

The Notion of Homology

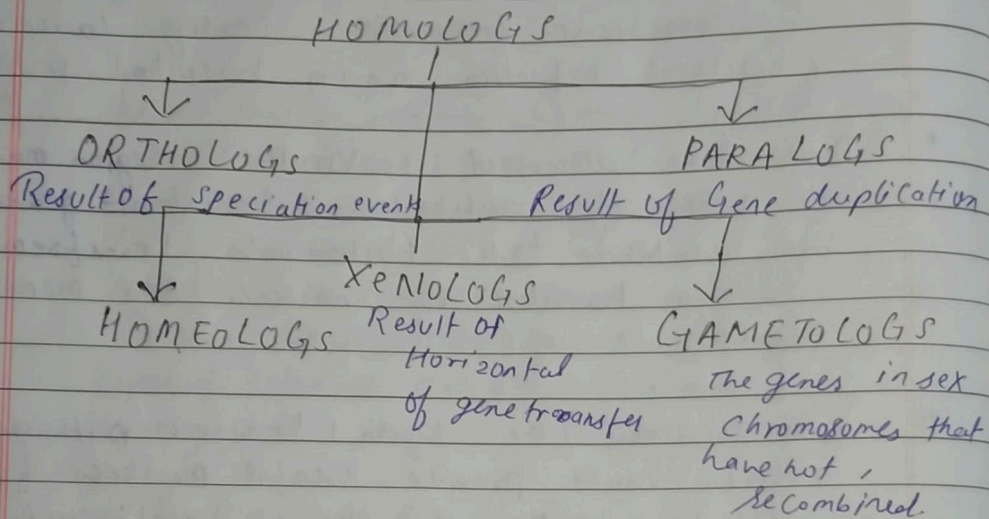
18



Date
Page

Introduction

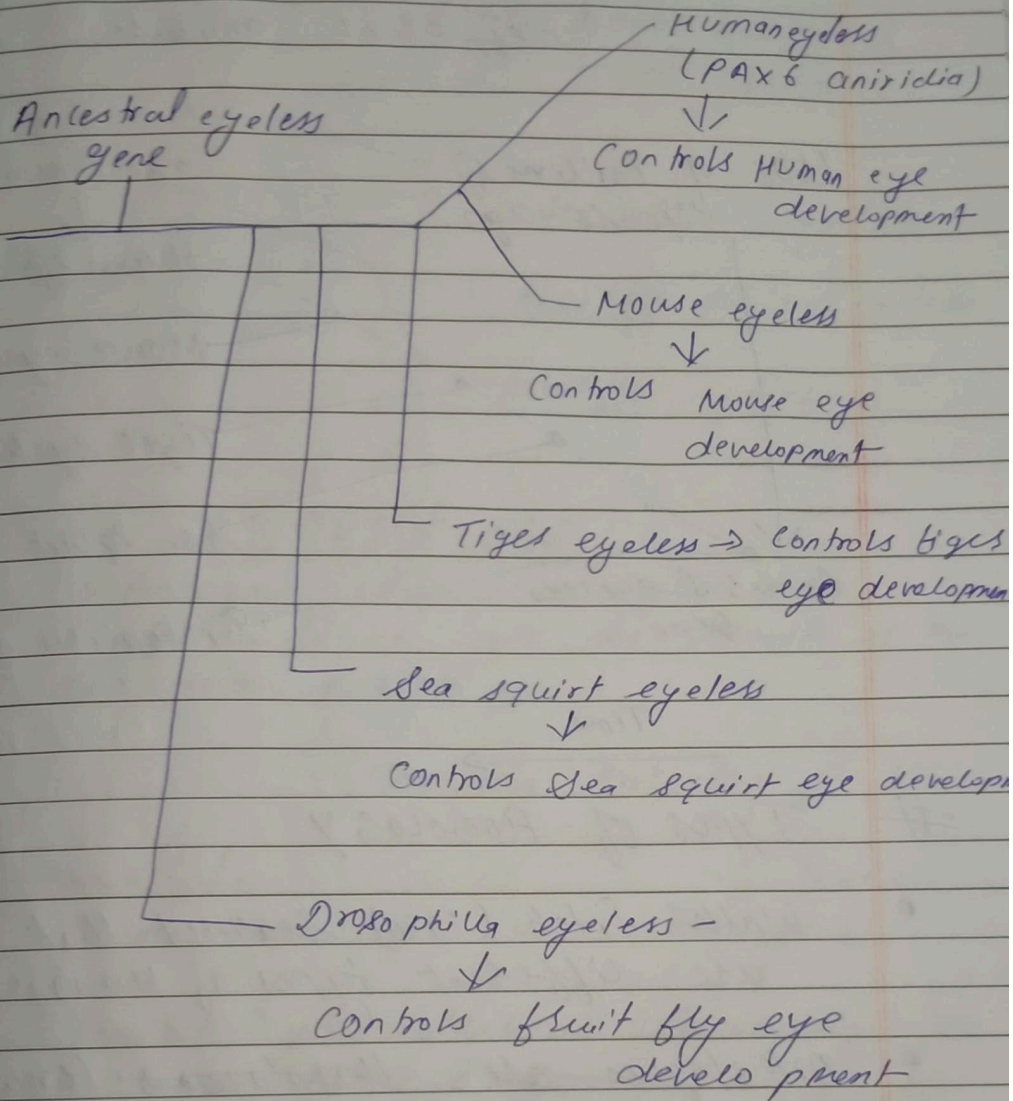
Homologs are two or more sequence that descend from a common ancestral sequence. Homologs are results of divergent evolution.



Slightly different versions of the eyeless gene control eye formation in many animals.

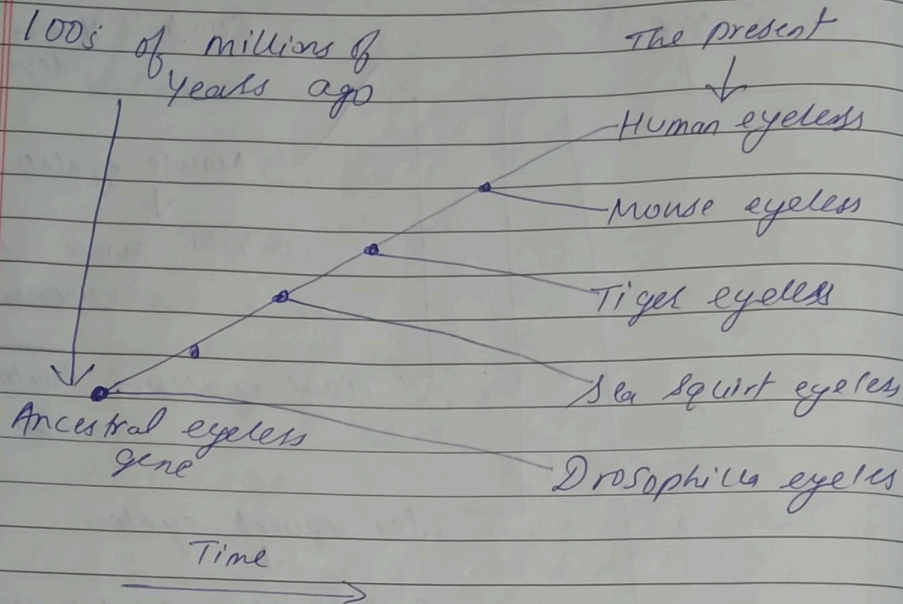
eyeless genes in different animals are Homologues.

i.e - they are Homologous (related) genes that descended from an ancestral gene in the ancestor of all these animals.



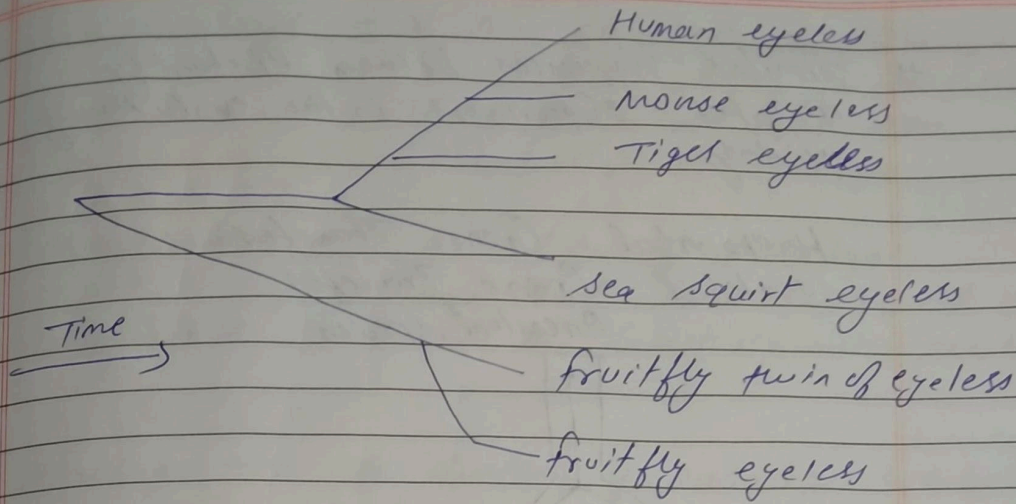


Aside: this is phylogenetic tree of eyeless genes in different animals.



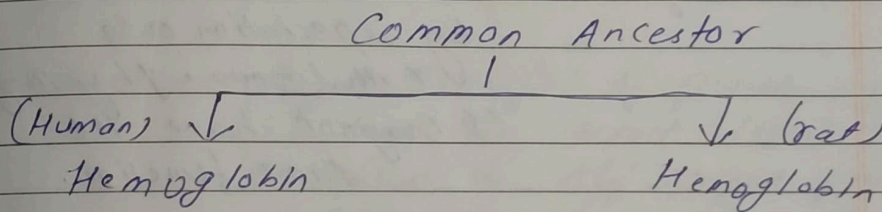
Types of Homology

- Walter Fitch (1970) realised that there are different types of Homologues.
- Orthologues are ~~homologous~~ Homologues in different species that arise due to the speciation event.



ORTHOLOGS

orthologous sequences are those sequences which are the results of the speciation process.



Similar or identical sequences in different species.

XENOLOGS

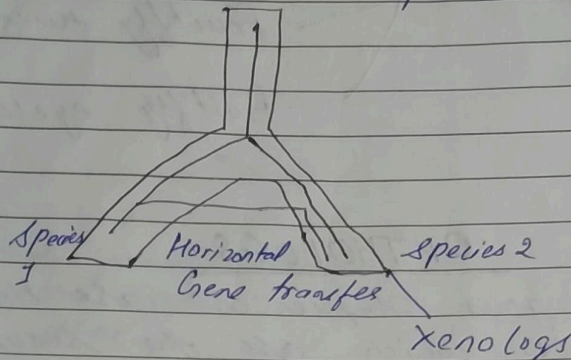
22



Date
Page

- # Similar sequences between distantly related organisms in the evolutionary history.

Horizontal Gene Transfer
Lateral Gene Transfer
Ancestral species



- #

Class	Description
Orthologs	→ Genes that diverged by speciation are orthologous. They may or may not have the same function.

- # Paralog → Genes that originate from an ancestral duplication and reside within the same species are paralogous.



- # In-paralogs → A subtype of paralogs. Genes that have duplicated after the last speciation event are in-paralogs.
- # Out-paralogs → A subtype of paralogs.
or
alloparalogs Genes that have duplicated before the latest speciation event.
- # Ohnologs → A subtype of paralogs. A whole genome duplication results in ohnologous genes. They have a common origin and have diverged an identical length of time.
- # Gametologs → The genes in sex chromosomes that have not recombined.
- # Xenologs → A transfer of genes, by means other than direct descent from parents to offspring from one organism to another, results in xeno-xenologous genes.
- # Homoeologs → A subtype of xenologs. The hybridization of genes from two separate species produces homoeologous species.

HOMEOLOGS

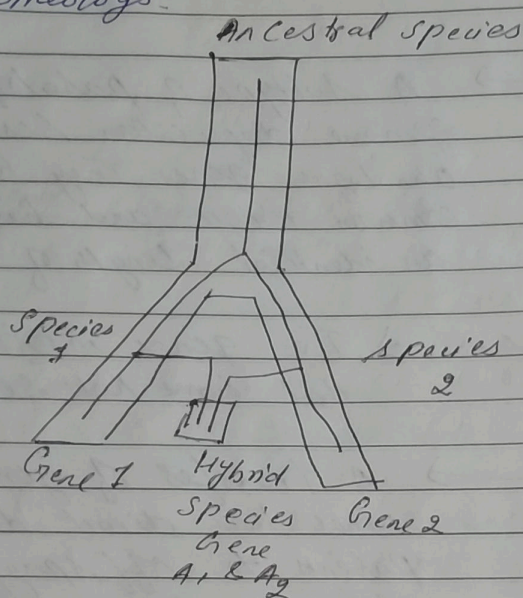
24



Date
Page

Homeologs are special cases of xenologs.

The genes which are separated by a speciation event when hybridized together via lateral gene transfer the gene sequences are known as Homeologs.



Homeologs



SUMMARY

- # Homologs :- Two or more sequence that descend from a common ancestral sequence.
- # Orthologs :- Sequence which are the result of the speciation process.
- # paralogous :- Sequence which are the result of the gene duplication process.
- # Xenologs :- Similar sequence b/w Distantly related organisms in the evolutionary history.
- # Analogs :- Sequence that show similar structure or function but don't share any common ancestral sequence.

#

SYMPARALOGS

Genes that have duplicated after the latest speciation event are known as In-Paralogs or Symparalogs.



Alloparalogs

Genes that have duplicated before the latest speciation event are known as out-paralogs or Alloparalogs.

OHNOLOGS

These are genes that are a result of a duplication of a whole genome.

ANALOGS

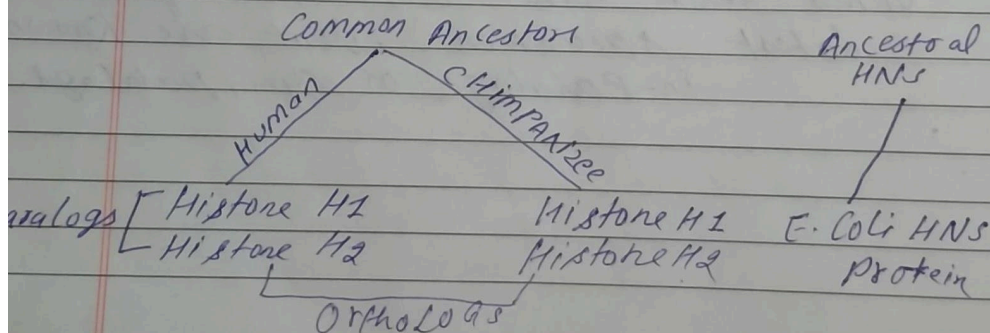
Analogous genes are genes that have identical or similar functions but don't share a common ancestor.

Convergent evolution

(i) Squid eye (ii) Human eye

Analogous have Homologous activity but heterologous origin

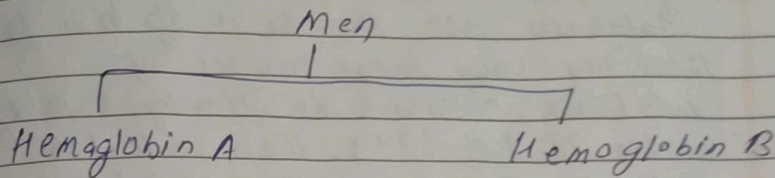
Homologs VS. ANALOGS





PARALOGUES

paralogous Sequence are Homologous
that have diverged within one
species i.e. arising during gene



Similar sequences within
SAME species.

The End

Sequence Information Sources

28

Date
Page

Databases :-

A database is a vast collection of data pertaining to a specific topic e.g. - nucleotide sequence, protein sequence etc. in an electronic environment.

Databases are the heart of bioinformatics. There is a very large number of databases, which is growing rapidly. At the end of 1999, there were 226 databases. In 2000, 55 new databases were created, raising the total to 281.

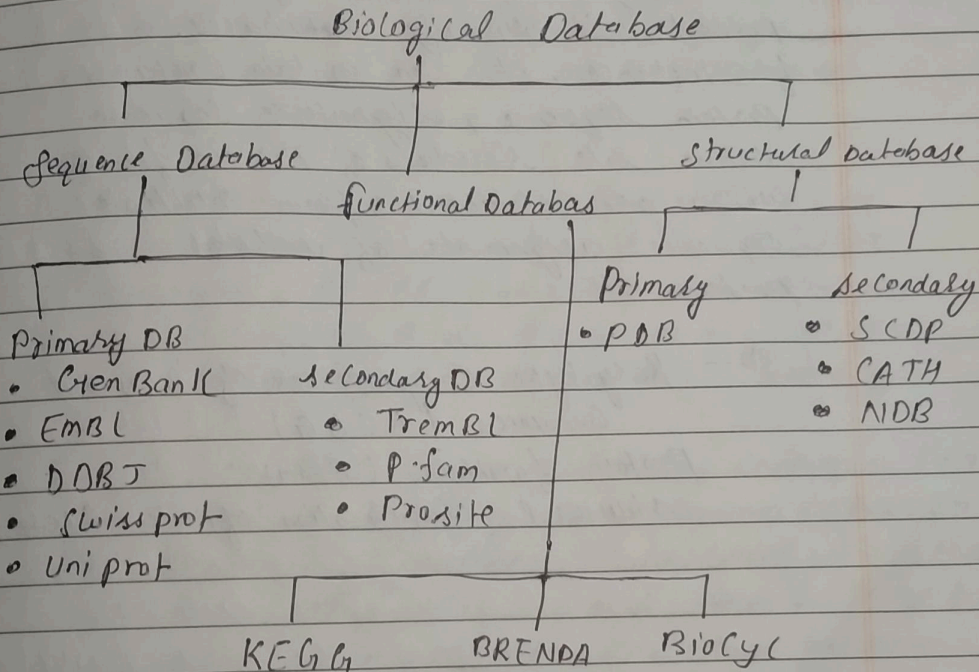
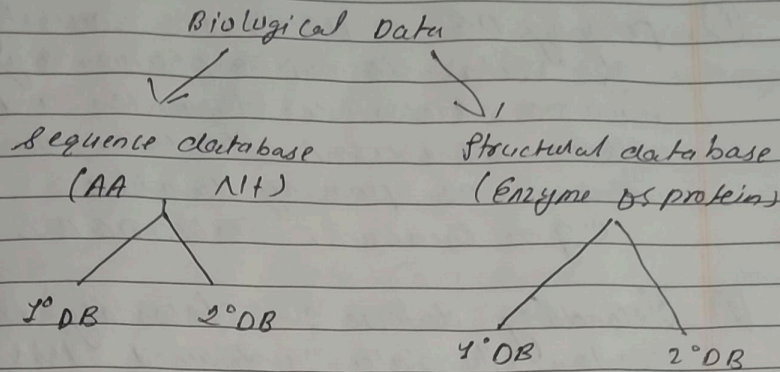
"Nucleotide Sequence databases :-

The major nucleotide sequence databases are Gene-bank held by NCBI, U.S.A DNA Database of Japan (DDBJ) and the Nucleotide Sequence database maintained by EMBL.

In addition to other nucleotide sequence databases have been created.



Classification of Databases.





Types of Databases

① Primary D.B. — Databases consisting of data derived experimentally such as nucleotide sequences and three dimensional structures are known as primary database.
e.g — Genbank, EMBL, DDBJ, NDB

② Secondary database — Secondary databases contain the data is obtained through the analysis or treatment of data present in primary databases.
for instance, it can contain conserved protein sequence, signature sequence active site residues of protein families which are obtained from multiple sequence alignment of related protein etc.

ex — Karyotype Encyclopedia of Genes and Genomics (KEGG)

Protein families (Pfam)

structural classification of protein (SCOP)



functional Database

Development of metabolic databases derived from the comparative study of metabolic pathways cater the industrial needs in more efficient manners to further the growth of systems biotechnology.

KEGG :— The Kyoto Encyclopedia of Genes and Genomics (KEGG) is the primary resource for the Japanese Genome Net Service that attempts to define the relationships b/n the functional meanings and utilities of the cell of the organisms and its genome information.

BRENDA :— It is main collection of enzyme functional data available to the scientific community. It is maintained and developed at Biochemistry and Bioinformatics at the technical University of Braunschweig, Germany.

Biocyc :— The Biocyc Database collection is a compilation of pathway and genome information for different organisms. It includes two other databases, ~~Eco~~ EcoCyc, which describe E. coli-like metaCyc, which describe pathways for more than 3000 organisms.



Sequence Database

Nucleotide and protein sequence databases represent the most widely established biological databases.

- # Serve as repositories for wet lab result and primary source for experimental result.
- # Genbank :— The Genbank nucleotide database is maintained by the NCBI which is part of the National Institute of Health (NIH), a federal agency of the US agency.
- # DDBJ :— DNA Data protein bank of Japan is a biological database that is run by the National Institute of Genetics, Japan.
- # PIR :— The protein information Resource an integrated public bioinformatics that support genomic and proteomic Research. and scientific. PIR has provided many protein databases and analysis tool including PSD of the functional annotated protein sequences.



- # Swiss Port:- Swiss port is a protein databank sequence and knowledge database. It is well known for its minimal redundancy, high quality of annotation, use of standard nomenclature. It contains amino acid sequence database, taxonomic group data and citation information.

- # TrEmbl:- This data consist of computer annotated entries derived from the translation of all coding sequence in the nucleotide databases.

Structure Database

Protein Data Bank:-

The PDB is the single world wide archive of structural data of Biological Macromolecules established in the Germany. National Lab in 1971.

- # It contains structural info. of the macromolecules determined by X-ray crystallographic, NMR methods. PDB is maintained by the Research Collaboratory for Structural Bioinformatics (RCSB).

Secondary:-

- # SCOP - structural classification of protein)
CATH - Class, Architecture Topology Homologous family
NDB - Nucleic Acid data base.



Entrez:-

Entrez is a molecular biology database system that provides integrated access to nucleotide and protein sequence, gene-centered and genomic-mapping information, 3D structure data, PubMed Medline.

The system is produced by NCBJ and available via Internet.

Entrez covers over 20 databases including the complete protein sequence data from PIR-International, PRT, Swiss, PROT, and PDB and nucleotide sequence data.

The Entrez Retrieval System uses an intuitive user interface for rapidly searching sequence and bibliographic data.



PROTEIN INFORMATION SOURCES

protein database can be sequence database or structural database.

Protein Sequence databases:-

The protein sequence database was developed at National Biomedical Research Foundation by Margaret Margreth Dayhoff in 1960's.

The protein sequence database was collaboratively maintained by PIR.

The Protein Information Resource (PIR) Database:-

It is main protein sequence database.

This database is classified into 4 classes:-

PIR1:- Classified and annotated entries.

PIR2:- Preliminary entries

PIR3 — Unverified entries

PIR4 — Conceptual translation of the sequence that are not transcribed, that are genetically engineered etc.

# SWISS PROT:-

It is a protein sequence database maintained collaboratively by medical biochemistry at the University of Geneva, Germany.

The d.b. endeavours to provide high level annotation description of the function of the protein and structure of the domains post translational modifications variants and so on.

They are interlinked to many source and have minimal redundancy.

TrEMBL:-

It was created in 1966 as a computer annotated supplement to Swiss Prot.

The d.b. contains translation of all coding sequences.

Two main sections:-

SP - TrEMBL - contain entries that not been annotated but they are eventually incorporated into Swiss Prot.



REM - Trembl :- Contains entries that are not included into Swiss prot
e.g - Ig seq - synthetic seq.

NRL - 3D :-

This d.B is produced by PIR from sequences extracted from PDB.
It is used both for similarity searches and keyword interrogation.

Structural Database

They store a collection of 3 dimensional biological macromolecular structures of proteins.

The last established database for protein structures is protein data bank (PDB).

PDB : It contains following information :-

Name of the protein

The species

Structural determination

Amino Acid sequence

Additional information.



Introduction to BLAST

An important goal of genomics and proteomics is to determine if a particular sequence is like another sequence. This is accomplished by comparing the new sequence with sequence that have already been reported and stored in database.

This process is principally one that using alignment procedures to uncover the "like" sequence in the database.

The alignment process will uncover those regions that are identical or close similar and those regions with little similarity.

Two alignment types are used
Global and local.

BLAST

BLAST - stands for Basic Local Alignment Search tool.

It is local alignment tool.

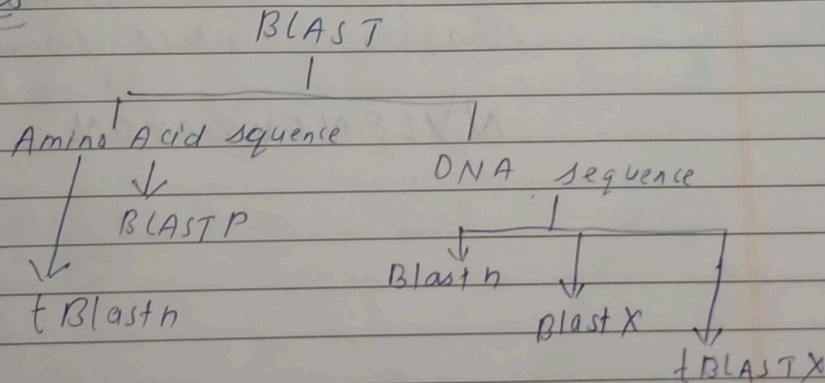
It helps to find regions of local similarity between sequences.

It is a program compares nucleotide or protein sequence to sequence databases and calculates the statistical significance of matches.

BLAST can be used to infer functional and evolutionary relationships b/w sequence as well as help identify members of gene families.

Google ^{search} → NCBI Homepage → Sequences

Types





Steps

Specifying A sequence of Interest



Selecting BLAST program



Selecting Database



Selecting optional parameter



Selecting formatting parameters

Process

The first step of the BLAST algorithm is to break the query into short words of specific length.

e.g - 12 amino acid near the terminal of the Arbidopsis thaliana protein phosphoglucosylase sequence are -

NYLENFVQATFN



This sequence is broken down into three character words by selecting the first Amino Acid characters.

N I Y L V L E L E N E N F N I F V F V Q
V Q A Q A T A T F T F N

These words are then compared against a sequence in a database.

for ex — word match with rabbit against a sequence in a database, muscle
Phosphoglucosylase

Query ENF

SSTNYAEMTIQSRNSTVEPAQR

Application of BLAST

BLAST can be used for several purposes these include.

- * Identifying species
- * Establishing phylogeny
- * DNA mapping
- * Locating domains.

the end



Multiple Sequence Alignment

- # A multiple sequence alignment (MSA) is a basic tool for the sequence alignment of two or more biological sequences.
- # Generally protein, DNA or RNA.
- # In many cases, the input set of query sequences are assumed to have an evolutionary relationship.
- # By which they share a lineage and are descended from a common ancestor.
- # Compare all sequence pairwise
- # Perform cluster analysis on the pairwise data to generate a hierarchy for alignment.
- # This may be in the form of a binary tree or a simple ordering.
- # Build the multiple alignment by first aligning the most similar pair of sequence.
- # Then the next most similar pair and so on



Once an alignment of two sequence 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 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.

EB — VTISCTGSSNAIG — NHUKWYALPL
VTISCTGSSING — ITVNWYALPH

Application of MSA: —

Detecting similarities b/w sequences (closely or distinctly related.)

Detecting conserved regions or motifs in sequences.

Detection of structural homologues.

Thus, assisting the improved prediction of secondary and tertiary structures of proteins.

Making pattern used to predict new pattern sequences falling in a given family.

Inferring evolutionary trees or linkages.



Progressive Alignment Method.

The most widely used approach to multiple sequence alignment.

Also known as the Hierarchical or Tree method.

Progressive Alignment Algorithms.

(i) Clustal W

(ii) T-Coffee

Iterative Refinement method —

A set of methods to produce MSA while reducing the errors inherent in progressive methods.

TREE Alignment —

In computational phylogenetic, it is used to analyse a set of sequences with evolutionary relationship using a fixed tree.

Essentially, Tree alignment is an algorithm for optimizing phylogenetic Tree.

STAR Alignment

Another form of Tree alignment.

Tree



FASTA format allows the sequence naming and comments to introduce the sequences. The use of FASTA format has become a standard for biologist to analyze the sequencing.

The format of FASTA codes is no longer than 120 characters.

Features of FASTA

Rather than trying to find out the best alignment between your data, it finds the patches of regional similarity.

It is rapid program - you can run the program locally or you can also send queries to an email server.

The alignments of FASTA can contain gaps. The sequence which contain the gap FASTA highlight these codes with red color.

Another features of FASTA it ignores the complete sensitivity and provide information about the expected matched alignments.



USES OF FASTA

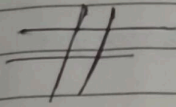
- # FASTA can be used for the Alignment of all types of proteins and DNA.
- # It can also use for the translation of algorithms which handle frame shift errors.
- # It can used for Calculating the similarity which can help biologists to decide whether the alignment is occurred by chance or it is due to infer homology.
- # You can also use FASTA for calculating the optimal score for alignment.
- # Help in identify the members of gene family.

DATA STRUCTURE

- # Data in FASTA is presented in a single code sequences. It has got a different search methods which help in searching the proteins.
Ex - Smith - Waterman type of algorithm
FASTA help you to find out the potential matches and save your time as well.



- # While the results of FASTA are reported in the form of histogram where the expected value are compared to random search. While the lower part of the histogram contain information about the matches of interest.



DATA SUBMISSION

NCBI, Bethesda Maryland.

Info

Houses series of databases relevant to biotechnology and biomedicine.

- # mainly genbank for DNA and pubmed, a bibliographic database for biomedical literature, epigenomics database.

- # Director: David Lipman

- # founded in 1988 by claud pepper.

NCBI

To submit a sequence in NCBI we need certain tools, which are easily found in NCBI page itself.

Database

Genbank

Sequence Read Archive

dbSNP (single nucleotide polymorphism)

dbVar (genome variant)

Geo (Gene expression omnibus)

BankIt

We have a single sequence, a simple set of sequences. (e.g. - 16S rRNA, mtDNA, ITS / rRNA) or a small batch of different sequences.

We prefer to use a web-based submission tool.

Sequin

- # We prefer to work on our submission off-line.
- # We have a sequence or sequence Polase complex
- # We would like graphical viewing and editing options, including an alignment editor.
- # Network access analysis tool.

Submission follow:

- # Contact information
- # Release date information
- # Reference information
- # Submission category and type
- # Nucleotide sequence
- # Organisms name
- # Feature of the sequence
- # Sequin (Registration)
 - (i) Author form
 - (ii) submission page
 - (iii) Contact page
 - (iv) Authors page
 - (v) Affiliation page
 - (vi) sequence format form
 - (vii) submission type
 - (viii) Seq. data format
 - (ix) submission category
 - (x) Organisms and sequence form
 - (xi) Nucleotide page.

due end



Genome Annotation pattern and Repeat Finding

Also called as DNA Annotation.

After DNA sequence DNA annotation is done. The DNA sequence will not make any sense without annotation in this:-

- # Location of genes are identified.
- # All coding Regions of the gene are identified.
- # Start and stop point of the gene are identified.
- # Function of the gene is identified.

Annotation Means:-

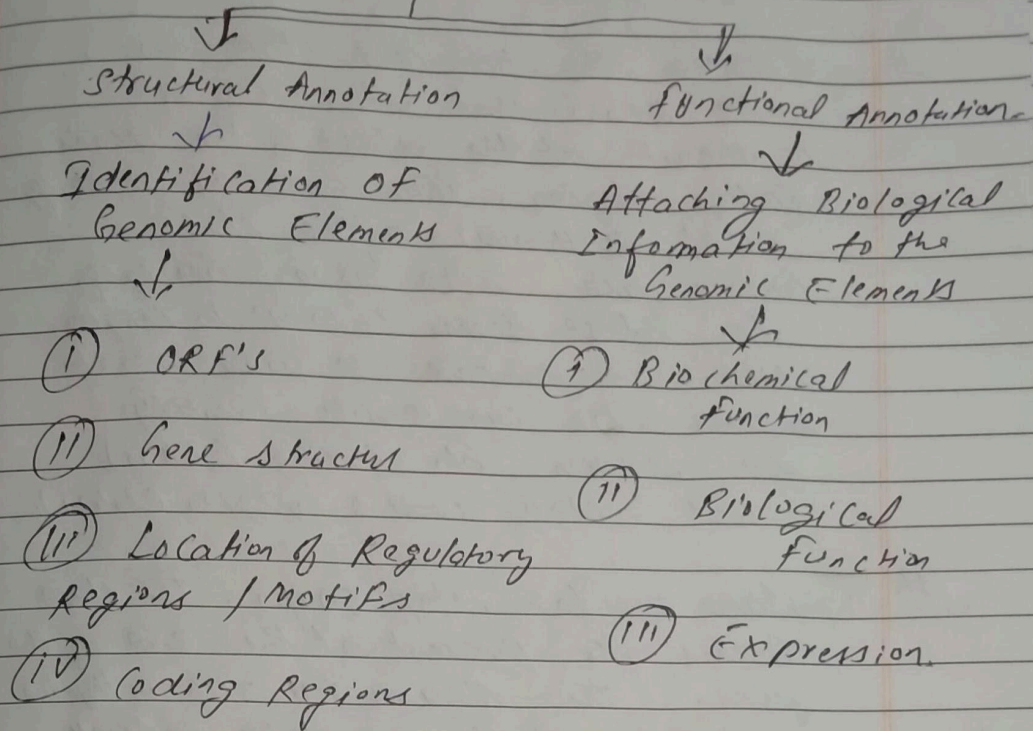
Extra information associated with any document is called as it annotation.

Includes comment or explanation.

Highlighting of key words, underlining certain lines in your answers are examples of annotation.



Types of Annotation



Gene Structural Annotation tool:—

Repeat masker is a program that screens DNA sequence for interspersed repeats and low complexity DNA sequences.

The output of the program is a detailed annotation of the repeats that are present in the query sequences as well as a annotated repeats have been masked.



Identification of Non-coding Regions:—

Identify the start and stop codon,
we know that start codon is present at the start of the gene's coding region. This codon is generally 'AUG' (Methionine). Other start codons are also known. Generally the start region is 20bp in length. When we compare various start regions, we can easily identify the region b/w coding (Exon) and non-coding region (Intron).

Same is for stop codon. The stop codons are present at the end of the gene. This region is rich in thymine and is about 20-30bp in length.

Gene Finding:—

(1) Done with the help of Gene Knockout experiments. The gene of interest is knocked out from the animal body. Then the animal is made to grow in the laboratory, phenotypic changes are observed.



BLAST, FASTA

Find out regions similarity between the given of the sequences. If a Region has very high similarity to mRNA or protein product, then it means that this region is a protein coding gene.

Functional Annotation:—

Similarity functional Annotation has been useful in various fields of life science includes:

Cancer cell profiling.

Study of complex disorders like Alzheimer disease, Schizophrenia etc.

Gynaecological studies like preterm birth.

Genome wide association studies

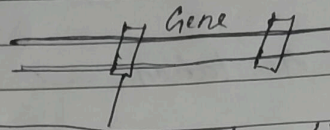
Brain related disorders like depression, Anxiety etc.

It is important in area of study genomics & proteomics.



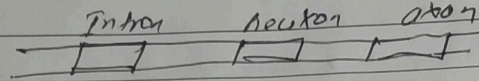
useful annotation

- (i) Automated annotation pipelines
- (ii) Human Genome project
- (iii) TIGR

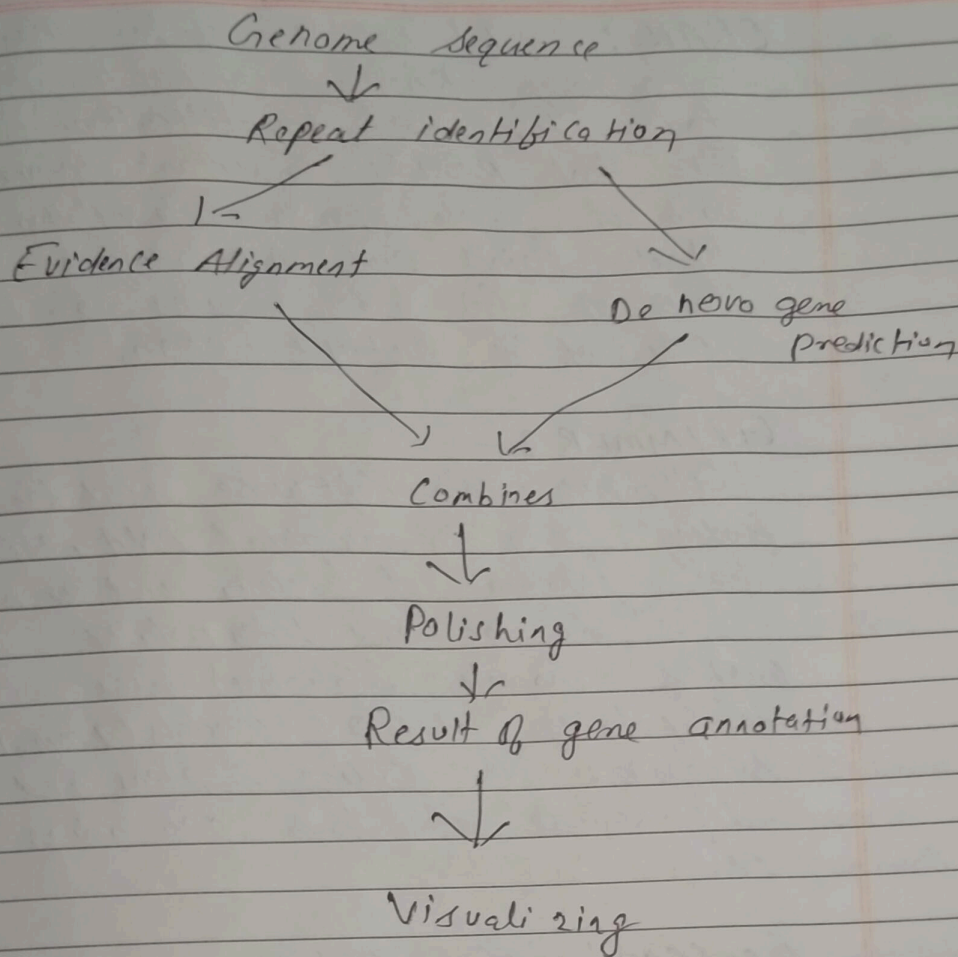


[GGC ~~GGA~~ ATG GCA AGT] ... [TTT TAG]

Open Reading frame



A G C T C G A T C G T A C G G C T A A T



Gene identification Tools

56

Date
Page

CRAIL:- It is one of the most commonly known computational tools for ORF identification. This tool provides important information such as splice junctions, translation start points and non-coding scores of 60 base regions on both sides of the putative exon.

GLIMMER:-

Glimmer is a software used for finding genes in microbial DNA, especially the genomes of bacteria and archaea. Gene locator and Interpolated Markov models (Glimmer) uses interpolated Markov models (IMMs) to recognize the coding regions and differentiate them from non-coding DNA.

Genscan:- This tool is used for the identification of complete gene structures in genomic DNA for various organisms. It can predict exon-intron structures of genes as well as locations in genomic sequences.

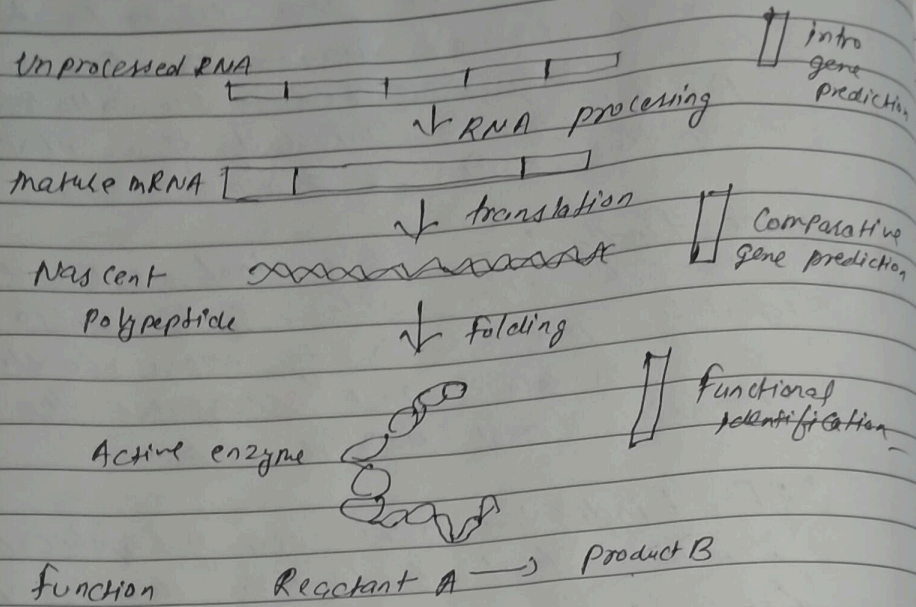


Genie :- This gene finder is based on generalized hidden Markov models. Genie was developed as a collaborative project by the University of California Computational Biology group.

Gene Finder :- This tool is used to predict splice sites. It can also identify protein-coding exons, construct gene models, and recognize the promoter and Poly-A region.

ORF Finder :- This is a graphical analysis tool that can detect open reading frames along with their protein translation from sequence already in the database. This program is used to search new DNA sequences for potential protein-encoding segments.

Easy Gene :- This tool is used to identify genes in prokaryotes, the current version of which includes 128 different organisms. Each gene identified by easy gene is attributed with a significance score (P-value), which reveals the probability of a sequence to be a non-coding open reading frame rather than a real gene.



Thank you

- Author - Amit Kumar Raaj
- Published - Biotechnologyme.blogspot.com
- Subject - Bioinformatics Notes