest in and use of Phenetics has been declining in recent years.

2.4.3 Multiple Alignment Method

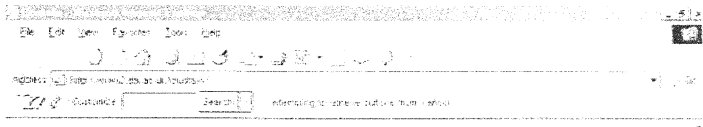
The most practical and widely used method in multiple sequence alignment is the hierarchical extensions of pair wise alignment methods.

- Multiple alignments are built by successive application of pair wise methods:
- Compare all sequences pair wise; (for N sequences there are $N(N-2)/2$ pairs or scores)
- Perform cluster analysis on pair wise scores to generate a hierarchy for alignments;
- Build the multiple alignment by first aligning the most similar pair of sequences, then the next most similar pair and so on;

Once an alignment of 2 sequences has been made, then this is fixed. Thus for a set of sequences A, B, C, D having aligned A with C and B with D the alignment of A, B, C, D is obtained by comparing the alignments of A and C with that of B and D using averaged scores at each aligned position.

2.5 Clustalw

ClustalW is a sequence alignment tool, which is used to produce Multiple alignments of a protein sequence/nucleotide sequence. This web tool available at <http://www.ebi.ac.uk/clustalw/>



2.5.1 upload a file

TREE TYPE:

TREE GRAPH INSTANCES:

```

1497 p001 agppatgtrtrh agppatgtrtrh
1498 p001 agppatgtrtrh agppatgtrtrh
1499 p001 agppatgtrtrh agppatgtrtrh
1741 agppatgtrtrh agppatgtrtrh
1861 agppatgtrtrh agppatgtrtrh
1862 agppatgtrtrh agppatgtrtrh
1863 agppatgtrtrh agppatgtrtrh
1864 agppatgtrtrh agppatgtrtrh
1865 agppatgtrtrh agppatgtrtrh
1866 agppatgtrtrh agppatgtrtrh
1867 agppatgtrtrh agppatgtrtrh
1868 agppatgtrtrh agppatgtrtrh
1869 agppatgtrtrh agppatgtrtrh
1870 agppatgtrtrh agppatgtrtrh
1871 agppatgtrtrh agppatgtrtrh
1872 agppatgtrtrh agppatgtrtrh
1873 agppatgtrtrh agppatgtrtrh
1874 agppatgtrtrh agppatgtrtrh
1875 agppatgtrtrh agppatgtrtrh
1876 agppatgtrtrh agppatgtrtrh
1877 agppatgtrtrh agppatgtrtrh
1878 agppatgtrtrh agppatgtrtrh
1879 agppatgtrtrh agppatgtrtrh
1880 agppatgtrtrh agppatgtrtrh
1881 agppatgtrtrh agppatgtrtrh
1882 agppatgtrtrh agppatgtrtrh
1883 agppatgtrtrh agppatgtrtrh
1884 agppatgtrtrh agppatgtrtrh
1885 agppatgtrtrh agppatgtrtrh
1886 agppatgtrtrh agppatgtrtrh
1887 agppatgtrtrh agppatgtrtrh
1888 agppatgtrtrh agppatgtrtrh
1889 agppatgtrtrh agppatgtrtrh
1890 agppatgtrtrh agppatgtrtrh
1891 agppatgtrtrh agppatgtrtrh
1892 agppatgtrtrh agppatgtrtrh
1893 agppatgtrtrh agppatgtrtrh
1894 agppatgtrtrh agppatgtrtrh
1895 agppatgtrtrh agppatgtrtrh
1896 agppatgtrtrh agppatgtrtrh
1897 agppatgtrtrh agppatgtrtrh
1898 agppatgtrtrh agppatgtrtrh
1899 agppatgtrtrh agppatgtrtrh
1900 agppatgtrtrh agppatgtrtrh
1901 agppatgtrtrh agppatgtrtrh
1902 agppatgtrtrh agppatgtrtrh
1903 agppatgtrtrh agppatgtrtrh
1904 agppatgtrtrh agppatgtrtrh
1905 agppatgtrtrh agppatgtrtrh
1906 agppatgtrtrh agppatgtrtrh
1907 agppatgtrtrh agppatgtrtrh
1908 agppatgtrtrh agppatgtrtrh
1909 agppatgtrtrh agppatgtrtrh
1910 agppatgtrtrh agppatgtrtrh
1911 agppatgtrtrh agppatgtrtrh
1912 agppatgtrtrh agppatgtrtrh
1913 agppatgtrtrh agppatgtrtrh
1914 agppatgtrtrh agppatgtrtrh
1915 agppatgtrtrh agppatgtrtrh
1916 agppatgtrtrh agppatgtrtrh
1917 agppatgtrtrh agppatgtrtrh
1918 agppatgtrtrh agppatgtrtrh
1919 agppatgtrtrh agppatgtrtrh
1920 agppatgtrtrh agppatgtrtrh
1921 agppatgtrtrh agppatgtrtrh
1922 agppatgtrtrh agppatgtrtrh
1923 agppatgtrtrh agppatgtrtrh
1924 agppatgtrtrh agppatgtrtrh
1925 agppatgtrtrh agppatgtrtrh
1926 agppatgtrtrh agppatgtrtrh
1927 agppatgtrtrh agppatgtrtrh
1928 agppatgtrtrh agppatgtrtrh
1929 agppatgtrtrh agppatgtrtrh
1930 agppatgtrtrh agppatgtrtrh
1931 agppatgtrtrh agppatgtrtrh
1932 agppatgtrtrh agppatgtrtrh
1933 agppatgtrtrh agppatgtrtrh
1934 agppatgtrtrh agppatgtrtrh
1935 agppatgtrtrh agppatgtrtrh
1936 agppatgtrtrh agppatgtrtrh
1937 agppatgtrtrh agppatgtrtrh
1938 agppatgtrtrh agppatgtrtrh
1939 agppatgtrtrh agppatgtrtrh
1940 agppatgtrtrh agppatgtrtrh
1941 agppatgtrtrh agppatgtrtrh
1942 agppatgtrtrh agppatgtrtrh
1943 agppatgtrtrh agppatgtrtrh
1944 agppatgtrtrh agppatgtrtrh
1945 agppatgtrtrh agppatgtrtrh
1946 agppatgtrtrh agppatgtrtrh
1947 agppatgtrtrh agppatgtrtrh
1948 agppatgtrtrh agppatgtrtrh
1949 agppatgtrtrh agppatgtrtrh
1950 agppatgtrtrh agppatgtrtrh
1951 agppatgtrtrh agppatgtrtrh
1952 agppatgtrtrh agppatgtrtrh
1953 agppatgtrtrh agppatgtrtrh
1954 agppatgtrtrh agppatgtrtrh
1955 agppatgtrtrh agppatgtrtrh
1956 agppatgtrtrh agppatgtrtrh
1957 agppatgtrtrh agppatgtrtrh
1958 agppatgtrtrh agppatgtrtrh
1959 agppatgtrtrh agppatgtrtrh
1960 agppatgtrtrh agppatgtrtrh
1961 agppatgtrtrh agppatgtrtrh
1962 agppatgtrtrh agppatgtrtrh
1963 agppatgtrtrh agppatgtrtrh
1964 agppatgtrtrh agppatgtrtrh
1965 agppatgtrtrh agppatgtrtrh
1966 agppatgtrtrh agppatgtrtrh
1967 agppatgtrtrh agppatgtrtrh
1968 agppatgtrtrh agppatgtrtrh
1969 agppatgtrtrh agppatgtrtrh
1970 agppatgtrtrh agppatgtrtrh
1971 agppatgtrtrh agppatgtrtrh
1972 agppatgtrtrh agppatgtrtrh
1973 agppatgtrtrh agppatgtrtrh
1974 agppatgtrtrh agppatgtrtrh
1975 agppatgtrtrh agppatgtrtrh
1976 agppatgtrtrh agppatgtrtrh
1977 agppatgtrtrh agppatgtrtrh
1978 agppatgtrtrh agppatgtrtrh
1979 agppatgtrtrh agppatgtrtrh
1980 agppatgtrtrh agppatgtrtrh
1981 agppatgtrtrh agppatgtrtrh
1982 agppatgtrtrh agppatgtrtrh
1983 agppatgtrtrh agppatgtrtrh
1984 agppatgtrtrh agppatgtrtrh
1985 agppatgtrtrh agppatgtrtrh
1986 agppatgtrtrh agppatgtrtrh
1987 agppatgtrtrh agppatgtrtrh
1988 agppatgtrtrh agppatgtrtrh
1989 agppatgtrtrh agppatgtrtrh
1990 agppatgtrtrh agppatgtrtrh
1991 agppatgtrtrh agppatgtrtrh
1992 agppatgtrtrh agppatgtrtrh
1993 agppatgtrtrh agppatgtrtrh
1994 agppatgtrtrh agppatgtrtrh
1995 agppatgtrtrh agppatgtrtrh
1996 agppatgtrtrh agppatgtrtrh
1997 agppatgtrtrh agppatgtrtrh
1998 agppatgtrtrh agppatgtrtrh
1999 agppatgtrtrh agppatgtrtrh
2000 agppatgtrtrh agppatgtrtrh

```

The nucleotide was set to file type: NUCLEOTIDE [2] 14 29

Comments or suggestions:

© 2000

2.5.1 upload a file

Upload a file from the computer which containing a valid protein nucleotide sequence in any format (GCG, FASTA, PIR, etc.) by using this option. This option only works with Netscape Browsers!

2.5.2 Sequences

The sequences have different names, as the first 30 characters of the name are significant. If clustalw finds two or more sequences with the same name it will fail!

The following are the points to follow not to get the fatal error.

A set of sequences in different formats has been used as input. For example:

>Seq1

```
MRVLLAALGL.LFLGALRAFQDRPFEDTCIIGNP SHYYDKAVRRCCYRCP  
MGLFPTQQCPQRPTDCR KQCEPDYYL.DEADRCTACVTC SRDD
```

>Seq2

```
PVAEERGLMSQPLMETCHSVGAAYLES.LPLQDASPAGGPSSPRDLPEPRV  
STEHTNNKIEKIYIMKADTVIVGTVKAELPEGRGLAGPAEP  
EL.FEELEADHTPHYPEQETEPPL.GSCSDVMI.SVEEEGKEDPI.PTAASGK
```

>Seq3

```
MSALLTAAGL   LFI.GMLQAFP   TDRPLKTTCA   GDLSHYPGEA  
ARNCCYQCPS GLSPTQPCPR  
GPAHCRKQCA   PDYYVNEDGK   CTACVTCLPG   LVEKAPCSGN  
SPRICECPG MHCCTPAVNS  
CARCKLHCSG   EEVVKSPGTA   KKDITICELPS   SGSGPNCNSP  
GDRKTLTSHIA TPQAMPTLES
```

All input sequence must be in the same format. The sequence separators every 10 aa in seq3 are not allowed in the fasta format specification. The program does not allow this kind of input.

- Sequences with spaces (deletions) will have the same affect as above.
- Sequences should not have illegal characters as for example.

>Badsequence

```
*SNERNVCN*WVYVWADQL*WIELKCFA*HICNSG*LQRYNPKNNQPI  
WGGIKRRCSS*
```

- Need at least two sequences to generate an alignment.

Illegal characters in the input should be removed. This could be due to cutting & pasting from a word processing program..

Remove any empty space or empty lines from the beginning of the input. put format for the alignment.

Example:

In FASTA format:

```
>gi|21326110|gb|AF503433.1| Sorghum bicolor BAC 170F8, complete sequence
AAGC'ITTTGCTAATCGGAAGCGAGACTTTGCT'CCCATCAGAGCAAAGT
ATGATC'TTGAGAAGTACAAGTGGGAAAGTCAAACC'GCAAGTGAAG
ATTGTCATCAAGAAAACCATCACAGAGGGCCTAAGAGGTGCCATCCAT
GTGAGACAGCAAAGAATACCTTGAGAAAGTGAAAGAA'TCAG'ITTACT'
GGTTCAACCAAAGCCCATGCTTCCACCCTGATCCAGAAACTCACTAAC
ATGAGGTTACAGGGGGGAGTGTGAGAGAGCACATTCT
```

2.5.3 search title

To identify the database search results.

2.5.4 cpu mode

This option allows the choice the version of clustalw based on cpu mode(Central processing unit). There are two options:

Clustalw_mp this is a parallelised version of clustalw that runs on multiprocessor SGI systems.

Clustalw this is the generic (single CPU) version of the program. Use this version if we are having problems with the MP version of the program.

2.5.5 alignment

choose a full alignment or using a stringent algorithm for generating the tree guide or a fast algorithm.

2.5.6 Output

Output format for multiple alignment results display. The options are ALN, GCG, PHYLIP, PIR and GDE.

2.5.7 Jalview

A new experimental option has been added to the results page, which involved using a Java Applet called JalView. This is a fully featured MSA editor which allows not only to edit the alignment but also, to exchange the alignment formats. For documentation click on the JalView Hyperlink.

2.5.8 Outorder

Decide, which order the sequences, should be printed in the alignment.

2.5.9 color

To get the alignment in color this option should be selected. This option only works when we choose ALN or GCG as the output format.

2.5.10 Fast pairwise alignment options

- | | |
|----------------|--|
| Ktup | This option allows you to choose which 'word-length' to use when calculating fast pairwise alignments.(note: make sure you have chosen 'fast' in the ALIGNMENT). |
| Window | Use this option to set the window length when calculating fast pairwise alignments.(Note: make sure you have chosen 'fast' in the ALIGNMENT). |
| Score | This option allows you to decide which score to take into account when calculating a fast pairwise alignment. (Note: make sure you have chosen 'fast' in the ALIGNMENT). |
| Topdiag | Select here how many top diagonals should be integrated when calculating a fast pairwise alignment.(Note: make sure you have chosen 'fast' in the ALIGNMENT). |
| Pairgap | Select here to set the gap penalty when generating fast pairwise alignments. |

2.5.11 multiple sequence alignment options

- | | |
|---------------|--|
| <i>matrix</i> | This option allows you to choose which matrix series to use when generating the multiple sequence alignment. The program goes through the chosen matrix series, spanning the full range of amino acid distances. |
|---------------|--|

<i>blosum (henikoff)</i>	These matrices appear to be the best available for carrying out data base similarity (homology searches). The matrices used are: Blosum80, 62, 40 and 30.
<i>pam (dayhoff)</i>	These have been extremely widely used since the late '70s. We use the PAM 120, 160, 250 and 350 matrices.
<i>gonnet</i>	These matrices were derived using almost the same procedure as the Dayhoff one (above) but are much more up to date and are based on a far larger data set. They appear to be more sensitive than the Dayhoff series. We use the GONNET 40, 80, 120, 160, 250 and 350 matrices. We also supply an identity matrix, which gives a score of 10 to two identical amino acids and a score of zero otherwise.

2.5.12 Gapopen

Set the penalty for opening a gap.

2.5.13 Endgap

Set the penalty for closing a gap.

2.5.14 Gapext

Set the penalty for extending a gap.

2.5.15 Gapdist

Set the gap separation penalty.

2.5.16 Phylogenetic Tree

To use this option need to input a sequence alignment, if this alignment is in PIR or PHYLIP format. ALN and GCG MSF files are not supported so this has to convert MSF files to PIR format.

This option allows choosing the following output formats for the tree:

Neighbor-Joining

Phylip

Distance

2.5.17 kimura correction of distances

This option allows setting on distances correction (correction for multiple substitutions).

This is because, as sequences diverge, more than one substitution will happen at many

Node: represents a taxonomic unit. This can be either an existing species or an ancestor.

Branch: Defines the relationship between the taxa in terms of descent and ancestry.

Topology: The branching patterns of the tree.

Branch length: Represents the number of changes that have occurred in the branch.

Root: The common ancestor of all taxa.

Distance scale: scale that represents the number of differences between organisms or sequences.

Clade: a group of two or more taxa or DNA sequences that includes both their common ancestor and all their descendants.

Operational Taxonomic Unit (OTU): Taxonomic level of sampling selected by the user to be used in a study, such as individuals, populations, species, genera, or bacterial strains.

A phylogenetic tree is composed of nodes--each representing a taxonomic unit (species, populations, and individuals)--and branches, which define the relationship between the taxonomic units in terms of descent and ancestry. Only one branch can connect any two adjacent nodes. The branching pattern of the tree is called the topology and the branch length usually represents the number of changes that have occurred in the branch. This is called a scaled branch. Scaled trees are often calibrated to represent the passage of time. Such trees have a theoretical basis in the particular gene or genes under analysis.

Branches can also be unscaled, which means that the branch length is not proportional to the number of changes that has occurred, although the actual number may be indicated numerically somewhere on the branch. Phylogenetic trees may also be either rooted or unrooted. In rooted trees, there is a particular node, called the root--representing a common ancestor--from which a unique path leads to any other node. An unrooted tree only specifies the relationship among species, without identifying a common ancestor, or evolutionary path.

2.7 Methods For The Phylogenetic Studies Of Cereals

Cereals i.e., sorghum bicolor, Hordeum vulgare, Zea mays, Triticum aestivum, orizya sativum are taken and co – related the phylogenetic relationship between these cereals by taking a nuclear enzyme, a mitochondrial enzyme and a chloroplast enzyme.

2.8 Reasons For Taking Enzymes In Phylogenetic Studies

2.8.1 Nuclear Enzyme

- 1) Nuclear enzymes have bi-parental inheritance.
- 2) Nuclear enzymes are in abundance in a Genome.
- 3) Amplification and sequencings are easy.
- 4) They have mosaics of highly conserved variable region. The conserved regions have been informative for resolving relationship at higher taxonomic levels. Alignment of the variable region is often problematic.
- 5) They exhibit a wide range of evolutionary rate in phylogenetic utility.
- 6) Many nuclear genes may contain large Introns that necessitate reverse transcriptase PCR.
- 7) The nuclear enzyme which we are taken are alcohol-dehydrogenase due to its role in metabolism a key functional pathway

2.8.2 Mitochondrial Enzyme

- 1) Mitochondrial enzymes are maternally inherited.
- 2) These enzymes are used to construct a PhylogeneticTree to display the Evolutionary relationships between Individual sequences.
- 3) The structure of this gene tree contains information which in conjunction with a calibrated mutation rate for the DNA sequence under study, can be used to estimate a time-scale for events in evolutionary prehistory.
- 4) Sites that have frequently undergone mutations are less conserved among species compared to those where the consensus is more the sequence is highly

conserved. Evolutionary changes are found the non-conserved regions of the sequence.

- 5) These will provide the phylogenetic evolution of a given mitochondrial gene.
- 6) The mitochondrial enzyme investigated in this study is Malate-dehydrogenase a key enzyme in the respiratory pathway.

2.8.3 Chloroplast Enzyme

- 1) Chloroplast gene Restriction-Site Variation has been shown to be well suited for studies of genetic relationships at or below the family level.
- 2) The chloroplast genome consists of a large and a small region of Single-copy DNA separated by a pair of identical but inverted repeat sequences.
- 3) Restriction-pattern differences between taxa may be interpreted as site changes caused by single base substitutions or single insertion deletion events.
- 4) By relating variation in chloroplast DNA restriction-fragment patterns to specific mutations, either base substitution or indels data sets suitable for phylogenetic reconstructions using Parsimony analysis may be produced.
- 5) Restriction-site variation is used to estimate total sequence divergence between taxa. Such distance measures may used to reconstruct phylogenies.
- 6) Chloroplast Enzymes found in plants only. These enzymes are mostly related to C3 and C4 pathways of photosynthesis.
- 7) These enzymes are maternally inherited.
- 8) The chloroplast enzyme, which we taken is Phosphoenol pyruvate the key enzyme of the C4 photosynthetic pathway.

2.9 Multiple Alignment Method

The most practical and widely used method in multiple sequence alignment is the hierarchical extensions of pair wise alignment methods.

2.10 Steps Involved In Multiple Alignment Method

2.10.1 Select The Most Conserved Enzymes:

Select the most conserved and functional mitochondrial, chloroplast and nuclear enzymes of cereals. In cereals the most conserved functional enzymes are as follows: -

Mitochondrial Enzymes:

1. NADH dehydrogenase 1.6.5.3
2. ATP synthase 3.6.3.14
3. Succinate dehydrogenase 1.3.99.1
4. Malate dehydrogenase 1.1.1.82
5. Citrate synthase 4.1.3.28

Chloroplast Enzymes:

1. Fructose 1,6 bis phosphatase 3.1.3.11
2. Phosphoenolpyruvate carboxylase 4.1.1.31
3. Glyceraldehyde- 3-phosphate dehydrogenase 1.2.1.9

Nuclear Enzyme:

1. Methyl transferase 2.1.1.37
2. Alcohol dehydrogenase 1.1.1.1
3. Cysteine synthase 4.2.99.8

From the all above-mentioned enzymes we selected the following enzymes, which are more convenient for our practical purpose. i.e

Malate dehydrogenase (mitochondrial enzyme)

Phosphoenolpyruvate carboxylase (chloroplast enzyme)

Methyl transferase (nuclear enzyme)

2.10.2 Search The Sequences From NCBI:

Search the sequences of the Malate dehydrogenase (mitochondrial enzyme), Phosphoenolpyruvate carboxylase (chloroplast enzyme), and Methyl transferase (nuclear enzyme)

From NCBI.(National Centre For Biotechnology Information)

Then note the accession number of the above enzymes of interest of cereals.

Chloroplast Enzyme (Phosphoenol Pyruvate Carboxylase)

Cereal	Accession number
Maize	E17154
Rice	AF271995
Sorghum	AF399915

Paste accession numbers of the cereals for a Chloroplast enzyme (Phosphoenol pyruvate carboxylase) in the NCBI. The sequences for the accession number of the cereals for ex:- (maize, oryza sativa and sorghum bicolor) will display.



National Center for Biotechnology Information

NIH Publication #97-1472 NLM Catalog #97-1472

Search by

What does NCBI do?

Established in 1988 as a national resource for the United States, the National Center for Biotechnology Information (NCBI) provides public domain access and online retrieval of computational biology, genetics, bioinformatics, and molecular biology databases. NCBI also disseminates biotechnology information and the Center's understanding of molecular biology to the scientific community at large and to the general public through our health and disease

[Home](#) [About](#) [Contact Us](#)

[What does NCBI do?](#)

[Database](#)

[Data Files](#)

[Genetics](#)

[Molecular Biology](#)

[Publications](#)

[Reference](#)

[Software](#)

[Training](#)

[Webinars](#)

[Workshops](#)

[Zotero](#)

[Database](#) [Data Files](#) [Genetics](#)

[Molecular Biology](#) [Publications](#) [Reference](#)

[Software](#) [Training](#) [Webinars](#)

[Workshops](#) [Zotero](#)

RefSeqs for viral genomes!



There are now over 1,000 RefSeqs for viral genomes. Visit RefSeq and analyze a viral genome of interest and compare to RefSeq.

[Viral genomes weekly](#) [More](#)

NCBI in the News

NCBI now offers a wide range of new resources through [PubMed Central](#). This new iteration of the Entrez database system expands the original relevance of NCBI's mission and commitment by allowing scientists to find related primary literature such as taxonomy, gene expression, and drug discovery in the [Scientific Method](#).

[PubMed Central](#) [More](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

[NCBI in the News](#)

File Edit View Favorites Tools Help

Address: http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=pubmed

Advanced Search

Showing 1 to 3 of 3 items

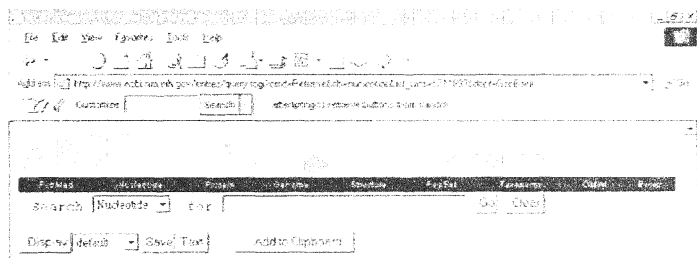
Pubmed ID	Article Title	Author	Journal	Pubmed ID	Article Title	Author	Journal
157773	Methylphenidate hydrochloride controlled-release tablets	g657112773&g657104	NEUROLOGY	157773	Methylphenidate hydrochloride controlled-release tablets	g657112773&g657104	NEUROLOGY
157774	Oral versus transdermal methylphenidate patch: a randomized, controlled trial	g65720444&g6571995	JAMA	157774	Oral versus transdermal methylphenidate patch: a randomized, controlled trial	g65720444&g6571995	JAMA
157775	Stability of the methylphenidate controlled-release tablets (MPH) over 12 months	g65720444&g6571995	JAMA	157775	Stability of the methylphenidate controlled-release tablets (MPH) over 12 months	g65720444&g6571995	JAMA

Revised: October 24, 2011

Page 1 of 1

2.10.3 Accession Number

First click on accession number of maize then oryza sativa followed by sorghum bicolor Sequences of Maize, oryza sativa and Sorghum bicolor phosphoenol pyruvate gene will obtain.



I: E17154. Maize phosphoenol [gi:5711837]
LOCUS E17154 6781 bp DNA linear PAT 28-JUL-1999
DEFINITION Maize phosphoenolpyruvate carboxylase gene.
ACCESSION E17154
VERSION E17154.1 GI:5711837
KEYWORDS JP 1998248419-A 2
SOURCE unidentified
ORGANISM *Zea mays*
unclassified.
REFERENCE 1 (bases 1 to 6781)
AUTHORS Matsuoka,M., Tokutomi,M., Toki,S. and Moris,S.K
TITLE C3 PLANT BODY EXPRESSING PHOTOSYNTHETIC ENZYME OF C4
PLANT
JOURNAL Patent: JP 1998248419-A 2 22-SEP-1998;
NORIN SUISANSYU NOGYO SHIBUITSU SHIGEN KENKYU SHU
COMMENT OS Zea mays (maize)
PN JP 1998248419-A 2
PD 22-SEP-1998
PF 11-MAR-1997 JP 1997056742
PI MATSUOKA MAKOTO, TOKUTOMI MITSUE, TOKI SEIKI, PI
MORIS
SUN-BEN KU
PC A01H5 00,A01H1 00,C07H21 04,C12N5 10,C12N9 00,C12N15 09,PC
iC12N5 10,
PC C12R1 91,(C12N15 09,C12R1 91);
CC strandedness: Double.
CC topology: Linear.

FH Key Location/Qualifiers
 FH
 FT source 1..6781
 FT /organism='Zea mays'
 FT /cultivar='Golden Cross Bantam' FT intron
 1477..1586
 FT exon 1587..1981
 FT intron 1982..2091
 FT exon 2092..2176
 FT intron 2177..3073
 FT exon 3074..3296
 FT intron 3297..3745
 FT exon 3746..3849
 FT intron 3850..3974
 FT exon 3975..4065
 FT intron 4066..4202
 FT exon 4203..4357
 FT intron 4358..4453
 FT exon 4454..5452
 FT intron 5453..5549
 FT exon 5550..5936
 FT intron 5937..6024
 FT CDS
 FT join(1297..1476,1587..1981,2092..2176, FT
 3074..3296,
 FT 3746..3849,3975..4065,4203..4357, FT
 4454..5452,5550..5936,
 FT 6025..6318)
 FT /product='phosphoenolpyruvate coaboxylase'.
 FEATURES Location/Qualifiers
 source 1..6781
 /organism="unidentified"
 /db_xref="taxon:32644"

BASE COUNT 1687 a 1854 c 1548 g 1692 t

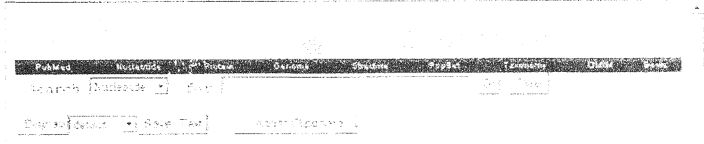
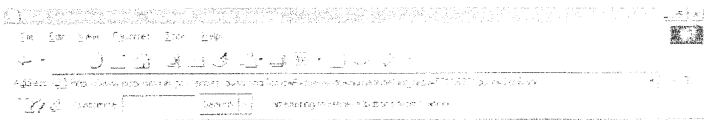
ORIGIN

1 tctagagatg taatggtgtt aggacacgtg gttagctact aatataaatg taaggcctaaa
 61 attcgatggt ttattttcta ttttcaacta cctagcatta tctcatttct aattgtgtga
 121 taacaaatgc attagaccat aattctgtaa atacgtacat ttaagcacac agtctatatt
 181 ttaaaattct tctttttgtg tggatatecc aacccaaate cacctctctc ctcaalccgt
 241 gtatcttcac cgtgccaag tccaacaac acatcgcac gtgcaaatct ttgttggtt
 301 gtgcacggte ggcgccaatg gaggagacac ctgtacgtg cctttggtag aacaacatcc
 361 ttatecctat atgtatggte cctctgtaga tggcaccctt atccctacaa tagccatgta
 421 ttgcatacca agaattaaat atacttttc ttgaaccaca ataatttat atageggcac
 481 ttctgttct ggttgaacac ttatttggaa caataaaate ccgagttcct aaccacaggt
 541 tcaacttttt tcttatcct cctaggaac taaattttaa atcataaat ttaattgaaa
 601 tgttaatgaa aacaaaaaaa ttactacaa agacgactct tagccacagc cgcctcactg

661 caccctcaac cacatctgc aaacagacac cctcggccaca tctctccaga ttcttccct
721 cgatgcagcc tacttgcata cagacgccct cteacatec tgaagaatc tctccaat
781 tcttgcate ccccgaatc agcattaact gctaaaggac gccctctcca catctteta
841 cccaattag caacggaata acacaagaag gcaggtgagc agtgacaaag caegtcaaca
901 gcaccgagcc aagccaaaag gagcaaggag gagcaagccc aagccgagc cgcagctcc
961 caggtccct tgcgattgcc gccagcagta gcagacaccc ctctccatc cctctcggc
1021 cgtcaacagc agcaagccaa gccaanaaga agcctcagcc acagccgggt cctgttgggt
1081 taccgccat cacatgccc aagcccgccc ttccaacag ccgagggccc cctgttccc
1141 tgcacagcca cacacacac ccgccgccc egaclccca tctctatgt aaceccccc
1201 cgcactgat tgatccaaa tgcctcgca gcagcacgag cagcacgccc tgccttcca
1261 accgtctgc ttccctgctt agcttcccgc cgcgccatgg cgtcgaccaa ggcacccggc
1321 cctggcgaga agcaccactc catcgacgcg cagctccgct agctgttccc aggcaaggct
1381 tccgaggac acaagctcat cagatcagat gcgctgctc tegaccgct cctcaacatc
1441 cctcaggacc tccagggccc cagccttgc gaattgtaa ctaaccacc ccgcccacca
1501 ttcttctt gaccggttgc cgcctggcgc cggcactgct cgtgtctgt gctagtctt
1561 agtgcttact actgtaatgc atgcaggtcc aggagtgcta ctaggtgtcgc gccgactatg
1621 agggcaaaag agacacgacg aagctggggc agctcggcgc caagctcagc gggcttggcc
1681 ccgccagcc catctctgtg gcgagctcca tcttgcacat gctcaacctc gccaaccttg
1741 ccgaggagggt gcagatcgcc caccgcccgc gcaacagcaa gctcaagaaa ggtgggttcg
1801 ccgacgaggg cctccgccc accgagtcg acatcgagga gcgctcaag cgcctcgtt
1861 ccgaggtcgg caagtcccc gagggggtgt tggagggcct caagaaccag accgtcgacc
1921 tcttctcac cgcgcatccc acgcagtcg cccgcccctc gctctcgag aaaaagcca
1981 ggtatatatt tctaatggc ttgatcgata tctactcac gttatatacc cttaagtct
2041 aaccattatt atatttttg ataaataaaa algtcggct tctcgtgca ggcctggaa
2101 ttgttgacc cagctgaatg ccaaggacat cactgacgac gacaagcag cgtctgatga
2161 gctctgcac agagagggtac gtacatatta catttcacac cagggaatgc aagaacttta
2221 tcaagagaca ttactctt gatagagata gaatagaaca catgcacagt acacgtggac
2281 tcatgagctt gcaagacatc gacgacgaca cgtgtaagt agtgcgccag aaaaacttc
2341 aatttatatg tcaagtcagg tcaggtctc cattaanaac acatataat aatattcat
2401 tattatcaag ctaaggtaat aaacaacca actttccac tatttaact gctttgcaa
2461 actcaaatg aaaaactaac ctatcagga aagaactaga ctgcacatt atgtttac
2521 aatgcaatga gagaactgct acatgtataa cagaattat tatatgagc cgaattgact
2581 taagattcaa tgttgaagac cacttgatga aaactacact gaattatta tatgtatc
2641 tccagctgtg ctcaaaagc ttctttta cttaaaaaag atcatttgt acaaaagatc
2701 ctactcata tagaccatf ttagtagaac ttcggtacca cagatgcat aatgttttag
2761 ttgtaatcaa gttgtgtac tcatcattat atttctca caagtaggc atcagttc
2821 tcttgatga ggaatcaac ctgatagcc ttaactccac accctcaat agtaggcta
2881 tctcaagtt ctagtgtta caaattcag acgagtcata atgcatcac tgagcactc
2941 gtaaagagcg tctctctat ggtgcatata tatgatgag accacctgag aagtttactg
3001 ctcaagcca ccaagtgtt attttgtt tttgggtgt ttagtctaa ttctttct
3061 tgggtgtca cagatccaag cagcctcag aaccgatgaa accaggggg cacaacccc
3121 cccccaggac gaaatgcgct atgggatgag ctacatccat gagactgtat ggaagggcgt
3181 gcctaagttc ttgcgccgtg tggatacagc cctgaagaat atcggcata atgagcctc
3241 tccetacaat gtttctca ttcggttctc ttctggatg ggtggtagc gcatgggtac
3301 atttgcct accctttca ataaagtgc aggagctctc tcttttcag ctgagagaa
3361 acctctctg ttaacttga ctgcaataga tttcagaaa aactagtcta tcaattcgag

3421 ctctcaggag ctagaatttt aaaattgaaa ttatttagta cacctcacta alaaaaattt
3481 atcatccata catgctagca caacatataa gcataattta atcaaatctt tatattgcaa
3541 cctggaacct aactgttga atttttata tcacagaat atacctgtag taltatttta
3601 tatacaaaag agtgccttata ttatcactg actgtctctg tcaatattca agcctaactg
3661 tttttttt tcgccagaaa atlatatata cagaattata tgtttttt aagcctgtat
3721 atctttgcaa tctatcgcta tataggaaat ccaagagtta cccccggaggt gacaagagat
3781 gtafcttgc tggccagaat gatggctgca aactgttaca tcatcagat tgaagagctg
3841 atgtttgag tactgtacat ccatactgca gattgtttg atgaaatgt ctatgattt
3901 ttgttgcct tgtttttg tgtctccggt ccaaccaga acttcatgc atgcatctc
3961 tcatatct gtactctct atgtggcgt gcaacgatga gcttctgtt cgtgccgaag
4021 agtcccacag ttctcttgg tccaaagtta ccaagtatta cataggtaac cacaaacaga
4081 agcattatg ttgtttaat ttttccctgc cgtacagctt ttgcaaaag tctccactag
4141 tgttttcaa ttaattgg ggtcttttg gcatctttc tgaagtgat ttgctggcgc
4201 agaattctg aagcaaatc ctccaacga gccctaccg gtgataact gccatgtaag
4261 ggacaagctg tacaacacac gcgagcgtgc tgcctactg ctgcttctg gattttctga
4321 aatttcagcg gaatcgtcat ttaccagat cgaagaggta aatategica tglatattt
4381 atatatatt ataglatgac atcagcactg caactaaca aaaaaaatc actactgtc
4441 tcatgcatg cagtctctg agccactga gctgtgtct naatcactgt gtgactgccc
4501 cgacaaggcc atcgccgacg ggagcctctt ggacctctg cgcagagtt tcaagtccg
4561 gctctccctg tltgaagctgg acatccggca ggagtcggag cggcacaccg acgtgacga
4621 cggccalcacc acgcacctg gcatcgggtc gtaccgag tggccgagc acaagcggca
4681 ggagtggctg ctgtcggagc tgcgaggcaa gcgccctg ctgccccgg accctccca
4741 gaccgaggag atcgccgacg tcatcggcgc gttccactg ctgcccggag tcccggcca
4801 cagctgggc cctacatca tctcatggc gacggcccc tggagctgc tgccttggg
4861 gctctgcag cggagctgg cgtgtcggcc agccgtgcc tgggtgccg tgttgcgaag
4921 gctggccagc ctgagctgg cgcctcctg cgtggagcgc ctctctcgg tggactggta
4981 catggaccgg atcaaggcca agcagcaggt catgttggc tactccgact ccggcaagga
5041 cgcggccgc ctgtccgcg cgttgcagct gtacaggggc caggaggaga tggcgcaggt
5101 ggccaagcgc tacggcgta agctcactt gttccagcgc cgcggaggca ccgtggcag
5161 ggggtggcggg cccacgcacc ttgccactc gttccagcgc ccggacacca tcaacgggtc
5221 catccgtgtg acgggtcagg gcgaggtcat cgagtctgc ttcggggagg agcactgtg
5281 ctccagact ctgcagcgt tcaaggcgc cacgtggag cacggcatgc acccggcgt
5341 ctctcccaa cccgagtggc gcaagctcat ggacgagat gcggtctg tggccaggga
5401 gtaccctcc gtcgtgta aggaggccc ctctctgag tacttcat cggtatctg
5461 ceattgccca ttgcttgg acgatgaa tcatcctg cgtactt ttcatct
5521 tggagcttt gtgcgtact cactatcagg ctacaccgga gaccgagta gggaggatga
5581 acatggcag ccggccagcc aagaggagc ccggggcgg catcacgacc ctgcgcgca
5641 tcccctgga ctctctgtg accagacca ggttccact cccctgttgg ctgggagctg
5701 gcggccatt caagtgc atcgacaag acgtcaggaa ctccagctc ctcaaagaga
5761 tglacaacga gtggccaltc ttcagggtca cctggact gctggagat gttttgccca
5821 agggagacc ccgcattgcc gcttgtat acgagctgt tgtggcagaa gaactcaagc
5881 cttttgggaa gcagctcagg gacaatac tggagacaca gcagctctc ctccaggtac
5941 caaaaccagc actgcactgt acgatatgaa taaaagtct ttgtctgct cctgatcga
6001 gactgactac tccatttgg gcagatcgt gggcacaagg atattcttga agggatcca
6061 ttctgaagc agggactggt gctgcgcaac cctacatca ccaacctgaa cgtgttccag
6121 gcctacacg tgaagcggat aagggacccc aactcaagg tgaccgccc gcccgcctg

6181 tcaagagat tecegaaga gaacaagccc gccgaactgg taangctgaa cccggcgagg
 6241 gagtaccgce ccggcctgga agacaagctc atctctacca tgaaggcaat cgcgcgcggc
 6301 atgcagaaca ctggctaggg gctctctctt cactcaectg cagagtgcaac cgcacatc
 6361 agcttccgga tgggtggggt ttctcagtt tggatggaaa tgcgaactg gcagcgtctg
 6421 ttctccat gaatagtaa ttctgceet cttaatac aacctcttg tcaagtcctt
 6481 gtcgaaate ttggcatat atacatatt taataataa catcgaacn tctgcattg
 6541 gtttctaat aataaataa tctcccgacc catgttatgg acttcttcc catggtctta
 6601 ctccgcgcaac cctctcttag ttgtctaaa caattctga ttgctatt ttatctaga
 6661 gtaacctagt gcatttact aagagagatg atctctagtg gcactagtga ttgtttgca
 6721 agattagaaa ctgttaactc gctctctagag gtaaacata gcaatgtggt ggagctttag
 6781 g



1: AF271995 *Oryza sativa* phos [gi:9828444]
 LOCUS AF271995 3307 bp mRNA linear PLN 16-NUG-2000
 DEFINITION *Oryza sativa* phosphoenolpyruvate carboxylase mRNA, complete cds
 ACCESSION AF271995
 VERSION AF271995.1 GI:9828444
 KEYWORDS
 SOURCE *Oryza sativa* (japomea cultivar-group)
 ORGANISM *Oryza sativa* (japomea cultivar-group)
 Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
 Spermatophyta; Magnoliophyta; Liliopsida; Poales; Poaceae;
 Ehrhartoideae; Oryzaceae; *Oryza*
 REFERENCE 1 (bases 1 to 3307)
 AUTHORS Yamamoto,N., Kurita,A., Masumura,T., Sugimoto,T., Morita,S.,
 Shiraishi,N., Oji,Y. and Tanaka,K.
 TITLE Root type of phosphoenolpyruvate carboxylase in developing rice
 seeds
 JOURNAL Unpublished
 REFERENCE 2 (bases 1 to 3307)
 AUTHORS Yamamoto,N., Kurita,A., Masumura,T., Sugimoto,T., Morita,S.,
 Tanaka,K., Shiraishi,N. and Oji,Y.

TITLE Direct Submission
JOURNAL Submitted (24-MAY-2000) Lab. Plant Nutrition, Faculty of
Agriculture, Kobe University, Nada, Kobe 657-8501, Japan

FEATURES Location/Qualifiers

source 1..3307

/organism="Oryza sativa (japonica cultivar-group)"

/cultivar="Nipponbare"

/db_xref="taxon:39947"

/dev_stage="developing seeds"

CDS 188..3067

/codon_start=1

/product="phosphoenolpyruvate carboxylase"

/protein_id="AAGi00180.1"

/db_xref="GI:9828445"

/translation="MERHQSIDAQLRLLAPGKVS EDDKI.VEYDALLVDRFLDILQDLH
GPHLREFVQECYELSAEYENDRDEARLDEI.GRKL TSLPPGDSIVVSSSFHMLNLANI.
AEVQIAHRRRIKIKRGDFADEASAPTESDIEETLKRI.VTQLGKSREIEVFDALKNQTV
DLVFTAHTPTQSVRRSLIQKHGRIRNCLRQLYAKDITADDKQELDEALQREIQAAFRTD
EIRRPPTPQDEMRAGMSYFHETIWKGVPKFLRRIDTALKNIGINERI.PYNAPLIQFS
SWMGGDRDGNPRVTPEVTRDVCLLARMMAANLYFSQIEDIMFELSMWRCSDDELIRRAD
DLHCSSRKA AKHYIEFWKIQPPNEPYRVILGGVRDKLYYTRERTRHLITGVSEIPEE
ATFTNVEEFLEPLEICYRSLCACGDKPIADGSLLDFLRQVSTFGLALVKIDIRQESDR
HTDVLDIAITTYLGIGSYAEWSEEKRQDWI.LSELRGKRPI.FGPDLPQTEE IADVIGTFH
VI.AELPADCFGAYIISMATAPSDVIAVELLQRECHIKQPLRVVPLFEKLADIEAAPAA
VARLFSIDWYMNRINGKQEVMI GYSDSGKDAGRLSAAWQMYKAQEELVKVAKHYGVK
L
TMFHGRGGTVGRGGGPSHLAILSQPPDTHGSLRVTVQGEVIEHSFGEEHLCFRTLQR
FTAATLEHGMHPPISPKPEWRALMDEM AVVATKEYRSIVFKEPRFVEYFRSATPETEY
GRMNIGSRPSKRKPSGGIESLRAIPWIFAWTQTRFHLPVWLGFGGAFKHIMQKDIRNI
HTLKEMYNEWPFVRVTL DLEMVFAKGDPGIAALYDKLLVAGDLQSFGEQLRNNFEET
KQLLQVAGHKDILEGDPYLKQRLRLRESYITTLNVCQAYTLKRIRDPSFEVMSQPAL
SKEFVDSNQPAELVQLNAASEYAPGLEDTLLITMKGIAAGMQNTG"
BASE COUNT 816 a 814 c 880 g 797 t

ORIGIN

1 cccacactet ctaccccaca acttgcctc caatctatct ccgcgcgcgc cgcgcctcc
 61 ccgcgcgagc ggcgtgcgcc agctgcgcgc gaggcctccc tcccgcagca tetggccggga
 121 ttctcgtag tcaaggcagc tctaagcgc agagaggcag atggccgcgc ccgcgcggaa
 181 ggcggcgatg gageggcacc agtcgatcga cgcgcagctg cggctgcctc cgcgcgggaa
 241 ggtatccgag gacgacaagc tctgcgagta cgcgcctcct ctcgtgcacc gcttctctga
 301 catctctcag gacctgcacg gccctcactt ccgcgaaatc gtgcaggagt gctacgagct
 361 gtgcgcggag tacgugaacg acagggcagc ggcgcggctc gacgagctgg ggagggaagct
 421 caccagcctc ccgcgcgggg actecatctg cgtctccagc tctctcgc acatgctcaa
 481 ctgcgcaac ctgcgcgagg aggtgcagat cgcgcaccgc cgcgggatac agctcaagcg
 541 cgggggacttc gccgacgagg cctccgcgcc caccgagctc gacatcgagg agacgtctaa
 601 ggcctctct acccagctcg gcaagtcgag gggaggaggt ttgatgcgc tcaagaacca
 661 gaccgttagc ctctgttca ccgcgcacc aacgcagtc gtcaggaggt cctctgcca
 721 gaagcacggg aggatccgga attgcctgag gcagttatat gcgaagaca tcaactctga
 781 tgaacaagcag gacgttgatg aggcctcgca gaggggagat caagcagctt tcagaactga
 841 tgaatccgc aggactcctc ccactctca ggalgaaatg cgcgcgggga tgaattatt
 901 tcatgaaact atatggaagg gttaccaaa gttctgcgc cgcattgata ctctctgaa
 961 aaatattggg attaatgagc gttctctta caatgctct ctaaccagt tctatctg
 1021 gatgggtggg gaccgtgatg gaaacccaag agttaccaca gaggttaca gagatgtatg
 1081 ctgtttgca agaattgatg ctgtaacct gtaactctc cagatagaag atctgatgt
 1141 tgagctctct atgtggcgt gcagtgatga acttgaatt cgtgcagatg atttgcaatg
 1201 ctctccaga aaagctgca agcactatat agaattctgg aagcaaatc ctcaaatga
 1261 gccattact gtcaacttg gtgtgtcag gataaattg tactalacgc gtaaacgtac
 1321 tgcacatca ttgacaactg gagtttctga aattccagag gaggcaactt ttactaactg
 1381 tgaagagttt ctggagccgc ttgagctgtg ctacagatca ttatgtctt ttgggtgaca
 1441 acctatagct gatggaagcc ttctgattt ctgtctcaa gtaacactt ttgggtggc
 1501 tcttgaataa ctgacataa ggcaggagtc tgatcgacac actgatctc ttgatccat
 1561 aactacata ctggggattg gatcttacc tgaatggctt gaggagaaac gccaggattg
 1621 gctttgtcc gaactaaggg gcaaacgtcc ttgtttgtt cctgatctc ctacagctga
 1681 agaaattgct gatgtttag gaacattca cgtctctgt gactcccg gcagattgtt
 1741 tgggtcttac atcatctca ttgcaactgc accatctgat gtccttctg ttgagctct
 1801 acagcgggag tccatataa aacagcctct gcgagttgt ccctatttg agaaactgc
 1861 agacttgaa gcagctccag cagcggggc acgactatt tcaatagact ggtacatgaa
 1921 taggattaat gccaagcagg aggtcatgat tggatactca gactctgca aggatgtctg
 1981 tctgtttct gcagcatggc aaatgtacaa agcacaagag gactctgca aggttgcaaa
 2041 gcattatgtt glaaagtga caatgtcca tggaaagagt ggaacagtg gcagaggagg
 2101 tggccccagt catctgcca tattatctca accaccagac acgatacatg gatcaactc
 2161 tgltaacagta cagggcgagg ttattgaaca ctcttttga gaggaaact tgtctttag
 2221 aactctgac cgtttactg cagctactct tgagcatgga atgcactc caattccc
 2281 caaacagaa tggcgtgctc taatggatga gatggctgtt gttgcaaca aggaattcg
 2341 atcaactgc ttcaaggagc caeggttgt tgaatactc cgtcggca caactgagac
 2401 agaatatggc aggatgaaca ttgtagccg accatcaag agaaagccta gttgtggcat
 2461 tgaatcctc cgtgcaatc ctggatttt tcttggaca caaactaggt tcatctcc
 2521 tglatgctt ggatttgg ggcgcttca acatataat cagaaggata tcaggaaact
 2581 ccatcactg aaagaaatg acaatgagtg gccattctc aggttacc tggactctg
 2641 tgagatggtt ttgcctaagg gagatccagg aattgctct ttatgaca aattgctgt

2701 ggaggtgat ctgaatact ttggagaga gctgagatac aactttagg agnnaaaat
 2761 gctactcctt caggttccc gtcacaagg tttctggaa ggggactctt acctgaaga
 2821 gctcttcggg ctgctgtagt ctgacataac taacttgaac gtttgcacag cctacacct
 2881 gaagcggata agggaaacct gcttcggagt gattgcagc ccggccctgt cgttaggagt
 2941 cctctacaga aaccagctg ccgagctgtt caactgaa cctctcagct gctcagagc
 3001 ttgctctggg gacaccctc tctgacgat gaagggcatt gctgcggcga tgcagaaac
 3061 aggttagatt gtcacaaga atctctccag caggtgcatt caataatca taataact
 3121 ctccatgat gaggatgta atattgtaa ctatgactt cctccatgct agtgaattc
 3181 ggagacttt ttctctcc acattttt ttggttgtt actgatgta taatcaagg
 3241 ggggaattt ttgctttg ctgcctcgc caaataaaa tatgaccaa taattctag
 3301 aaaaaaa

1: AF399915 Sorghum bicolor p. [gi:15029404]
 LOCUS AF399915 1389 bp mRNA linear PLN 30-Jul-2001
 DEFINITION Sorghum bicolor phosphoenolpyruvate carboxylase kinase (PPCK)
 mRNA.

complete cds.

ACCESSION AF399915

VERSION AF399915.1 GI:15029404

KEYWORDS

SOURCE sorghum

ORGANISM [Sorghum bicolor](#)

Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;
 Spermatophyta; Magnoliophyta; Liliopsida; Poales; Poaceae; PACC;
 clade; Panicoideae; Andropogoneae; Sorghum

REFERENCE 1 (bases 1 to 1389)

AUTHORS Hartwell,J., Gib,A., Nimmo,G.A., Wilkins,M.B., Jenkins,G.L. and Nimmo,H.G.

TITLE Phosphoenolpyruvate carboxylase kinase is a novel protein kinase regulated at the level of expression

JOURNAL Plant J. 20 (3), 333-342 (1999)

MEDLINE 20049076

PUBMED 10571501

REFERENCE 2 (bases 1 to 1389)

AUTHORS Fontaine,V., Hartwell,J., Jenkins,G.I. and Nimmo,H.G.

TITLE Direct Submission

JOURNAL Submitted (12-JUL-2001) IBL, University of Glasgow, University
Avenue, Glasgow G12 8QQ, UK

FEATURES Location/Qualifiers

source 1..1389
/organism="Sorghum bicolor"
/db_xref="taxon:4558"
gene 1..1389
/gene="PPCK"
CDS 21..944
/gene="PPCK"
/note="protein kinase"
/codon_start=1
/product="phosphoenolpyruvate carboxylase kinase"
/protein_id="A_AK81871.1"
/db_xref="GI:15029405"

/translation="MSAE LKR DYE IGA EI GRG RFG VVH RCT SRATGEAF AVKSVDRSQ

LADDI DRELAELEPKLAQLAGAGNPGVVQTHAVYEDETWTHIVMDLCA GPDLLEWVR
L

RRGAPVPEPLAAAIVAQVAQALALCHRRGV AHRDVKPDNIIIDTAPAGSDDGDEDEED

SGEAEAEETAPRARIADFGSAAWVGAGGLGRAEGLVGTPHYVAPEVVGGGEYGA KAD
V

WSAGVVMYALLSGGALPFGGETA AAEVLA AVLRG SVRFP PRLFSGVSPA AKDILRRMIC
RDEWRRFTA EQVLAHPWIVS GGGARAIERPT"

BASE COUNT 263 a 393 c 507 g 226 t

ORIGIN

1 gcacgaggaa ttaactgaagt atgagtgccg agctgaagag ggactacgag ataggcgcgg
61 agatcggctc cggccgcttc ggggtggtcc accgctgcac gtccegcgcc accggcgagg
121 cgttccegt caagtccgtg gaccggtcgc agctggccga cgacctggac cgcgagctcg
181 cggagctgga gcccaagctg gcgcagctcg ccggcgcggg caaccccgcc gtgggtcaga
241 cgcacgcggt gtacgaggac gagacgtgga cccacacggt gatggacctg tgcgggggcc
301 cggacctgct cgagtgggtg cgcctccgcc gcggcgcgcc tgtgccggag cccctggcgg
361 ccgccatcgt cgcgcaggtc gcccaaggcc tcgcctctcg ccaccgcgcc gccgctgcgc
421 acccgacgt aaagcccgac aacatctca tcgacaccgc ccccgagggg agcgacgacg
481 gcgaggacga ggaagaggac agcggcgagg ccgaggaggc cgagacggcg ccgcgcgcgc
541 ggctggcgga ctccgggtcg gcggcgtggg tggggccggg tgggctgggc cgcgcggaag
601 gcttgggtgg gacgcccac tacgtggcgc ccgaggtggt gggcgccggc gagtacggcg
661 cgaagcgga cgtgtggagc gccggcgtgg tgatgtacgc gctgctctcc gccgcgcgc
721 tccgttcgg gcgcgagacc gcggcgaggg tctagcggc cgtgctcggg gccagcgtcc
781 ggttcccgcc cagctgttc agcggggtgt cccccgcgcc caaggacctg ttgagcgca

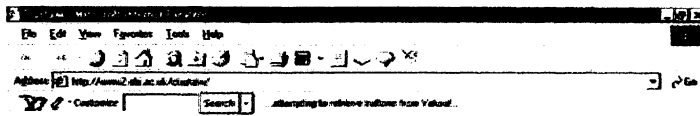
841 tgatctgccg cgacgagtgg cggaggttca ccgccgagca agtctctgct caccctgga
901 tcgtgagcgg cggaggagca cgggcaatcg agcggccaac ctgaggatga ggagggcata
961 ccataaccagc agaggaggag accatatggc gggcggacgt gtattgccgt agagacacag
1021 acacggcaag gaatcccgac gtcgttatgt atacggtttt gtatgtatgt atgtatgtac
1081 gtaeggacca gagagtagta gagtacggtg gctgctgccg aatgatgagt agcatcateg
1141 attccgcgct gtgtaggcgc agcagagcag atcgateggc aaagaagcft gtataggctt
1201 ctggttticc ttacttcggg gtttttggga tgcctttlac ctagcagtgc agagcaccct
1261 tttttggat gtccttgta agtggtaget ggatgtatat acactgaccg tctgtacaa
1321 atggagtgga gttggataaa taataaaaat cgagcccttg tacttccaaa aaaaaaaaaa
1381 aaaaaaaaaa

2.10.4 Exon Regions

Note the Exon regions of maize, *oryza sativa* and *sorghum bicolor* pyruvate carboxylase from the NCBI, which is given above.

2.10.5 Multiple Alignment

Make a Multiple alignment by putting the sequences (only Exons) of *Maize*, *oryza sativa* and *Sorghum bicolor* phosphoenol pyruvate carboxylase in Clustal W.



Clustalw

Run ClustalW
 Include bootstrap values

ClustalW options: NONE
 ClustalW options: clustalw_mpi
 ClustalW options: full
 ClustalW options: oh/nnumbers
 ClustalW options: clustmed

ClustalW options: no
 ClustalW options: del
 ClustalW options: del
 ClustalW options: percent
 ClustalW options: del
 ClustalW options: del
 ClustalW options: del
 ClustalW options: del

ClustalW options: none
 ClustalW options: off
 ClustalW options: on

TREE TYPE: cladogram
TREE GRAPH DISTANCES: hide

Some examples of a set of ClustalW input data files:

```

-BA11X
1587 gtcgc agagatggtgta agagaggtgtag ggcagctcttgc
1621 agagcawagag agagcagagag agagctggagag agctcggagc cagagctcagc agagctgagcc
1681 cagcagcagc cactcctctgag ggcagctccc tctcgcagctc gctcagcttc ggcagctcagc
1741 cagcagcagc gcagatctgag cagagctccc ggcagcagc gctcagagag agtgggtctg
1801 cagcagcagc ctctgagcagc agagagctcag agctcggagc ggcagctcagc agctcagctc
1861 cagcagctcagc agagcagcagc g
2082 gctcagcagc
2101 tctctgagc cagctgagctc cagagagctc cactcagcagc ggcagcagcagc agctcagctc
2161 gctcagcagc
  
```

Upload a file

This document was last modified on 12/06/2001 22:24:29

Comments or suggestions [E-mail us](#)

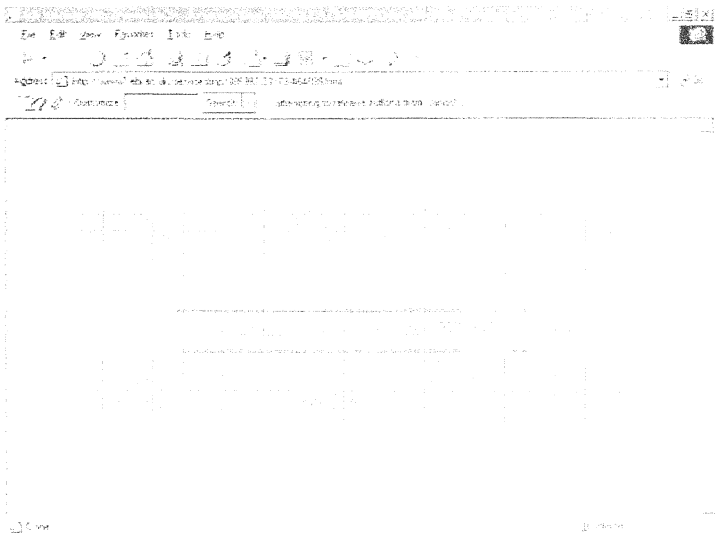
© EBI 2000



ClustalW Submission Form

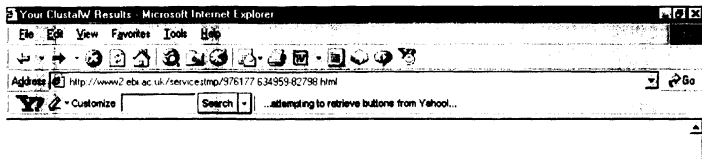
2.10.6 Run Clustalw

Click RUNCLUSTALW button, which is present on CLUSTALW. Job output will be displayed which is given below.



2.10.7 Jalview

Jalview will obtain after the job output by clicking on **RUN MULTISEQUENCE** button



Your ClustalW Results:

[View JalView:](#)

[JalView](#)

[SUBMIT ANOTHER JOB](#)

Pairwise Scores:

CLUSTAL W (1.01) Multiple Sequence Alignments

CLUSTAL W (1.81) Multiple Sequence Alignments

Sequence format is Pearson

Sequence 1: MAIZE 2339 bp

Sequence 2: ORYZA 2878 bp

Sequence 3: SORGHUM 924 bp

Start of Pairwise alignments

Aligning...

Sequences (2:3) Aligned, Score: 11

Sequences (1:2) Aligned, Score: 74

Sequences (1:3) Aligned, Score: 76

Guide tree file created: [/net/nfs0/vol1/production/w/nobody/tmp/976177634959-82798.dnd]

Start of Multiple Alignment

There are 3 groups

Aligning...

Group 1: Sequences: 2 Score: 35531

Group 2: Delayed

Sequence 3 Score: 8696

Alignment Score 16642

CLUSTAL-Alignment file created: [/net/nfs0/vol1/production/w/nobody/tmp/976177634959-82798.aln]

Your Multiple Sequence Alignment:

CLUSTAL W (1.81) multiple sequence alignment

```

MAIZE -----
ORYZA      ATGACCGCGGCAACAGTCGATCGACGGGACCTCGCGCTCCTCGCGGCGGGGAAGGTATCC 60
SORGHUM    -----

MAIZE -----
ORYZA      GAGGACGA/AAGCTGTGGAGTACGACGGCCCTCTGTGTCGACCGCTTCTCGACATCCTC 120
SORGHUM    -----

MAIZE -----
ORYZA      -----GTCACGAGTCTACGAGGTCGGCC 27
          CAGGACCTTGCACGGCCCTCACCTCCGCGAATTCGTGCAGGAGTCTCAAGCTGTGGCG 180

```


SORGHUM -----

MAIZE GACTATGAGGGCAAAGGAGACACGACGAAGCTGGGGGAGCTCGGCGCCAAGTTCACGGGG 87
 ORYZA GAGTACGAGAACGACAGGGACGAGGCGCGCTCGACAGCTCGGAGGAAAGTTCACGAGC 240
 SORGHUM -----

MAIZE CTGGCCCGCGCCGACGGCATCTCTCGTGGCGAGCTCCATCTCGCATGCTCAACCTCGCC 147
 ORYZA CTGGCCCGCGGGAACTCCATCTCGTCTCCAGCTCTTCTCGCATGCTCAACCTCGCC 300
 SORGHUM -----

MAIZE AACCTGGCCGAGGAGGTGCAGATCGCGCACCGCCGCCCAACAGCAAGCTCAAGAAAGGT 207
 ORYZA AACCTGCCCGAGGAGGTGCAGATCGCGCACCGCCGCCCAACAGCAAGCTCAAGAAAGGT 357
 SORGHUM -----

MAIZE GGGTTCGCGCAGGAGGCTCCGCCACCACCGAGTCCGACATCGAGGAGAGCTCAAGGCG 267
 ORYZA GACTTCGCGCAGGAGGCTCCGCCACCACCGAGTCCGACATCGAGGAGAGCTCAAGGCG 417
 SORGHUM -----

MAIZE CTCGTCTCGAGGTCGGCAAGTC----- 290
 ORYZA CTCGTCACTCAGCTCGGCAAGTCGAGGAGGAGGCTTCATGCGGCTCAAGAACGAGCC 477
 SORGHUM -----

MAIZE -----CCCCG----- 295
 ORYZA GTTGACTTCGTGTTGACCGGCATCCAACGAGTCCGTCAAGAGGTCCTGCTCCAAGAG 547
 SORGHUM -----

MAIZE -----GATCCGGAATTGCTGACCCAGCTGAATGCCAAGGATCACTGACGACGAC 347
 ORYZA CACGGGAGGATCCGGAATTGCTGAGGCAAGTATATGCGAAGACATCACTGCTGATGAC 597
 SORGHUM -----

MAIZE AAGCAGGAGCTCGATGAGCCTCTGCACAGAGATCCAAAGCAGCTTCAGAACCGATGAA 407
 ORYZA AAGCAGGAGCTTGATGAGGCCCTGCAGAGGAGATTCAAGCAGCTTCAGAACTGATGAA 657
 SORGHUM -----

MAIZE ATCAGGAGGGCACACCACCCCCCAGGACGAAATGGCCTATGGGATGAGCTACATCCAT 467
 ORYZA ATCCGAGGACTCCCGCACTCCTCAGGATCAATGGCGCCGGGATGAGTTATTTTCAT 717
 SORGHUM -----

MAIZE GAGACTGTATGGAAGGGGCTGCCCTAAGTTCTTGCGCCGTGTGGATACAGCCCTGAAGAAT 527
 ORYZA GAAACTATATGGAAGGGGTACCCTAAGTTCTTGCGCCGATGATACTGCTCGAAAAAT 777
 SORGHUM -----

MAIZE ATCGGCATCAATGAGCGCCTTCCCTACAATGTTCTCTCATCGGTTCTCTCTGGATG 587
 ORYZA ATTGGGATTAATGAGCGCTTCTCCTTACAATGCTCCTCTCATCGAGTTCTCATCTGGATG 837
 SORGHUM -----

MAIZE GGTGGTGAACCGGATGGAATCCAAGAGTTACCCCGGAGGTGACAAGAGATGTATGCTTG 647
 ORYZA GGTGGTGAACCGTGTGGAACCACAGAGTTACACCAGAGGTTACAAGAGATGTATGCTTG 897
 SORGHUM -----

MAIZE CTGGCCAGAAATGATGGCTGCAAACTTGTACATCGATCAGATTGAAGAGCTGATGTTTGAG 707
 ORYZA TTGGCAAGAAATGATGGCTGCTAACCTGTACTTCTCTCAGATAGAAGATCTGATGTTTGAG 957
 SORGHUM -----

MAIZE CTCTCTATGTGGCGCTGCAACGATGAGCTTCGTGTTCTGTGCCGAAGAGCTCCACAGTTCG 767
 ORYZA CTCTCTATGTGGCGCTGCAGTGATGAACCTCGAATTCGAATTCGTGCAGATGATTTGCATTGCTCC 1017
 SORGHUM -----

MAIZE	TCTGGTTCCAAAGTTACCAAGTATTACATAGAATTCCTGGAAGCAAATTCCTCCAAACGAG	827
ORYZA	TCCAG---AAAAGTCGCAAAGCACTATATAGAAATTCCTGGAAGCAAATTCCTCCAAATGAG	1074
SORGHUM	-----	
MAIZE	CCCTACCGGGGTGATACTAGGCCATGTAAGGGACAAAGCTGTACAAACACACGCCGAGCGTGCT	887
ORYZA	CCTTATCGTGCATACCTTGGTGGTGTCCGGGATAAAATTTGTACTATACGGCTGAACGTACT	1134
SORGHUM	-----	
MAIZE	CGCCATCTGCTGGCTTCTGGAGTTCTGAAATTCAGCGGAAATCGTCAATTTACCAGTATC	947
ORYZA	CGCCATCTATTGACAACCTGGAGTTCTGAAATTCAGAGGAGGCACCGTTTACTAATGTT	1194
SORGHUM	-----	
MAIZE	GAAGAGTTCCTTGAGCCACTTGAGCTGTGCTACAAACTACTGTGTGACTGCGGGGACAAAG	1007
ORYZA	GAAGAGTTTCTCGAGCCGCTTGAGCTGTGCTACAGATCATFATGTGCTTGTGGTGACAAA	1254
SORGHUM	-----	
MAIZE	GCCATCGCGGACGGGAGCCCTCCTGGACCTCTGCGCCAGGTTTTCAAGTTCCGGGCTCTCC	1067
ORYZA	CCTATAGCTGATGGAAGCCCTCTTGATTCTTGGCTAAGTATCAACTTTTGGGCTGGCT	1314
SORGHUM	-----	
MAIZE	CTGGTGAAGCTGGACATCGGCCAGGAGTCGGAGCGGCACACCGAGCTGATCGACGCCATC	1127
ORYZA	CTTGTAAAACCTTGACATAAGGCCAGGAGTTGATCGACACACTGATGTCCTTGATGCCATA	1374
SORGHUM	-----	
MAIZE	ACCACGCACCTCGGCATCGGGTCTACCGGAGTGTGCGGAGGACAAAGCCGAGGAGTGG	1187
ORYZA	ACTACATATCTTGGGATTGGATCTTACGCTGAATGCTCTGAGGAGAAACGCCAGGATTGG	1434
SORGHUM	-----	
MAIZE	CTGCTGTCCGAGCTGCGAGGCAAGCCGCCGCTGCTGCCCCGGACCTTCCACAGCCGAC	1247
ORYZA	CTTTTGTCCGAACTAAGGGGCAACGTCCTTTGTTTGGCTGATCTTCTCAGACTGAA	1494
SORGHUM	-----	
MAIZE	GAGATCCCGACGTCATCGGGCGGTTCCACGCTCCCGGGAGCTCCCGCCGACAGCTTC	1307
ORYZA	GAAATTTGCTGATGTTTTAGGAACATTTCACTCCCTTGCTGAGCTCCGGCAGATGTTTT	1554
SORGHUM	ATGAGTGGCGAGCTGAAGAGGGACTACGAGATAGGCGCGGAGATCGGTCCGGCCGCTTC	60
	* * * * *	
MAIZE	GGCCCCACATCATCTCCATGGCGACGGCCUCCCTCGGACGTGCTCGCCGTGGAGCTCTG	1367
ORYZA	GGTGCCTACATCATCTCAATGGCAACTGCACCATCTGATGTGCTTGTCTGGAGCTTCTA	1614
SORGHUM	GGGSGTGTCCACCGTGCACGTCGCGGUCACCGGCGAGGCGTTCGGCTCAAGTCCGTG	120
	** * * * *	
MAIZE	CAGCGGAGTG--CGGCGTGGCGGCAGCCGTGCCCG--TGCTCCCGCTTTCGAAAGGCT	1424
ORYZA	CAGCGGAGTG--CCATATAAACCAGCCTCTCGGAG--TTCTTCCCTATTTGAGAAGCTT	1671
SORGHUM	GACCGTTCAGCTGGCGGACGACTCGGACCGGAGCTCGCGGAGCTGGAGCCGAAAGCTG	180
	* * * * *	
MAIZE	GCCAGCCTCGACTCGGCGCCCGCGTCCGTGGAGCGCCTCTTCTCGGTGGACTGGTACATG	1484
ORYZA	GCAGATCTTGAAGCAGCTCCAGCAGCGGTGGCACGACTATTTTCAATAGACTGGTACATG	1731
SORGHUM	GCGCACCTC--GCCGCGCGGGCAACCCCGGCTGTGCGAGCGCACCGCGGTG--TACCAG	237
	** * * * *	
MAIZE	GACCGGATCAAGGGCAAGCAGCAGGTGATGTTGCGCTACTCCGACTCUGGCAAGGACGC-	1543
ORYZA	AATAGGATTAATGGCAAGCAGGAGGTGATGTTGATCTCAGACTCTGGCAAGGATGC-	1790
SORGHUM	GACGAGAC---GTGGACCCACACGGTGTGACCTGTGCGCGGGCCCGACCTGTCTCGAG	294
	* * * * *	
MAIZE	CGGCGCCTGTCCGCGCGGTGGCAGCTGTACAGGGCGCAGGACGAGATGGCCAGGTGGC	1603
ORYZA	TGGTCTCTTTCTGCAGCATGGCAAATGTACAAGGACAAAGAGGAGCTTGTCAAGGTGGC	1850
SORGHUM	TGGTGGCGCCTTCGCGCGCGCGGCTGTGCGGGA--GCCCTTGGCGGCGCCATCGTCCG	353
	** * * * *	

MAIZE	CAAGCGCTACGGCGCTCAAGCTCACCTTGTTCACGGCCGCGGAGGCACCGTGGCGAGGG- 1662
ORYZA	AAAGCATTATGGTFAAAGTTGACAATGTTCCATGGAAAGAGTTGGAACAGTTGGCAGAG- 1909
SORGHUM	---CGAGGTGCGCCAGCGCTCGCGCTTGCCACCGCGCGGGCTGGCCACCGGGACGT 410
	* * * * *
MAIZE	GTGGCGGGCCACCGACCTTGGCATCCTGTCCAGCCGCGGSACACCATCAACGGGTCCA 1722
ORYZA	GAGGTGGTCCAGTTCATCTTGCCATATTATCTCAACCACCAGACACGATACATGGATCAC 1969
SORGHUM	AAAGCCCGACAACATCCTCATCGACACCGCCCGCCGAGGGAGCGACCGCGGAGGACA 470
	* * * * *
MAIZE	TCCGTGTGACGGTGC--AGGGCGAGGTGATCGAGTTTCGTTCGGGGAGGAGCACCTGTG 1780
ORYZA	TTCCGTGTAACACATAC--AGGGCGAGGTTATTGAACACTTCGTTGGAGAGGAACTTTGTG 2027
SORGHUM	GGAAAGAGCACCGCGCCAGCGGAGGAGGCGAGACGGCGCGCGCGGGTGGCGGA 530
	* * * * *
MAIZE	CTTCAGACTCTGCAGCGCTTACGGCCGCGCCAGCTGGAGCACGGCATGCACCCGCGGT 1840
ORYZA	CTTTAGAACTCTGCAGCGCTTACTGCACTACTCTTGAGCATGGAATGCATCCTCCAAT 2087
SORGHUM	CTTCGGTTCGGCGGGCTGGTGGCGCCGCTGGCTGGGCCTGGGCCTGGCGGAAAGCCCTGGTGGG 590
	* * * * *
MAIZE	CTCTCCCAAGCCGAGTGGCCGAAGCTCATCGACGAGATGGCGGTCTGTGCCACGGAGGA 1900
ORYZA	TTCCCCCAACCAGAAATGGCGTCTTAATGGATGAGATGGCTGTTGTGGCAACAAAGGA 2147
SORGHUM	GACGCCCTACTAGTGGCGCCGAGGTGGTGGGGCG--CGCGGAGTACGGC-GCGAAGGC 647
	* * * * *
MAIZE	GTACCGTCCGTCGTCGTCGAAGGAGCCCGCTTCGTGGAGTACTTCAGATCGGTACACC 1960
ORYZA	ATATCGATCAATCGTCTTCAAGGAGCCACGGTTTGTGAATACTTCGGATCGGCAACACC 2207
SORGHUM	GGACGTGTGAGGCGCGCGGTGGTGATTTACCGCTGTCTCTCGCGCGCGCGTCCCGTT 707
	* * * * *
MAIZE	GGAGACCGAGTACGGGAGGATGAACATCGGCAGCCGCGCAGCCAAAGAGGAGGCCCGCGGG 2020
ORYZA	TGAGACAGAATATGGCAGGATGAACATGGTAGECGCACCATCAAGAGAAAGCCTAGTGG 2267
SORGHUM	CGCGGCGGAG-ACCRCGCGGAGGTGCTACCGCGCTGCTCGGGGACCGCTCCGGTTCC 766
	* * * * *
MAIZE	CGGCATCAGGACCCTCGCGCCATCCCTGGATCTTCTCGTGGACCGACAC-CAGGTTCC 2079
ORYZA	TGGCATTGAATCGCTCCGTGCAATTCCCTTGGATTTTTCGTTGACACAAAC-TAGGTTTC 2326
SORGHUM	CGCCGAGGCTGTTACGGCGGGTGTCCCGCGCGCCAAAGACCTGTTGAGGGCCGATGATC 826
	* * * * *
MAIZE	ACCPCCGGTGTGGTGGGAGTCGGCGCCGATTCAAGTTGCCATCGACAAGGAGCTCA 2139
ORYZA	ATCTTCTCTATGGCTTGGATTTGGTGGAGCGTTCAAAACATATAATGCAAGAAGGATATCA 2386
SORGHUM	CGCGGACGAGTGGCGGAGGTTACCGCGGAG--CAAGTCC-TGCGTACCCCTGGATCG 883
	* * * * *
MAIZE	GGAACTCCAGGTCCTCAAAGAGATGTACAACAGATGCCCATCTTTCAGGGTCCACCTGG 2199
ORYZA	GGAACTCCATACACTGAAAGAAATGTACAAATGAGTGGCCATTTTCAGGCTCACCTTGG 2446
SORGHUM	TGAGCGCGGAGGAGCACGGGCAAT-----CGAGCGGCCACCTGA----- 924
	* * * * *
MAIZE	ACCTGCTGGAGATGGTTTTCCCAAGGGAGACCCCGGATTTGCCGGTTCGTATGACGAGC 2259
ORYZA	ACTTTGCTTGAGATGGTTTTCGCTAAGGGAGATCCAGGAATGCTGCTTTATATGACAAAT 2506
SORGHUM	-----
	* * * * *
MAIZE	TGCTTGTGGCAGAAGAACTCAAGCCCTTTGGGAAGCAGCTCAGGGACAAATACGTGGAGA 2319
ORYZA	TGCTTGTGGCAGGTGATCTGCAATCCTTTGGAGAGCAGCTGAGAAACAACTTTGAAGAGA 2566
SORGHUM	-----
	* * * * *
MAIZE	CACAGCAGCTTCTCCTCCAG----- 2339
ORYZA	CAAAACAGCTACTCCTTCAGGTTGCCGTCACAAGGATATCTGGAAGGGGATCCTTACC 2626
SORGHUM	-----
	* * * * *
MAIZE	TGAGCAGCGTCTGCGGCTGCGTACGTACATCACTACCTTGAACCTTTGCCAAGCCT 2686
ORYZA	-----
SORGHUM	-----

MAIZE -----
 ORYZA ACACCCCTGAAGCGGATAAGGGACCCTAGCTTCGAGGTGATGTGGCAGCCGGCCCTGTCGA 2746
 SORGHUM -----

MAIZE -----
 ORYZA AGGAGTTCGTCGACAGCAACCACCCTGCGGAGCTGGTCCAACTGAACGCTGCGAGCGAGT 2806
 SORGHUM -----

MAIZE -----
 ORYZA ACCCACCTGGCCCTGGAGGACACCCTCATCTTGACGATGAAGGGCATTCCTGCCGGCATGC 2866
 SORGHUM -----

MAIZE -----
 ORYZA AGAACACAGGCT 2876
 SORGHUM -----

Your guide tree:

(
 MAIZE:0.03716,
 ORYZA:0.11465,
 SORGHUM:0.669551);

5

2.10.8 Alignment Graph

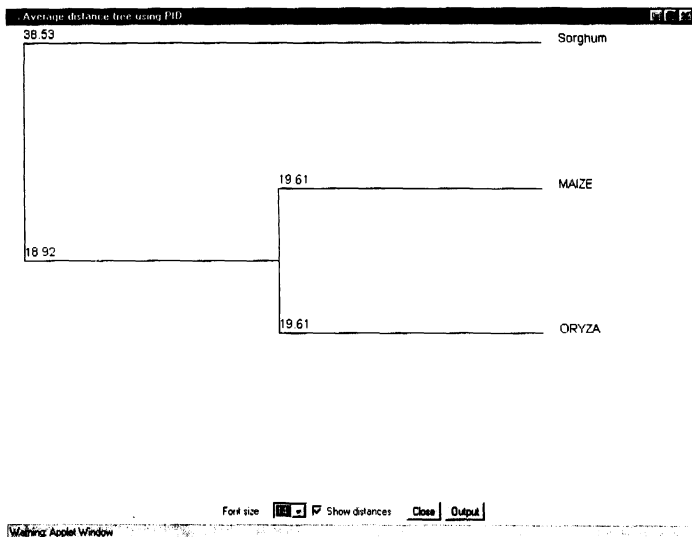
An alignment graph among Maize, oryza sativa and Sorghum bicolor phosphoenol carboxylase gene will obtain.

630 640 650 660 670 680 690 700 710 720 730 740 750 760 770 780 790
A
C
E
G
H
I
J
K
L
M
N
O
P
Q
R
S
T
U
V
W
X
Y
Z
[REDACTED]

790 800 810 820 830 840 850 860 870 880 890 900 910 920 930 940 950
[REDACTED]

950 960 970 980 990 1000 1010 1020 1030 1040 1050 1060 1070 1080 1090 1100 1110
[REDACTED]

1110 1120 1130 1140 1150 1160 1170 1180 1190 1200 1210 1220 1230 1240 1250 1260 1270
[REDACTED]



2.10.10 Tree Analysis

Check the two species which having maximum alignment By tree Analysis. Take the species, which have maximum alignment. By Tree Analysis it has been found that there is maximum alignment between Maize and oryza sativa phosphoenol pyruvate carboxylase gene.

2.10.11 Design Primers For The Sequences Of Maize:

Design two sets of Left and Right Primers for the exon regions of maize sequences by using PRIMER3. One set of primers are sufficient for PCR, but the second set is taken as additional primer in the event the first fails hybridize.

2.10.12 Selection Of First Set Of Primers

Design left primer from one Exon of the sequence and right primer from another Exon of the sequence. The distance between the two primers should be about 500 to 1000 base pairs. Take the sequence from the species, which has both Exons and introns. i.e we have taken Maize sequence because it has both Exons and introns. But rice sequence for this enzyme has only coding regions i.e.Exons.so we can't design primer from rice for sequencing of DNA but can be used as a marker.

2.10.12.1 Primer3 For Left Primer

Take one Exon sequence of the maize and then paste it in PRIMER3 and pick primer. The region of the exon from which the left primer is designed are 4454-5452. The left primer is picked from the Exons, which is right of the exon from which the right primer is picked. The parameters for both the primers should be same. the product size should be Min: 100, Opt: 200, Max: 1000. Primer size should be Min: 18bp, Opt: 20bp, Max: 27bp. Primer annealing temperature should be Min: 57, Opt: 60, Max: 63 degree celcius. The GC% should be Min: 20 and Max: 80.

REPORT OF THE...

Printers

... ..

... ..

Back	Pick	Back
1500	1500-2000	1500
1505	1505	1505
1510	1510	1510
1515	1515	1515
1520	1520	1520
1525	1525	1525
1530	1530	1530
1535	1535	1535
1540	1540	1540
1545	1545	1545
1550	1550	1550
1555	1555	1555
1560	1560	1560
1565	1565	1565
1570	1570	1570
1575	1575	1575
1580	1580	1580
1585	1585	1585
1590	1590	1590
1595	1595	1595
1600	1600	1600

Back Printer

... ..

Min	Max	Min	Max
1500	1600	1500	1600
1505	1605	1505	1605
1510	1610	1510	1610

General Primer Picking Conditions

Min	Max	Min	Max
1500	1600	1500	1600
1505	1605	1505	1605
1510	1610	1510	1610
1515	1615	1515	1615
1520	1620	1520	1620
1525	1625	1525	1625
1530	1630	1530	1630
1535	1635	1535	1635
1540	1640	1540	1640
1545	1645	1545	1645
1550	1650	1550	1650
1555	1655	1555	1655
1560	1660	1560	1660
1565	1665	1565	1665
1570	1670	1570	1670
1575	1675	1575	1675
1580	1680	1580	1680
1585	1685	1585	1685
1590	1690	1590	1690
1595	1695	1595	1695
1600	1700	1600	1700

Other Pec-Sequences Inputs

... ..

Objective Function Penalty Weights for Primers:

Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq

Objective Function Penalty Weights for Primer Pairs

Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq

Hyb Oligo (Internal Oligo) Per-Sequence Inputs

Seq	Seq
-----	-----

Hyb Oligo (Internal Oligo) General Conditions

Max	Seq	Max	Seq
Max	Seq	Max	Seq
Min	Seq	Min	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq
Seq	Seq	Seq	Seq

Objective Function Penalty Weights for Hyb Oligos (Internal Oligos)

Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq
Seq	Seq	Seq

2.10.12.2 Primer3 Output For Left Primer.

Put the sequence of Exon of the maize for the left primer in the primer3 output and by clicking pick primer the left primer and right primer both will be obtained. But, left primer only has to be taken into consideration. Primer3 output for left primer are given below.

Primer3 Output

WARNING: Numbers in input sequence were deleted.

```
No mispriming library specified
Using 1-based sequence positions
OLIGO
LEFT PRIMER      663 20  60.83  55.00  7.00  1.00  cgtcaagctcaaccttcttc
RIGHT PRIMER     843 20  60.33  55.00  7.00  0.00  ctgagagatctggaagcaca
SEQUENCE SIZE: 999
INCLUDED REGION SIZE: 999
```

PRODUCT SIZE: 181, PATR ANY COMPL: 4.00, PATR 3' COMPL: 3.00

```
1  ttccttgagccacttgagctgtgtctacaactcactgtgtgacttcgggcacaagggccatc
61  ggggacgggagcctcctggacctcctggccaggctttccacttgggctctccctggtg
121  aagctggacatccggcaggagtccggagcggcacccagaagtgatecagcccatcaccag
181  cacctcggcatccgqctcgtaccgcagtggtccaggacaagcggcaggagtggtgctg
241  tcggagctgcaggcaagcgcccgcctgtgcctccggaccttcctccagaccaggagatc
301  gcgcagctcatcggcgcttcaccgtcctcggaggctcccgcccagacgctcggccc
361  tacatcatctccatggcgacggccccctcggagctgctcgcctggagctcctgcagcgc
421  gagtgcggcgctgcggccagccgtgcctcgggtgcctggttcgaaaggtcgccagcctg
481  cagtcggcgcccgcctccgtggagcgctctctctcggtgactggtacatggaccggatc
541  aaggccaagcagcaggctcatggtcngtacctccgactccggcaaggcagcggcggcctg
601  tcggcggcgtggcagctgtacaggcgcaggaggagatggcgagggtggcacaagcgtac
661  ggctcaagctcaccttgctccacggccggaggcaccgtggcgagggtagggggccc
    >>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>
721  accgaccttgacatcctgtcccagccgcggacaccatcaacgggtccatccgtgtgac
```


Form 101 (Rev. 1-15-60)

Part I - Personal Information

Name (Last, First, Middle Initial)
 Social Security Number
 Date of Birth (Month, Day, Year)
 Place of Birth (City, State, Country)
 Present Address (Street, City, State, Zip)
 Present Telephone Number (Area Code, Number)
 Present Employer (Name, Address, City, State, Zip)
 Present Position (Title, Grade, Step)

Education (School, Degree, Date)
 Experience (Employer, Position, Dates)
 Other Information (Languages, Skills, etc.)

Remarks (Additional comments, references, etc.)

Part II - General Position Picking Conditions

1. Position Title: _____
 2. Position Grade: _____
 3. Position Step: _____
 4. Position Series: _____
 5. Position Family: _____
 6. Position Branch: _____
 7. Position Division: _____
 8. Position Office: _____
 9. Position Location: _____
 10. Position Status: _____
 11. Position Type: _____
 12. Position Category: _____
 13. Position Subcategory: _____
 14. Position Grade Range: _____
 15. Position Step Range: _____
 16. Position Series Range: _____
 17. Position Family Range: _____
 18. Position Branch Range: _____
 19. Position Division Range: _____
 20. Position Office Range: _____
 21. Position Location Range: _____
 22. Position Status Range: _____
 23. Position Type Range: _____
 24. Position Category Range: _____
 25. Position Subcategory Range: _____

Other Per-Sequence Inputs

1. Position Title: _____
 2. Position Grade: _____
 3. Position Step: _____
 4. Position Series: _____
 5. Position Family: _____
 6. Position Branch: _____
 7. Position Division: _____
 8. Position Office: _____
 9. Position Location: _____
 10. Position Status: _____
 11. Position Type: _____
 12. Position Category: _____
 13. Position Subcategory: _____
 14. Position Grade Range: _____
 15. Position Step Range: _____
 16. Position Series Range: _____
 17. Position Family Range: _____
 18. Position Branch Range: _____
 19. Position Division Range: _____
 20. Position Office Range: _____
 21. Position Location Range: _____
 22. Position Status Range: _____
 23. Position Type Range: _____
 24. Position Category Range: _____
 25. Position Subcategory Range: _____

Objective Function Penalty Weights for Primers

Primer 1	Min	Max	Primer 2	Min	Max
GC Content	40	60	GC Content	40	60
GC Bias	0	10	GC Bias	0	10
Primer Length	20	30	Primer Length	20	30
Primer GC Content	40	60	Primer GC Content	40	60
Primer GC Bias	0	10	Primer GC Bias	0	10
Primer Melting Temp	55	65	Primer Melting Temp	55	65
Primer Self-Complementarity	0	20	Primer Self-Complementarity	0	20
Primer Dimerization	0	10	Primer Dimerization	0	10
Primer Hairpins	0	10	Primer Hairpins	0	10
Primer Secondary Structure	0	20	Primer Secondary Structure	0	20
Primer Specificity	0	10	Primer Specificity	0	10
Primer Stability	0	10	Primer Stability	0	10

Objective Function Penalty Weights for Primer Pairs

Primer 1	Min	Max	Primer 2	Min	Max
GC Content	40	60	GC Content	40	60
GC Bias	0	10	GC Bias	0	10
Primer Length	20	30	Primer Length	20	30
Primer GC Content	40	60	Primer GC Content	40	60
Primer GC Bias	0	10	Primer GC Bias	0	10
Primer Melting Temp	55	65	Primer Melting Temp	55	65
Primer Self-Complementarity	0	20	Primer Self-Complementarity	0	20
Primer Dimerization	0	10	Primer Dimerization	0	10
Primer Hairpins	0	10	Primer Hairpins	0	10
Primer Secondary Structure	0	20	Primer Secondary Structure	0	20
Primer Specificity	0	10	Primer Specificity	0	10
Primer Stability	0	10	Primer Stability	0	10
Primer Pair GC Content	40	60	Primer Pair GC Content	40	60
Primer Pair GC Bias	0	10	Primer Pair GC Bias	0	10
Primer Pair Melting Temp	55	65	Primer Pair Melting Temp	55	65
Primer Pair Self-Complementarity	0	20	Primer Pair Self-Complementarity	0	20
Primer Pair Dimerization	0	10	Primer Pair Dimerization	0	10
Primer Pair Hairpins	0	10	Primer Pair Hairpins	0	10
Primer Pair Secondary Structure	0	20	Primer Pair Secondary Structure	0	20
Primer Pair Specificity	0	10	Primer Pair Specificity	0	10
Primer Pair Stability	0	10	Primer Pair Stability	0	10

Hyb Oligo (Internal Oligo) Per-Sequence Inputs

Parameter	Min	Max	Unit
GC Content	40	60	%
GC Bias	0	10	%
Primer Length	20	30	bp
Primer GC Content	40	60	%
Primer GC Bias	0	10	%
Primer Melting Temp	55	65	°C
Primer Self-Complementarity	0	20	%
Primer Dimerization	0	10	%
Primer Hairpins	0	10	%
Primer Secondary Structure	0	20	%
Primer Specificity	0	10	%
Primer Stability	0	10	%

Hyb Oligo (Internal Oligo) General Conditions

Parameter	Min	Max	Unit
GC Content	40	60	%
GC Bias	0	10	%
Primer Length	20	30	bp
Primer GC Content	40	60	%
Primer GC Bias	0	10	%
Primer Melting Temp	55	65	°C
Primer Self-Complementarity	0	20	%
Primer Dimerization	0	10	%
Primer Hairpins	0	10	%
Primer Secondary Structure	0	20	%
Primer Specificity	0	10	%
Primer Stability	0	10	%

Objective Function Penalty Weights for Hyb Oligos (Internal Oligos)

Parameter	Min	Max
GC Content	40	60
GC Bias	0	10
Primer Length	20	30
Primer GC Content	40	60
Primer GC Bias	0	10
Primer Melting Temp	55	65
Primer Self-Complementarity	0	20
Primer Dimerization	0	10
Primer Hairpins	0	10
Primer Secondary Structure	0	20
Primer Specificity	0	10
Primer Stability	0	10

2.10.12.4 Primer3 Output For Right Primer

By clicking pick primer we will get the right primer and left primer both. but right primer only has to be consider. primer3 output for right primer are given below.

Primer3 Output

WARNING: Numbers in input sequence were deleted.

No mispriming library specified
Using 1-based sequence positions

OLIGO							
LEFT PRIMER	1	20	60.28	60.00	4.00	3.00 gctacacccggagaccggagta	
RIGHT PRIMER	197	20	59.70	50.00	7.00	3.00 tggaaattctctgacgtcctt	

SEQUENCE SIZE: 387
INCLUDED REGION SIZE: 387

PRODUCT SIZE: 197, PAIR ANY COMPL: 3.00, PAIR 3' COMPL: 1.00

```

1 gctacacccggagaccggagtaaacatccggaccgaccagcraaaagagaaq
>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>>
61 cccggcggcgggcaccagaccctgcgggcacatccctcggaatctctctgtgaccccaaac

121 aggttccacctcccgggtgtggctgggagtcggcgcgcgcattcaagttgccatgacaag

181 gacgtcaggaacttccagttccraaaagaaatgtacaacgagtgcgcacattcttcagggtc
<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<<

241 acctcggaccctgctggagatgggtttttcgcacagggagaccaccggcatttccctgcttat

301 gacgagctgcttctggccaaagaaactcaagccctttgggaagcagctcagggacaataac

361 gtggacacacagcagcttctctccag

```

KEYS (in order of precedence):
>>>>> left primer
<<<<< right primer

ADDITIONAL OLIGOS

1 LEFT PRIMER	1	20	60.28	60.00	4.00	3.00 gctacacccggagaccggagta
RIGHT PRIMER	196	20	59.70	55.00	7.00	1.00 ggaagtctctgacgtccttg
PRODUCT SIZE: 196, PAIR ANY COMPL: 3.00, PAIR 3' COMPL: 2.00						
2 LEFT PRIMER	174	20	60.30	55.00	6.00	1.00 cgacaaggacgtcaggaact
RIGHT PRIMER	377	20	59.65	55.00	5.00	2.00 agctgctgtgtctccacgta
PRODUCT SIZE: 204, PAIR ANY COMPL: 4.00, PAIR 3' COMPL: 2.00						
3 LEFT PRIMER	1	20	60.28	60.00	4.00	3.00 gctacacccggagaccggagta
RIGHT PRIMER	193	20	60.30	55.00	6.00	2.00 agttccctgagctccttgctg
PRODUCT SIZE: 193, PAIR ANY COMPL: 3.00, PAIR 3' COMPL: 2.00						
4 LEFT PRIMER	186	20	60.23	55.00	7.00	3.00 caggaacttccaggtcctca
RIGHT PRIMER	377	20	59.65	55.00	5.00	2.00 agctgctgtgtctccacgta
PRODUCT SIZE: 192, PAIR ANY COMPL: 4.00, PAIR 3' COMPL: 0.00						

```

Statistics
      con   too   in   in   no   tm   tm   high   high
      sid  many  tar  excl  bad  GC  GC  low  high  poly  high
      ered  Ns  get  reg  GC  clamp  low  high  compl  compl  X  stab  ok
Left   2541  0   0   0   86  0   144  1972  0   10  0   41  288
Right  2438  0   0   0   7   0   214  1728  1   8   0   52  418
Pair Stats:
considered 224, unacceptable product size 159, high end compl 13, ok 52
primer3 release 0.9

```

```
(primer3_www_results.cgi v 0.2)
```

Note down the right primer from the above PRIMER3 output.

2.10.12.5 Calculation Of Product Size

Left primer sequence size is 999

Left primer starting point is 663

Left primer total length is 336

Introns length is 98

Right primer sequence size is 387

Right primer ending point is 197

Total right primer length is 190

Product size = Left primer total length is (336) +introns(98) +Right primer ending point is (197) = 631

2.10.13 Selection Of Second Set Of Primers

Repeat the same process as explained in the first set of primers. Set the parameters for designing the primers for PCR.

2.10.13 .1 Primer3 For Left Primer

The parameters set to design the left primer for second set of primers are: the product size should be Min: 100, Opt: 200, Max: 1000. the primer annealing temperature is Min: 49, Opt: 51, Max: 54. the GC% is Min: 20 and Max: 80. the primer size is Min: 18bp, Opt: 20bp, Max: 27bp. The exon regions which are taken to design the left primer are 3074-3296.

The following information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

This information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

Code	Description	Value
01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20

The following information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

This information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

General Printer Picking Conditions

The following information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

This information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

Other Post-Sequence Inputs

The following information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

This information is provided for your information only. It is not intended to be used as a substitute for professional advice. The information is provided for your information only. It is not intended to be used as a substitute for professional advice.

2.10.13.2 Primer3 Output For Left Primer.

Primer3 Output

WARNING: Numbers in input sequence were deleted.

No mispriming library specified
Using I-based sequence positions

OLIGO

	len	len	TM	GC	TM	TM	high	high	high
LEFT PRIMER	12	18	50.69	38.89	5.00	2.00	cttccagaacccgatgaaat		
RIGHT PRIMER	191	20	51.23	35.00	5.00	3.00	cgaaatgagagaacattgta		

SEQUENCE SIZE: 223
INCLUDED REGION SIZE: 223

PRODUCT SIZE: 180, PAIR ANY COMPL: 3.00, PAIR 3' COMPL: 2.00

```
1 atccaaagcagcccttcagaacccgatgaaatcaggagggcacaacccccccaggaacaa
    >>>>>>>>>>>>>>>>
61 atgcqctatgggatgagctacatccatgagactgtatgaaqqqgcctgccaaagtctcttq
121 cgccgtgtggatcacgcctgaagaatatacgccatcaatgagcgccctccctacaattgt
181 tctctcattcgggtctctctctttqqatgggtggtgacccgcatg
    <<<<<<<<<<<<<<<<
```

KEYS (in order of precedence):

>>>>> left primer
<<<<< right primer

ADDITIONAL OLIGOS

	len	len	TM	GC	TM	TM	high	high	high
1 LEFT PRIMER	12	18	50.69	38.89	5.00	2.00	cttccagaacccgatgaaat		
RIGHT PRIMER	185	20	51.05	40.00	5.00	0.00	agagaaacattgtaggggaq		
PRODUCT SIZE:	174								
PAIR ANY COMPL:			4.00						
PAIR 3' COMPL:					0.00				
2 LEFT PRIMER	13	18	51.55	38.89	5.00	3.00	ttcagaacccgatgaaatc		
RIGHT PRIMER	191	20	51.23	35.00	5.00	3.00	cgaaatgagagaacattgta		
PRODUCT SIZE:	179								
PAIR ANY COMPL:			3.00						
PAIR 3' COMPL:					1.00				
3 LEFT PRIMER	13	18	51.55	38.89	5.00	3.00	ttcagaacccgatgaaatc		
RIGHT PRIMER	185	20	51.05	40.00	5.00	0.00	agagaaacattgtaggggaq		
PRODUCT SIZE:	173								
PAIR ANY COMPL:			3.00						
PAIR 3' COMPL:					1.00				
4 LEFT PRIMER	12	18	50.69	38.89	5.00	2.00	cttccagaacccgatgaaat		
RIGHT PRIMER	188	20	51.69	40.00	5.00	0.00	atgagagaacattgtagggg		
PRODUCT SIZE:	177								
PAIR ANY COMPL:			3.00						
PAIR 3' COMPL:					1.00				

Statistics

	con	too	in	in	no	tm	tm	high	high	high	.		
	sid	many	tar	excl	bad	GC	too	too	any	3'	poly	end	
	ered	NS	get	reg	GC:	clamp	low	high	compl	X	stab	ok	
Left	1069	0	0	0	0	0	18	983	4	27	22	2	13
Right	1206	0	0	0	0	0	17	1107	0	1	0	4	77

Pair Stats:
considered 303, unacceptable product size 149, high end compl 4, ok 150
primer3 release 0.9

(primer3_www_results.cgi v 0.2)

note the left primer from the above PRIMER3 output.

2.10.13 .3 Primer3 For Right Primer

The exon region which are taken to design the righth primer are 3746-3849.the parameters are similar to that of left primer.

Primer3

Primer3 is a program that takes as input a DNA sequence and returns a list of primers that are suitable for PCR amplification. The program is designed to be used as a web service, but it can also be run locally. The program is written in Perl and is available for Windows, Linux, and Mac OS. The program is distributed under the GNU General Public License. For more information, see the Primer3 website at http://www.genome.gov/cgi-bin/primer3_www.cgi.

Primer3 is a program that takes as input a DNA sequence and returns a list of primers that are suitable for PCR amplification. The program is designed to be used as a web service, but it can also be run locally. The program is written in Perl and is available for Windows, Linux, and Mac OS. The program is distributed under the GNU General Public License. For more information, see the Primer3 website at http://www.genome.gov/cgi-bin/primer3_www.cgi.

Primer3 is a program that takes as input a DNA sequence and returns a list of primers that are suitable for PCR amplification. The program is designed to be used as a web service, but it can also be run locally. The program is written in Perl and is available for Windows, Linux, and Mac OS. The program is distributed under the GNU General Public License. For more information, see the Primer3 website at http://www.genome.gov/cgi-bin/primer3_www.cgi.

General Primer Picking Conditions

Primer3 is a program that takes as input a DNA sequence and returns a list of primers that are suitable for PCR amplification. The program is designed to be used as a web service, but it can also be run locally. The program is written in Perl and is available for Windows, Linux, and Mac OS. The program is distributed under the GNU General Public License. For more information, see the Primer3 website at http://www.genome.gov/cgi-bin/primer3_www.cgi.

Other Seq-Sequence Inputs

Primer3 is a program that takes as input a DNA sequence and returns a list of primers that are suitable for PCR amplification. The program is designed to be used as a web service, but it can also be run locally. The program is written in Perl and is available for Windows, Linux, and Mac OS. The program is distributed under the GNU General Public License. For more information, see the Primer3 website at http://www.genome.gov/cgi-bin/primer3_www.cgi.

Sequence Quality

Sequence Quality

<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
--------------------------	--------------------------	--------------------------	--------------------------

Objective Function Penalty Weights for Primers

<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>

Objective Function Penalty Weights for Primer Pairs

<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>

Hyb Oligo (Internal Oligo) Per-Sequence Inputs

<input type="checkbox"/>

Hyb Oligo (Internal Oligo) General Conditions

<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

Objective Function Penalty Weights for Hyb Oligos (Internal Oligos)

<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/>	<input type="checkbox"/>

2.10.13 .5 Calculation Of The Product Size

Left primer sequence size is	223
Left primer starting point is	12
Left primer total length is	211
Introns length is	450
Right primer sequence size is	104
Right primer ending point is	101
Total right primer length is	101

Product size = Left primer total length is (211) +introns(450) +Right primer ending point is (101) = 762

2.10.14 Mitochondrial And Nuclear Enzyme

Mitochondrial and nuclear enzyme selected are Malate dehydrogenase and methyl transferase. The procedures of multiple alignments for these two enzymes are similar to that of chloroplast enzyme, which are explained above.

The summary sheet of the three enzymes is given in the table below

Sl. No.	Enzyme	Accession number	Cereal name	Exon number
1	malate dehydrogenase	M31965	sorghum vulgare	603-783,873-963, 1461-1547,1643-1727, 1984-2165,2420-2490, 2570-2660,2749-2823, 2912-3013,3100-3222, 3303-3371,3456-3524, 3619-3672,3765-4638
	malate dehydrogenase	AF353203	Oryza sativa	90-1088
	malate dehydrogenase	AF007581	zea mays	78-1076
2	methyl transferase	AF063403	zea mays	1854-3316,3481-4960, 5107-5271,5361-5553, 5634-5745,5819-5966, 6069-6266,6364-6642, 6723-6884,6960-7178, 7259-7417
	methyl transferase	AF063403	zea mays	5107-5271,5361-5553
	methyl transferase	AB028870	hordeum vulgare	83-3349
	methyl transferase	AF387790	sorghum bicolor	110-1198
	methyl transferase	U76384	Triticum aestivum	52-1173
	methyl transferase	AF042332	Oryza sativa	187-1227
3	Phosphoenolpyruvate carboxylase	E17154	maize	1587-1981,2092-2176,3074-3296,3746-3849,3975-4065,4203-4357,4454-5452,5550-5936
	Phosphoenolpyruvate carboxylase	AF399915	sorghum bicolor	21-844
	Phosphoenolpyruvate carboxylase	AF271995	oryza sativa	188-3065

2.11 Results

For our study we have taken Phosphoenol pyruvate carboxylase a chloroplast enzyme, Malate dehydrogenase a mitochondrial enzyme in sorghum bicolor, maize and oryza sativa, and methyl transferase a nuclear enzyme in zea mays, Hordeum vulgare, sorghum bicolor, Triticum aestivum and oryza sativa. In case of Phosphoenol pyruvate carboxylase, Malate dehydrogenase and methyl transferase, there is a maximum alignment is between oryza sativa and maize, maize and sorghum bicolor, zea mays and Hordeum vulgare. the maximum alignment can be studied by alignment graph. The distance between the two species can be studied by phylogenetic tree. The exons of the above mentioned enzymes are used for designing primers for PCR. The summary sheet of the primers for the above enzymes are given in the table 1.

Table1:

Accession number	Gene	Species	Exons region	Left primer	Left temp	Right primer	Left	Right
E17154	Phosphoenolpyruvate carboxylase	zea mays	4454-5452, 5550-5936	cgctcaagctcacctgtcc	60.83	tggaagttcctgacgtcct	59.7	631
E17154	Phosphoenolpyruvate carboxylase	zea mays	3074-3296, 3746-3849	cttcagaaccgatgaaat	50.69	aaacatcagctctcaatct	51.53	762
M31965	malatedehydrogenase	sorghumvulgare	2912-3013, 3100-3222	gttcctgattcttgaatgc	54.82	ccctgtggagaaccagtc	56.33	312
AF063403	methyl transferase	zea mays	5107-5271, 5361-5553	tgtatgtgcctgtggatga	59.94	tcgagccctctgataaacc	60.21	417
AF063403	methyl transferase	zea mays	5634-5745, 5819-5966	tatcttttacaaaatggcgattg	61.30	acacggaggaccaccattta	60.23	317

2.12 Discussion:

Phylogeny is about evolution and is used to reconstruct evolutionary events. It is now possible to construct phylogenetic evolution at a molecular level through analysis of molecular sequences, namely proteins & nucleic acids.

To construct phylogenetic tree among grass family, the sequences of conserved enzymes from mitochondria, chloroplast and nucleus are probed using bio-informatics tools. The scheme for such study is the following

- Identify exon regions for the enzyme to be investigated.
- An exon region of the particular enzyme is used to design the primers.
- Confirm the presence the particular sequence of the enzyme (exon) in the species of interest using wet lab techniques.
 - ✓ Isolation of chloroplast,mitochondrial and nuclear DNA
 - ✓ Amplification of DNA by using PCR
 - ✓ Hybridization techniques(southern blotting)
 - ✓ DNA sequencing by chemical and enzymatic methods.
 - ✓ Analysis of sequence based on mitochondrial and chloroplast to determine maternal inheritance.
 - ✓ Analysis based on nucleus to determine paternal inheritance.
 - ✓ Comparison of sequences using multiple alignment tools

Determine the relationship among the species under study

References

- Benson,G.(1999)Tandem repeats finder program to analyze DNA sequences *Nucleic acids Research* 27:573-580.
- GerardF.B.(2001).The use of the Monsanto draft rice genome sequence research. *Plant physiology* 125:1164-1165.
- Katti, V.M.;Sami-Subbu,R.;Ranjekar,P.K. and gupta. V.S.(2000). Amonoacid repeat patterns in protein sequences:Their diversity and structural and functional implications. *Protein science* 9:1203-1209.
- Lench,N.J.;Norris,A;Bailey,a.;Booth,A. and Markham,A.F.(1996).Vectoreete PCR isolation of microsatellite repeat sequences using anchored dinucleotide primers.*Nucleic acids research* 24:2190-2191.
- Macas,J.;Meszaros.T. and Nouzova,M.(2002).Plantsat:a specialized database for plant satellite repeats .*Bioinformatics* 18:28-35.
- Mahalakshmi,V. and Ortiz, R.(2001). Plant genomics and agriculture :from model organisms to crops, the role of data mining for gene discovery. *Electronic journal of biotechnology*.
- Dinakar bhatramakki,jjanmindong,ashok K.chhabra, and Gary E.Hart. An integrated SSR and RFLP linkage map of *Sorghum bicolor*(L.) Moench
- Pearson,C.E,nad Sinden,R,R.(1998). Trinucleotide repeat DNA structures:Dynamic mutations from dynamic DNA. *Current openion in Structural biology* 8:321-330
- Yuan. q.; Quackenbush,J.;Sultana ,R.;Pertea,m.;Salzberg,S.L and Buell, C.R.(20010). Rice bioinformatics. Analysis of rice sequence data and leveraging the data to other plant species. *Plant physiology*,125:1166-1174.

References

BASU,D.,K.DEHESH,H.J.SCHNEIDER-

POETSCH,S.E.HARRINGTON,S.R.McCOUCH and P.H.QUAL.2000.rice PHYC gene: structure,expression.map position and evolution. plant mol bio44:27-42.

DJE, Y., M. HEUERTZ, C. LEFEBVRE and X.VENKEMANS.2000.Assessment of genetic diversity within and among germplasm accessions in cultivated sorghum using microsatellite markers. Theoretical and applied genetics.100:918-925.

DOEBLEY, J. 1989. Isozymic evidence and the evolution of crop plants. Pp. 165-191 in D. E. Soltis and P. S. Soltis, eds. Isozymes in plant biology. Dioscorides press, Portland, Oregon.

DGGETT,H. 1976.Pp. 112-117.Evolution of crop plants. Longman, Essex, UK.

GAUT,B.S. and M. T. CLEGG. 1993a.Molecular evolution of the Adh1 locus in the genus Zea.proc Natl Acad Sci U S A 90: 5095-5099.

HARLAN, J. R. and A. B. L. STEMLER. 1976. The races of Sorghum in Africa. Pp. 465-478 in J. R. Harlan, J. M. de wet, and A. B. L. STEMLER, eds.origins of African plant Domestication: Mouton Press, The Hague.

HILTON,H. and B. S. GAUT. 1998.Speciation and domestication in maize and its wild relatives: evidence from the globulin-1 gene. Genetics 150:863-872.

HUDSON, R. R., M. KREITMAN and M.AGUADE.1987.Atest of neutral molecular evolution based on nucleotide data. Genetics 116: 153-159.

KIMURA, M. 1980. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences Mol Evol 16: 111-20.

KOLUKISA OGLU, H.U., S. MARX., C. WIEGMANN, S. HANELT and H. A.

SCHNEIDER-POETSCH.1995. Divergence of the phytochrome gene family predates angiosperm evolution and suggests that selaginella and Equisetum arose prior to psilotum. J Mol Evol 41: 329-337.

KUMAR, S., K. TAMURA, I. B. JAKOBSEN and M. NEI, 2001 MEGA2: Molecular Evolutionary Genetics Analysis software, Pp., Arizona state University, temp, Arizona, USA.

Nei,m.1987. molecular Evolutionary Genetics.Columbia University Press,New York.