

SEMESTER V

17BTU512A

BIOINFORMATICS PRACTICAL

4H - 2C

Total hours/week: L:0 T:0 P:4

Marks: Internal: 40 External: 60 Total: 100

Practical

1. Sequence information resource
2. Understanding and use of various web resources: EMBL, Genbank, Entrez, Unigene, Protein information resource (PIR).
3. Understanding and using: PDB, Swissprot, TREMBL
4. Using various BLAST and interpretation of results.
5. Retrieval of information from nucleotide databases.
6. Sequence alignment using BLAST.
7. Multiple sequence alignment using Clustal W.

References

1. Ghosh, Z., & Bibekanand M. (2008). *Bioinformatics: Principles and Applications*. Oxford University Press.
2. Pevsner, J. (2009). *Bioinformatics and Functional Genomics* (2nd ed.). Wiley-Blackwell.
3. Campbell ,A. M., & Heyer, L.J. (2006). *Discovering Genomics, Proteomics and Bioinformatics* (2nd ed.). Benjamin Cummings.

KARPAGAM ACADEMY OF HIGHER EDUCATION, COIMBATORE
DEPARTMENT OF BIOTECHNOLOGY

LIST OF PRACTICAL

Name of the Faculty: Dr. S. Selvakumar	Department: Biotechnology
Course name: Bioinformatics Practical	Course code: 17BTU512A
Academic Year: 2019-2020	Semester: Fifth
Class: III year B.Sc.,	Section: A

1. Sequence information resource
2. Understanding and use of various web resources: EMBL, Genbank, Entrez, Unigene, Protein information resource (PIR).
3. Understanding and using: PDB, Swissprot, TREMBL
4. Using various BLAST and interpretation of results.
5. Retrieval of information from nucleotide databases.
6. Sequence alignment using BLAST.
7. Multiple sequence alignment using Clustal W.

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

Exp No: 1 Sequence information resources

AIM: To explore unique features and applications of various sequence resources.

Procedure: The National Center for Biotechnology Information (NCBI) is part of the United States National Library of Medicine (NLM), a branch of the National Institutes of Health (NIH). The NCBI houses a series of databases relevant to biotechnology and biomedicine and is an important resource for bioinformatics tools and services.

Gene database has been implemented at NCBI to characterize and organize the information about genes. It serves as a major node in the nexus of genomic map, expression, sequence, protein function, structure and homology data. A unique GeneID is assigned to each gene record that can be followed through revision cycles. Protein database is an important protein resource at NCBI. It maintains the text record for individual protein sequences, derived from many different resources such as NCBI Reference Sequence (RefSeq) project, GenBank, PDB and UniProtKB/SWISS-Prot. Protein records are present in different formats including FASTA and XML and are linked to other NCBI resources. Protein provides the relevant data to the users such as genes, DNA/RNA sequences, biological pathways, expression and variation data and literature.

- The amino acid sequence of proteins, nucleic acids and SNP would be downloaded in the FASTA format from the website <https://www.ncbi.nlm.nih.gov/>.
- Once the FASTA file is downloaded, it can be opened in a word pad file and stored in txt format.

UniProt is a freely accessible database of protein sequence and functional information, many entries being derived from genome sequencing projects. It contains a large amount of information about the biological function of proteins derived from the research literature. UniProt Archive (UniParc) is a comprehensive and non-redundant database, which contains all the protein sequences from the main, publicly available protein sequence databases.

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

There are many publicly available protein sequence databases.

FlyBase: the primary repository of genetic and molecular data for the insect family Drosophilidae (FlyBase)

H-Invitational Database (H-Inv)

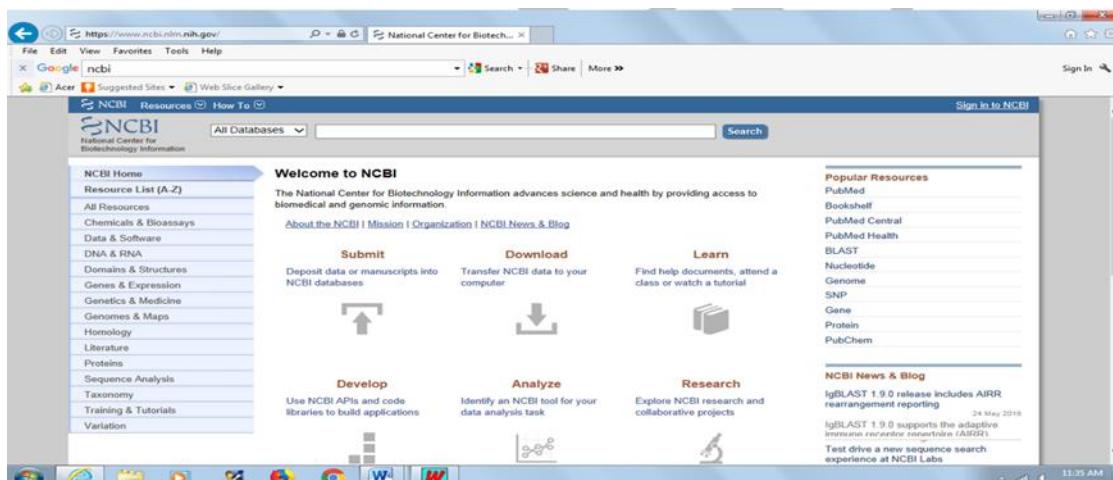
International Protein Index (IPI)

Saccharomyces Genome Database (SGD)

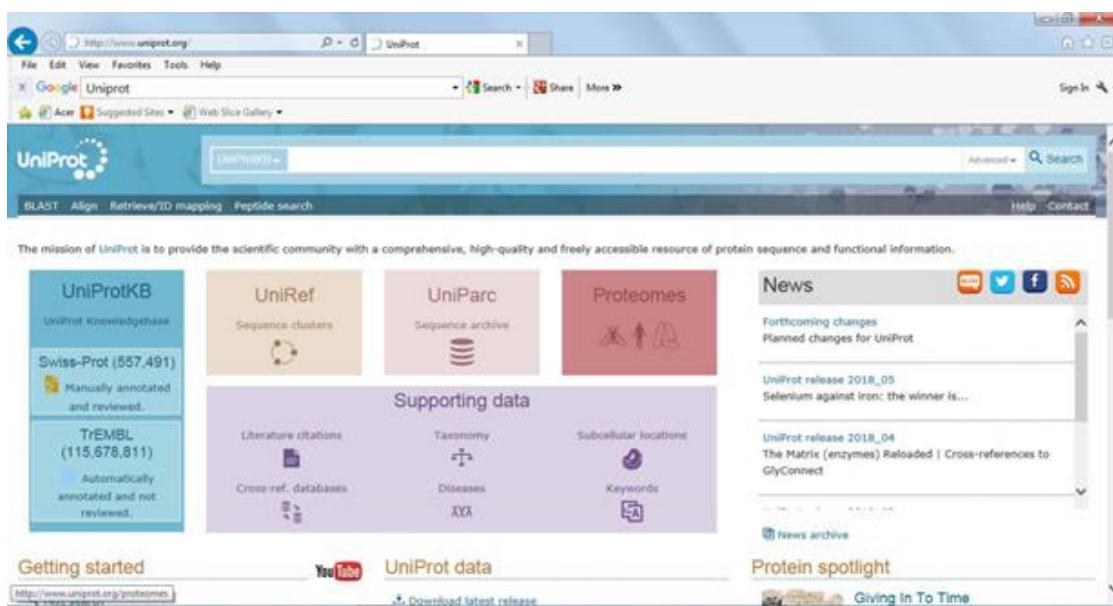
The Arabidopsis Information Resource (TAIR)

Vertebrate and Genome Annotation Database (VEGA)

WormBase



**CLASS : III B. Sc., BT
 COURSE NAME : BIOINFORMATICS PRACTICAL
 BATCH : 2017 – 2020
 COURSE CODE : 17BTU512A**



Information pertaining to primary, secondary and tertiary structures of proteins/nucleic acids can be inferred from the annotation webpages of the corresponding sequences. For instance, annotation data of myoglobin are herein listed out for the quick understanding.

Protein : Myoglobin

Gene : MB

Organism : *Homo sapiens* (Human)

Status : Reviewed

Annotation score : Function

Serves as a reserve supply of oxygen and facilitates the movement of oxygen within muscles.

Sites

Feature key	Position(s)	DescriptionActions	Graphical view	Length
-------------	-------------	--------------------	----------------	--------

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Metal binding ⁱ	65	Iron (heme distal ligand)		1
Metal binding ⁱ	94	Iron (heme proximal ligand)		1

GO - Molecular function

[View the complete GO annotation on QuickGO](#)

GO - Biological process

Keywords

Molecular function	Muscle protein
Biological process	Oxygen transport, Transport
Ligand	Heme, Iron, Metal-binding

Enzyme and pathway databases

Reactome ⁱ	R-HSA-8981607 Intracellular oxygen transport
-----------------------	--

Protein family/group databases

TCDB ⁱ	1.A.107.1.3 the pore-forming globin (globin) family
-------------------	---

Names & Taxonomy

Protein names	Recommended name: Myoglobin
Gene names	Name: MB
Organism	Homo sapiens (Human)
Taxonomic identifier	9606 [NCBI]
Taxonomic lineage	Eukaryota › Metazoa › Chordata › Craniata › Vertebrata › Euteleostomi › Mammalia › Eutheria › Euarchontoglires › Primates › Haplorrhini › Catarrhini › Hominidae › Homo
Proteomes	UP000005640 Component: Chromosome 22

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Organism-specific databases

EuPathDB	HostDB:ENSG00000198125.12
HGNC	HGNC:6915 MB
MIM	160000 gene
neXtProt	NX_P02144

Subcellular locationⁱ

Extracellular region or secreted Cytosol Plasma membrane Cytoskeleton Lysosome Endosome Peroxisome ER Golgi apparatus Nucleus Mitochondrion Manual annotation Automatic computational assertionGraphics by Christian Stolte; Source: COMPARTMENTS

- GO - Cellular component
 - Cytosol
 - Cytosol Source: Reactome
 - Extracellular region or secreted
 - Extracellular exosome Source: UniProtKB

View the complete GO annotation on QuickGO

Pathology & Biotech

Organism-specific databases

DisGeNET	4151
OpenTargets	ENSG00000198125
PharmGKB	PA30658

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Chemistry databases

ChEMBL	CHEMBL2406892
DrugBank	DB02671 1-Methylimidazole DB03385 4-Methylimidazole DB02379 Beta-D-Glucose DB02073 Biliverdine Ix Alpha DB03399 Ethyl Isocyanide DB04337 Methyl Isocyanide DB02396 Methylethylamine DB01826 N-Butyl Isocyanide DB04050 N-Propyl Isocyanide DB02646 Nitrosoethane DB01710 Porphyrin Fe(III) DB02528 Tetrazolyl Histidine

Polymorphism and mutation databases

BioMuta	MB
DMDM	127661

Results:

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

Exp No: 2. Understanding and use of various web resources: EMBL, GenBank, Entrez, Unigene and Protein information resources

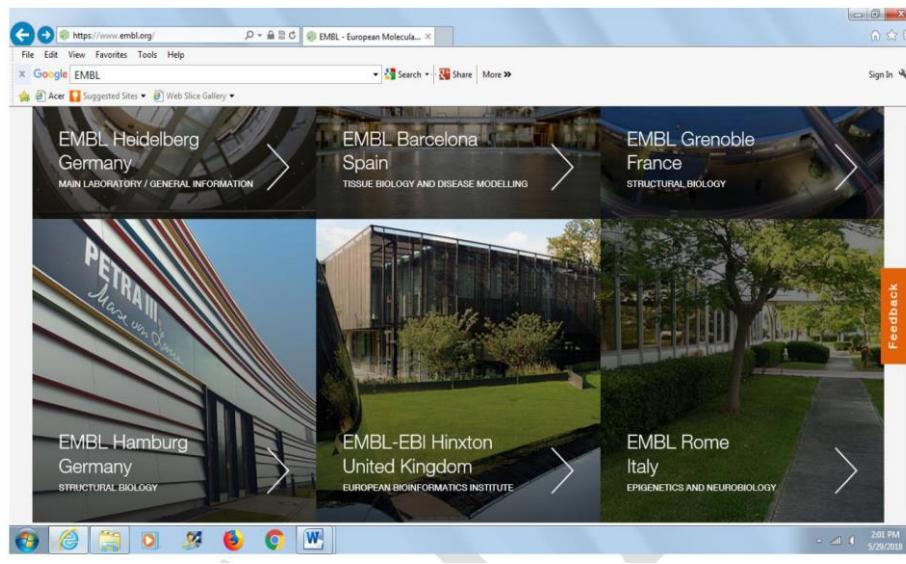
Aim: To understand and exploit the biological data available in various web resources.

Procedure: The European Molecular Biology Laboratory (EMBL) is a molecular biology research institution supported by 22 member states, four prospect and two associate member states. EMBL was created in 1974 and is an intergovernmental organisation funded by public research money from its member states. Research at EMBL is conducted by approximately 85 independent groups covering the spectrum of molecular biology. The Laboratory operates from five sites: the main laboratory in Heidelberg, and outstations in Hinxton (the European Bioinformatics Institute (EBI), in England), Grenoble (France), Hamburg (Germany), Monterotondo (near Rome) and Barcelona (Spain). EMBL groups and laboratories perform basic research in molecular biology and molecular medicine as well as training for scientists, students and visitors. The organization aids in the development of services, new instruments and methods, and technology in its member states. Israel is the only Asian state that has full membership. In March 2010, the EMBL Advanced Training Centre (ATC) was inaugurated on the main campus in Heidelberg. Shaped in the form of a double helix, it hosts conferences and provides training. EMBL also runs an active Science and Society Programme which offers activities and events on current questions in life science research for the general public and the scientific community.

GenBank is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences. GenBank is part of the International Nucleotide Sequence Database Collaboration, which comprises the DNA DataBank of Japan (DDBJ), the European Nucleotide Archive (ENA), and GenBank at NCBI. These three organizations exchange data on a daily basis.

KARPAGAM ACADEMY OF HIGHER EDUCATION

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**



Access to GenBank: There are several ways to search and retrieve data from GenBank.

Search GenBank for sequence identifiers and annotations with Entrez Nucleotide, which is divided into three divisions: CoreNucleotide (the main collection), dbEST (Expressed Sequence Tags), and dbGSS (Genome Survey Sequences).

Search and align GenBank sequences to a query sequence using BLAST (Basic Local Alignment Search Tool). BLAST searches CoreNucleotide, dbEST, and dbGSS independently; see BLAST info for more information about the numerous BLAST databases.

Search, link, and download sequences programmatically using NCBI e-utilities.

The ASN.1 and flatfile formats are available at NCBI's anonymous FTP server: <ftp://ftp.ncbi.nlm.nih.gov/ncbi-asn1> and <ftp://ftp.ncbi.nlm.nih.gov/genbank>.

Entrez is a molecular biology database system that provides integrated access to nucleotide and protein sequence data, gene-centered and genomic mapping information, 3D structure data, PubMed MEDLINE, and more. The system is produced by the National Center for Biotechnology Information (NCBI) and is available via the Internet

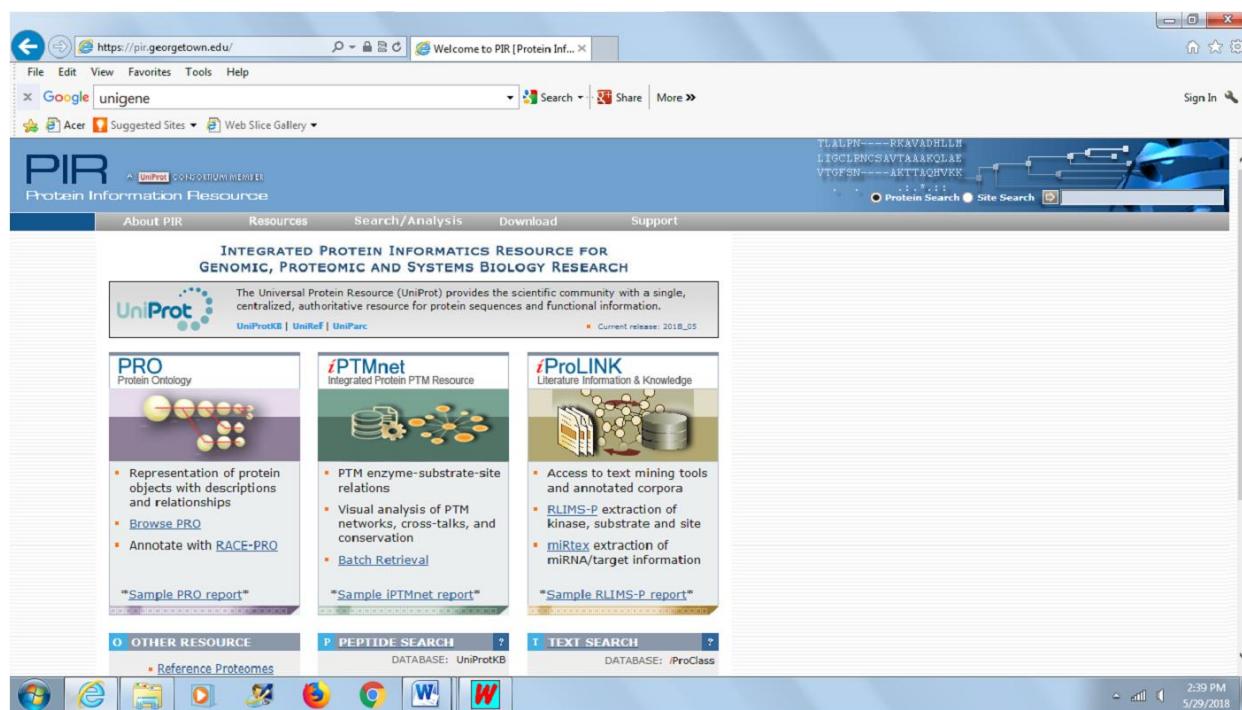
UniGene is an NCBI database of the transcriptome and thus, despite the name, not primarily a

Prepared by: Dr. Selvakumar, S. Assistant Professor, Department of Biotechnology, KAHE.

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

database for genes. Each entry is a set of transcripts that appear to stem from the same transcription locus (i.e. gene or expressed pseudogene). Information on protein similarities, gene expression, cDNA clones, and genomic location is included with each entry.

The Protein Information Resource (PIR) is an integrated public bioinformatics resource to support genomic, proteomic and systems biology research and scientific studies. The PIR continues to offer world leading resources to assist with proteomic and genomic data integration and the propagation and standardization of protein annotation.



Result:

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

Exp No: 3. Understanding and use of PDB, Swissprot and TrEMBL

Aim: To analyze various sequence and structural information provided in the PDB, Swiss-Prot and TREMBL databases.

Procedure: SWISS-PROT is a curated protein sequence database which strives to provide a high level of annotation (such as the description of the function of a protein, its domains structure, post-translational modifications, variants, etc.), a minimal level of redundancy and high level of integration with other databases. Recent developments of the database include format and content enhancements, cross-references to additional databases, new documentation files and improvements to TrEMBL, a computer-annotated supplement to SWISS-PROT. TrEMBL consists of entries in SWISS-PROT-like format derived from the translation of all coding sequences (CDSs) in the EMBL Nucleotide Sequence Database, except the CDSs already included in SWISS-PROT. SWISS-PROT is available at: <http://www.expasy.ch/sprot/> and <http://www.ebi.ac.uk/swissprot/>

The Protein Data Bank (PDB) is a crystallographic database for the three-dimensional structural data of large biological molecules, such as proteins and nucleic acids. The data, typically obtained by X-ray crystallography, NMR spectroscopy, or, increasingly, cryo-electron microscopy, and submitted by biologists and biochemists from around the world, are freely accessible on the Internet via the websites of its member organisations (PDBe, PDBj, and RCSB). The PDB is overseen by an organization called the Worldwide Protein Data Bank, wwPDB. The PDB is a key resource in areas of structural biology, such as structural genomics.

KARPAGAM ACADEMY OF HIGHER EDUCATION

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

The image contains two screenshots of protein databases. The top screenshot shows the SWISS-PROT database homepage, which features links for SWISS-PROT Home, Information, Access, Submissions, Tools, FTP, Group info, Documents, and Contact. It also highlights the TREMBL database and provides access to the SWISS-PROT Database. The bottom screenshot shows the RCSB PDB homepage, featuring a search bar, a map of the world, and sections for A Structural View of Biology, RCSB PDB Services and Impact, and May Molecule of the Month (a 3D model of a protein structure).

For instance, the following is the structural information provided in the PDB file of a cardiotoxin from *Naja naja atra*. The PDB ID of the protein is 2CRT.

HEADER CARDIOTOXIN 12-MAR-94 2CRT

TITLE CARDIOTOXIN III FROM TAIWAN COBRA (NAJA NAJA ATRA

TITLE 2 DETERMINATION OF STRUCTURE IN SOLUTION AND COMPARISON WITH

TITLE 3 SHORT NEUROTOXINS

Prepared by: Dr. Selvakumar, S. Assistant Professor, Department of Biotechnology, KAHE.

Page 11

KARPAGAM ACADEMY OF HIGHER EDUCATION

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

COMPND MOL_ID: 1;

COMPND 2 MOLECULE: CARDIOTOXIN III;

COMPND 3 CHAIN: A;

COMPND 4 ENGINEERED: YES

SOURCE MOL_ID: 1;

SOURCE 2 ORGANISM_SCIENTIFIC: NAJA ATRA;

SOURCE 3 ORGANISM_COMMON: CHINESE COBRA;

SOURCE 4 ORGANISM_TAXID: 8656

KEYWDS CARDIOTOXIN

EXPDTA SOLUTION NMR

AUTHOR R.BHASKARAN,C.C.HUANG,K.D.CHANG,C.YU

REVDAT 3 24-FEB-09 2CRT 1 VERSN

REVDAT 2 01-APR-03 2CRT 1 JRNL

REVDAT 1 01-NOV-94 2CRT 0

JRNL AUTH R.BHASKARAN,C.C.HUANG,D.K.CHANG,C.YU

JRNL TITL CARDIOTOXIN III FROM THE TAIWAN COBRA (NAJA NAJA

JRNL TITL 2 ATRA). DETERMINATION OF STRUCTURE IN SOLUTION AND

JRNL TITL 3 COMPARISON WITH SHORT NEUROTOXINS.

JRNL REF J.MOL.BIOL. V. 235 1291 1994

JRNL REFN ISSN 0022-2836

KARPAGAM ACADEMY OF HIGHER EDUCATION

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

JRNL PMID 8308891

JRNL DOI 10.1006/JMBI.1994.1082

REMARK 1

REMARK 2

REMARK 2 RESOLUTION. NOT APPLICABLE.

REMARK 3

REMARK 3 REFINEMENT.

REMARK 3 PROGRAM : X-PLOR

REMARK 3 AUTHORS : BRUNGER

REMARK 3

REMARK 3 OTHER REFINEMENT REMARKS: NULL

REMARK 4

REMARK 4 2CRT COMPLIES WITH FORMAT V. 3.15, 01-DEC-08

REMARK 100

REMARK 100 THIS ENTRY HAS BEEN PROCESSED BY BNL.

REMARK 210

REMARK 210 EXPERIMENTAL DETAILS

REMARK 210 EXPERIMENT TYPE : NMR

REMARK 210 TEMPERATURE (KELVIN) : NULL

The following is atomic coordinates of the protein. It provides information pertaining to atom

KARPAGAM ACADEMY OF HIGHER EDUCATION

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

types, atom nomenclature, atom number, atom position and atom occupancy.

ATOM	1	N	LEU	A	1	12.825	6.867	0.006	1.00	1.76
N										
ATOM	2	CA	LEU	A	1	12.032	6.006	0.919	1.00	0.90
C										
ATOM	3	C	LEU	A	1	11.849	4.631	0.268	1.00	0.73
C										
ATOM	4	O	LEU	A	1	12.036	4.471	-0.924	1.00	0.90
O										
ATOM	5	CB	LEU	A	1	10.658	6.639	1.165	1.00	1.16
C										
ATOM	6	CG	LEU	A	1	10.813	8.094	1.633	1.00	1.61
C										
ATOM	7	CD1	LEU	A	1	9.429	8.662	1.957	1.00	2.23
C										
ATOM	8	CD2	LEU	A	1	11.676	8.150	2.900	1.00	2.57
C										
ATOM	9	H1	LEU	A	1	13.292	6.273	-0.707	1.00	2.15
H										
ATOM	10	H2	LEU	A	1	12.194	7.547	-0.470	1.00	2.36
H										
ATOM	11	H3	LEU	A	1	13.547	7.380	0.551	1.00	2.21
H										

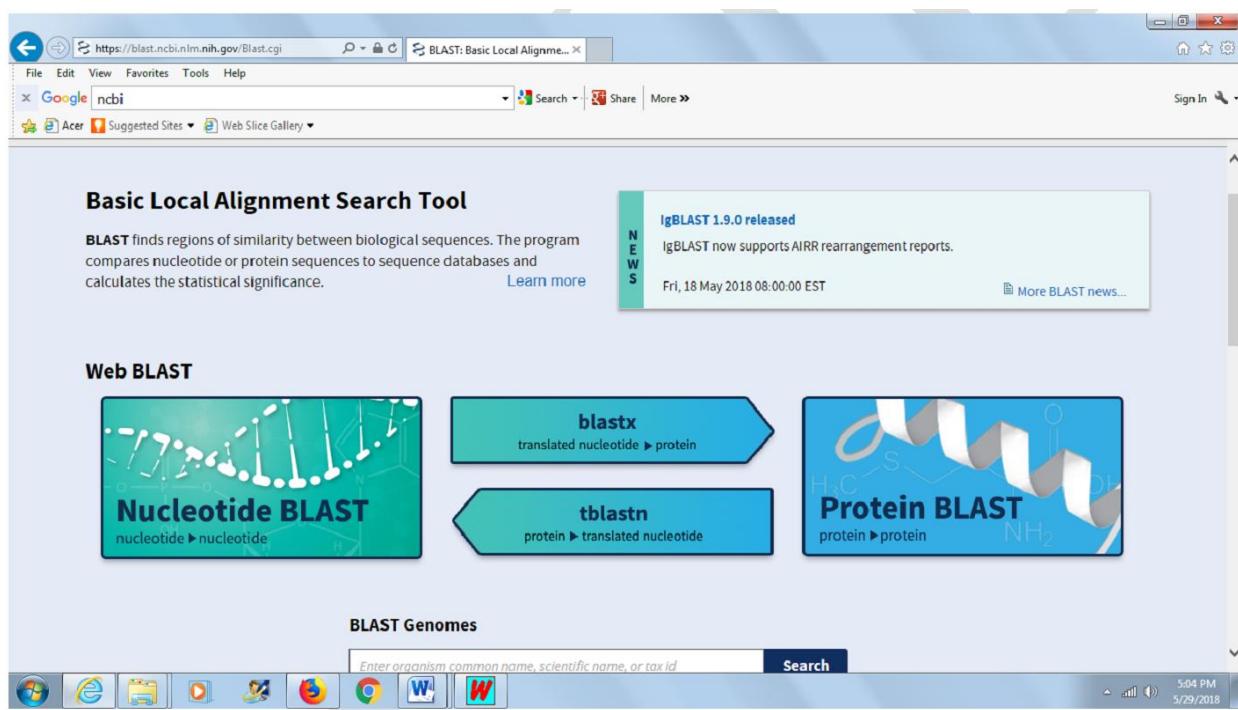
Results:

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

Exp No: 4. Using various BLAST and interpretation of results

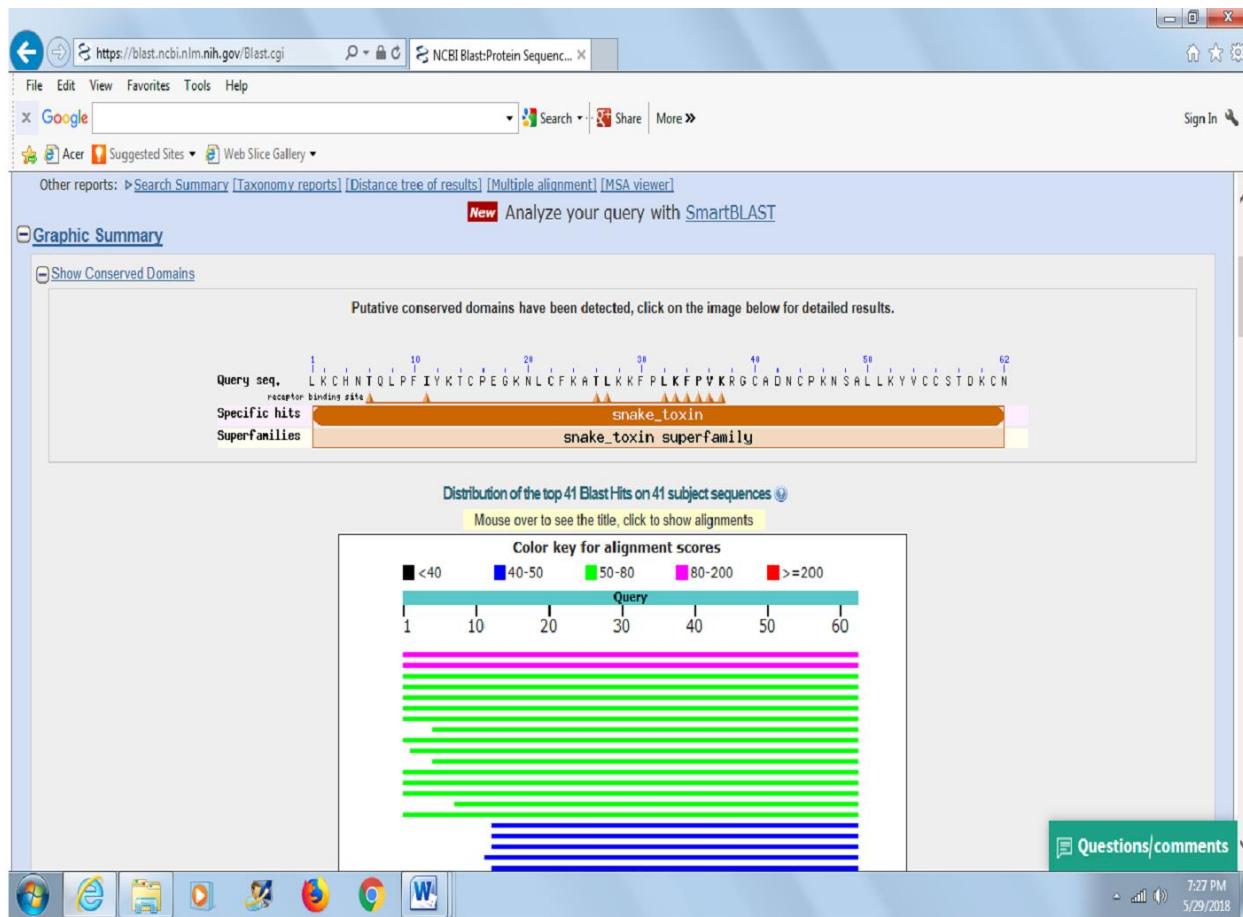
Aim: To retrieve similar sequences for a given protein/nucleotide primary structure using BLAST and rationalizing the outcomes

Procedure: Retrieve a few numbers of protein/nucleotide sequences from NCBI/UniProt databases (<https://www.ncbi.nlm.nih.gov/>) and save the retrieved sequences in FASTA format and accession IDs for all the sequences as well. The ‘sequence similarity search’ can be executed using various BLAST algorithms.



Using protein blast, similar sequences can be obtained either from primary or from structural databases. In general, for distant related and close related sequences, PSI-BLAST, BLASTP algorithms are preferred, respectively. For instance, a cardiotoxin (1CVO) was subjected to BLASTP similarity search and the outcomes were as shown herein below.

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**



The sequences would be further analyzed on the basis of query coverage, total score, percentage of identities and E-values. The score values and other statistical parameters for a few hits of the query sequences are depicted in the following illustration.

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

The screenshot shows a Microsoft Internet Explorer browser window displaying the NCBI Blast results for protein sequences. The URL in the address bar is <https://blast.ncbi.nlm.nih.gov/Blast.cgi>. The search term is "NCBI Blast:Protein Sequenc...". The results table has columns for Description, Max score, Total score, Query cover, E value, Ident, and Accession. The results include entries for various cardiotoxins and their structures.

Description	Max score	Total score	Query cover	E value	Ident	Accession
Chain A_ Structure Of Cytotoxin Homolog Precursor	124	124	100%	1e-39	100%	1QXI_A
Chain A_ Crystal Structure Of A Three Finger Toxin From Snake Venom	86.3	86.3	100%	2e-24	68%	3VTS_A
Chain A_ CARDIOTOXIN IV/I FROM NAJA MOSSAMBICA MOSSAMBICA: THE REFINED CRYSTAL STRUCTURE	78.2	78.2	100%	3e-21	66%	1CDT_A
Chain A_ NMR structure with tightly bound water molecule of cytotoxin I from Naja oxiana in aqueous solution (major form)	76.3	76.3	100%	2e-20	65%	1RL5_A
Chain A_ Minor form of the recombinant cytotoxin-I from N. oxiana	76.3	76.3	100%	2e-20	65%	5LUE_A
Chain A_ X-RAY STRUCTURE AT 1.55 Å OF TOXIN GAMMA, A CARDIOTOXIN FROM NAJA NIGRICOLLIS VENOM. CRYSTAL PACKING REVEALS A	75.5	75.5	100%	4e-20	61%	1TGX_A
Chain A_ CARDIOTOXIN II FROM TAIWAN COBRA VENOM. NAJA NAJA ATRA: STRUCTURE IN SOLUTION AND COMPARISON AMONG HOMOLOG	75.1	75.1	100%	6e-20	65%	1CRE_A
Chain A_ SOLUTION STRUCTURE OF CARDIOTOXIN IV_NMR_1-STRUCTURE	74.3	74.3	93%	1e-19	64%	1KBS_A
Chain A_ Nmr Structure Of Ctx A3 At Neutral Ph (20 Structures)	73.6	73.6	100%	2e-19	63%	1D2_A
Chain A_ Crystal Structure Of Beta-cardiotoxin, A Novel Three-finger Cardiotoxin From The Venom Of Ophiodophagus Hannah	73.6	73.6	98%	3e-19	57%	3PLC_A
Chain A_ DETERMINATION OF THE NUCLEAR MAGNETIC RESONANCE SOLUTION STRUCTURE OF CARDIOTOXIN CTX IIB FROM NAJA MOSSA	72.8	72.8	93%	5e-19	59%	2CX_A
Chain A_ STRUCTURE OF COBRA CARDIOTOXIN CTX I AS DERIVED FROM NUCLEAR MAGNETIC RESONANCE SPECTROSCOPY AND DISTANCE	71.2	71.2	100%	2e-18	61%	2CDX_A
Chain S_ Elucidation Of The Solution Structure Of Cardiotoxin Analogue V From The Taiwan Cobra (Naja Naja Atra) Venom	70.1	70.1	100%	5e-18	60%	1CHV_S
Chain A_ Crystal Structure Of Cardiotoxin VI From Taiwan Cobra (Naja Atra) Venom	70.1	70.1	100%	6e-18		

Similar types of data analyzes can be carried out for nucleotide sequences and as well for other BLAST algorithms.

Results:

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Exp No: 5. Retrieval of information from nucleotide databases

Aim: To retrieve nucleotide sequences and their annotations from various nucleotide databases

Procedure: The nucleotide sequences under interest can be retrieved from any one of the suitable databases such as Gene, UniGene, EST, Nucleotide and Species-specific Genome repository. Using name of the protein, a general search can be carried out at the beginning and then species-specific sequences can be obtained in a systematic search options. Following are the data of myoglobin (*Homo sapiens*) as taken from the ‘Nucleotide (NCBI database)’.

Homo sapiens myoglobin (MB) gene, complete cds

GenBank: AH002877.2

```
LOCUS          AH002877                6889 bp      DNA      linear    PRI 01-AUG-
2016
DEFINITION    Homo sapiens myoglobin (MB) gene, complete cds.
ACCESSION     AH002877 M10090 M14602 M14603
VERSION        AH002877.2
KEYWORDS       Fok family repetitive sequence; direct repeat; myoglobin; repeat
region; tandem repeat.
SOURCE         Homo sapiens (human)
ORGANISM       Homo sapiens
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
Catarrhini; Hominidae; Homo.
REFERENCE     1 (bases 1 to 6889)
AUTHORS       Akaboshi, E.
TITLE         Cloning of the human myoglobin gene
JOURNAL       Gene 33 (3), 241-249 (1985)
PUBMED        2989088
COMMENT        On or before Aug 1, 2016 this sequence version replaced M10090.1,
M14602.1, M14603.1, AH002877.1.
FEATURES      Location/Qualifiers
source        1..6889
              /organism="Homo sapiens"
              /mol_type="genomic DNA"
              /db_xref="taxon:9606"
              /map="22q11.2-qter"
repeat region 1100..1750
              /note="Fok family repetitive element"
prim transcript 1895..>2552
              /note="myoglobin mRNA and intron"
gene          1965..5924
              /gene="MB"
```

KARPAGAM ACADEMY OF HIGHER EDUCATION

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

```
CDS          join(1965..2059,4067..4289,5778..5924)
             /gene="MB"
             /note="myoglobin"
             /codon_start=1
             /protein_id="AAA59595.1"


```

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

/note="myoglobin; G00-119-378"
/number=3

ORIGIN 2 bp upstream of RsaI site.

1 gtactgtatt ttcatttcctc ttagttatct ccctaaaaag actctgagtt ccttgaacac
61 aggaagggtgt tttatggat tttgttatcc tcagcatgtc gcagtgtctg acacacagta
121 ggtgctctat cactgtgaga gggatggatg gatgggtgga gttacagatg gatagaagga
181 tagatggagg gatgggtgga tgatggatgg atagatggat ggagggggga tggatgaatgg
241 agggataatg agtggatgaa tgagggatg ggtggatgaa tggatggagg gatggaggaa
301 cagatagata gatggaggaa tgggtgggt atggatggat agatggatgg agggagggat
361 gatgaatgaa gggataatgaa atggatgaa gagggatgg gtggatggat gaatggaggg
421 atgatgggtg gatgaatgaa ttgagggatg gatggatgaa cacatggatg gatggataga
481 tggatagatg gaggaactgg tggatttgg atggatgggt ggtggatag atgaatgaat
541 gcctggatag acaaagagat gatggataga tgaatagatg aattaaggaa tgcggatag
601 atggagggat tgatagatgt tggatggatg ggtggatggat ggtatagatgatgatgcat
661 ggatagacaa agagatgatg gatggatgaa ttaagggatg acatggatg gatggatgaa
721 gtaactggat ggacaagtgg ataaatggat agatggatgaa atacctgatgatgatgaa
781 aggatgcatt gatgtatgaa aaggctaact atcctccact ctctttctt gcaaaaccat
841 ccaccattt actcaataaaa catttattca gttcaaaactt ggcacaaaagc accatgtgag
901 gccccaaagaga tacgtgggtt aataaaaacag agctcctgcc ctcctgaaaa ctgcaaaagaa
961 aggggcgtgg cttccgttgt tcaaattccca actctgcac cgactagctg tacatcagtgt
1021 atgtttccctt actttctctc aattaaatag ggataatgtc agtacccatc acattgggag
1081 gtcttgccc gattaaatgaa gttaccaat gccaagtgtt tggacaggg cctggcaccc
1141 agcaaagtct cttgtgagtg ctggctgta ttatcctaatt ggagaagatg gcatgaaaac
1201 cagaaatag gatgccctt gggaaagcaat gcaacaggaa cttacacaaa gaaaggaaag
1261 gagaaagcaa ttagtgggtg ctcaaaaggag tatgtcaaga aaaactttc agagggaaac
1321 ctttgagcag ggtcatgaaa acaggatgtc tctaagagat tggacttg cctggacca
1381 cctggctata agcacaaaac catccggttc cttctgtca cttctggcg gtgaggggtc
1441 tctggcaaaag gggcagaagg tgcgtgagag gttgcataatg gccaggactg tctggggcc
1501 agccggggca cctggctggcc aagcttagaa acatgacagg tcccttggg aggctgacc
1561 gcagggagcg ttgggtttca ggctgctggc gtccggctt gtggccct ttctgtggc
1621 tatgagagtc cagacagtgc ccaacccctt ccccttccat ccacacgcac aaccacccca
1681 cccctgtgg cctgagctgt cctgcctcgc cacaatggca cctgcctaa aatagcttcc
1741 catgtgaggg ctagagaaag gaaaagatta gaccctccct ggatgagaga gagaaagtga
1801 aggagggcag gggagggggaa cagcgagca ttgagcgatc ttgtcaagc atcccaag
1861 gtataaaaac gcccggggaa ccaggcagcc tcaaacccca gctgtgggg ccaggacacc
1921 cagtgagccc atactgctc ttttgttctt cttcagactg cgccatgggg ctcagcgacg
1981 gggaaatggca gttgggtgtc aacgtctggg ggaagggtgaa ggctgacatc ccaggccatg
2041 ggcaggaagt cctcatcagg taaaaggaag agattccatt gcccgtccca cccacacccct
2101 aagatcaagg gtgttcagct gcaagggtgaa aagtttgcac gtgggttgg tcagttggct
2161 gcattagtta aggggtttag aacggtcact tgctttctt ttgcattaa gtgtcaggaa
2221 ttggactcag gagagggaaa ggagccattt caggtgtatg tcagcagctg gagaaagcat
2281 gagaatcaaa cctaggatgc tcagagtccaa ccaggaagaa ttttagaatt atagacagtc
2341 agagttaaca agggcctga gagattttgt acagccacct ctcttacagg atgaggacaa
2401 aaagcgactg agaagggggag gacatttcca gagtcacagc tcattaaatg ctctttaaagt
2461 gtcaaggtta agacatgctc ttcaagggaa gacagatctg gttctagact tggctctgccc
2521 actgagccac tgggtgaccc ttggaaaggt ac
[gap 100 bp] [Expand Ns](#)

2653 gaattctg agaattgtct aaacccagga ggtggagggtt gcggtgagca
2701 gagattgcac cactgcactc cagcctggc aacagagcc gattccatct caaaaaagaa
2761 aaacaaaaaa caaaaaagcc atgaactcat tttcaggttg aggagctcag catcctgggt
2821 gtgaaatacc ctcctcataa aaccctggaa tggagactac ggggatcagg tgcttccttgg
2881 tgacaacttc tgggtgaccc ggctcaggc gcaactggaa gtgtggccac aatacataact
2941 gtgtactttt acaaggatgt cacagagcc gggatcataa aagaggagc tttcaagga

KARPAGAM ACADEMY OF HIGHER EDUCATION

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

3001 actgaaacca ttagacagga gagagagccc tggcagaca gggttgcccg tgccaaacat
3061 ttcatgtgc gcacaaggga aagggtggg gttatgaaac tggccattt tggtttagg
3121 tctggctct gcccttagct agccaagtga cttggccac ttatcttgt ggtcttccat
3181 gagtaaaagg cgaaactca ctccatcca gagggcaggt ctgactccc ttaaccagca
3241 cccacctgct cacagcagga aggactgagg tctaaagctg gaggtggca ggaaggactg
3301 aggtctaaag ctggaggtgg gcaggaagga ccgaggctca aagctggagg tgggcaggaa
3361 ggaccgagg tctaaagctgg aggtggctgc tcagagtccc agcagaggcc tctggggcac
3421 ctcactgagt gcctggcagg agtgggtgcc tgcgtcaggg ctgggttgag ttgctccac
3481 caggaccctt cgtcatctgc acagtggagg gactggagg ttcatgaggt cacagcttg
3541 gctaaaaca agcaagaggt ttctgttgt gaggattgt ctggagtgga atggccctca
3601 caggttaggag tgaggctctt gttagctagag gtatthaagc agctgaagga caatccctgg
3661 gcaggaagct gcagagatgg tcgcagcgtg gactagaact gctgtttgg taactcagac
3721 ctcattccag cctggcttct ctggacagca cccctgcaat agtgagctgg tgactttacg
3781 cctcagaacc tcgggttcta catctgtaaa atggaaatta tatgacactc actatgtgcc
3841 agacaccctg ttggtacata gcacacacta tctacttaa tcctcaagt agggacaagt
3901 tatccccatc ctttatatga ggaagctgag gcacagagag gtgaagtgaa tgcccaagg
3961 tcacacagct gggaaagacag ggagctaaac ttgaactcta gtctggctgc cccagacact
4021 cacaccgcac ctccatgccc gactccagcc ttccctgtgc ccacaggtc tttaagggtc
4081 acccagagac tctggagaag tttgacaagt tcaagcacct gaagtcaagag gacgagatga
4141 aggcatctga ggacttaaag aagcatgtg ccactgtgct caccgcctg ggtggcatcc
4201 ttaagaagaa gggcatcat gaggcagaga ttaagccctt ggcacagtgc catgccacca
4261 agcacaagat cccctgtaa tacctggagg taggaggcag agcctggca ggtgggagga
4321 tgcggggaaag gcctcgggtg gggcaatggg atctgggtc gagtccaagc tcagccacta
4381 acttgtgggta tgacctatgc cactcttctc tgtgccccag gttctcatt tgtaaagggg
4441 actgccaccc actttgcctt cctcctgggaa ttgttgagaa tgaacacatt tagcattttt
4501 aatttagtat gccaaattca catcttatta ccaaagagga aaggagagg gatatattggg
4561 tgcaaaattt gcatttcctc catgggtagg taccattatc atatccactt gatagatggg
4621 gaaactgagg ctcacagagg ttaagcagct tgccacggc cacaggaggt ggataatggc
4681 agacccaaga ttcaaacgca ggtctctatt actacagaac cccagccct aactgctgtg
4741 ccactgggag tctgtacat gcaggactta tgtggcagga gctcagcaag tggggctcaa
4801 ttgggggtgg gggtagccag caggctggct ctattgggtc cagcatctc acagatgaag
4861 agacaggacc tcgggttcca gcacaagacaa ttgtgttggaa cctcctgaga tgggttggaa
4921 agttgggtgg atcagggttg gggcaggag cctggcttc aggttgtgt tctataactg
4981 gtgggaggag gcgatttggg gagaggaggg agctgggat gaaggaccac agggacaggt
5041 gcatcccccg agggtagaaa cagcaggaag tctgggtcag ccatgaggat taggatgtgg
5101 tgatagctac cccctggat gggccacagt gagatttc tgccatgcct agcacatgca
5161 tccatcctca aagttgcctc atggccaaaa tgactgcaag agctccagcc agctcttcta
5221 tattcccaac tggaagcagg agaaagagag gaatgctctc ttttgaggag tttaaggagt
5281 cccagaaatc tcatccaaca attttattt catctcatttgc [gap 100 bp] [Expand Ns](#)
5424 gagctct cacctggttt cagtgggtc tacatcctga
5461 tggagtggag gggctgtga gtaagagctg gggctccgga gccggccctc ctgggtccaa
5521 atgtcccttc cattcaaccc cccctcgcc cagttctgc atctgtaaat cgagggcagt
5581 tgttagtatct atctcacagt gttgtgggg atcaaagggg ttcatccgtg gagatcacac
5641 agactctcac ctgggccta gcaagtgtc aatacacgt cctggataa agagaaggta
5701 ggaggacaac tgactccat ctggccctg gctgtccca ccctggtgac cattttctct
5761 cctcaccctc cctgcagttc atctcggaaat gcatcatcca gttctgtcag agcaagcatc
5821 cccggggactt tgggtctgtat gccgaggggg ccatgaacaa ggccctggag ctgtccgg
5881 aggacatggc ctccaactac aaggagctgg gcttccaggg ctggccct gccgctccca
5941 ccccccaccca tctggggccc gggttcaaga gagagcgggg tctgatctcg ttagccata
6001 tagagttgc ttctgttgt ctgctttgtt tagtagaggt gggcaggagg agctgggggg
6061 ctggggctgg ggtgttgaag ttggcttgc atgcccagcg atgcgcctcc ctgtggatgc
6121 tcatcaccctt gggaaaccggg agtggccctt ggctcactgt gttctgtcatg gttggatct

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

6181 gaatttaattg tcctttcttc taaatcccaa ccgaacttct tccaacctcc aaactggctg
6241 taaccccaaa tccaagccat taactacacc tgacagtagc aattgtctga ttaatcactg
6301 gccccttcaa gacagcagaa tgtcccttg caatgaggag gagatctggg ctgggcgggc
6361 cagctgggaa agcattgac tatctggAAC ttgtgtgtgc ctccctcaggat atggcagtga
6421 ctcacctgggt tttataaaaa caacctgcaa catctcaggat tctgcctggc attttcatc
6481 tccttagagta aatgatgccc ccaccagcac cagcatcaag gaagaaatgg gaggaaggca
6541 gaccctgggc ttgtgtgtgc agagagcctc aggaaagagg agaaggggag gaggaaaggc
6601 aggagggtga gaggacagg agcccacccct ccctgggcca ccgctcagag gcaggcccag
6661 tgcagggcat gggaaaatgg aagggacagg cttggcccca gccttggag caccttctct
6721 tcgggggagg tgggaggcag cgaacagacc tctgcaatac gaggagagag tgacaggtgc
6781 gccaggctgt gggAACCCAG aggagagggg aagccatcat catcatggct gcaataacctt
6841 cagtaacgtg ggaaggtcac cctgctagta agtggcagag ctgggactc

//

Similarly, the nucleotide sequences from other databases can also be analyzed in a forthright manner.

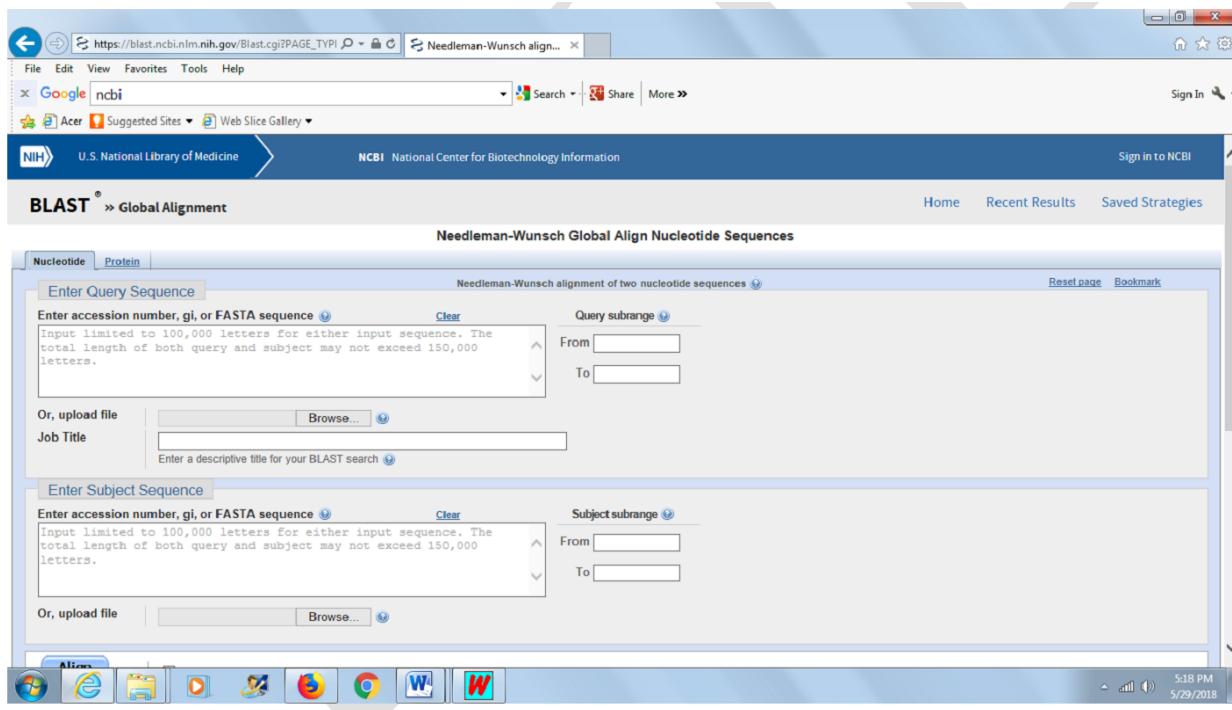
Results:

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

Exp No: 6. Sequence alignment using BLAST

Aim: To examine pair-wise alignments of protein sequences using BLAST

Procedure: Retrieve the two protein sequences under interest either from primary or structural databases and save them in ‘Fasta’ format. Either the sequences or accession IDs of the sequences can be pasted in the space provided under categories of ‘Query sequence’ and ‘Subject sequence’ of the algorithm. Upon selecting appropriate alignment algorithm and gap penalty options, the program can be run for aligning the sequences.



For example, two cardiotoxins bearing PDB IDs of 1CRF and 1CVO were aligned and the results are shown below herein.

Seq2

Sequence ID: Query_206921Length: 62Number of Matches: 1

Related Information

Prepared by: Dr. Selvakumar, S. Assistant Professor, Department of Biotechnology, KAHE.
Page 23

**CLASS : III B. Sc., BT
 COURSE NAME : BIOINFORMATICS PRACTICAL
 BATCH : 2017 – 2020
 COURSE CODE : 17BTU512A**

Range 1: 1 to 62 [Graphics](#) [Next Match](#) [Previous Match](#) [First Match](#)

Alignment statistics for match #1

NW Score	Identities	Positives	Gaps	Frame
213	40/62(65%)	46/62(74%)	2/62(3%)	

Features:

Query 1 LKC-NKLVPLFYKTCPAGKNLCYKMFMVS-NLTVPVKRGCIDVCPKNSALVKYVCCNTDR 58

LKC N +P YKTCP GKNLC+K + L PVKRG D CPKNSAL+KYVCC+TD+

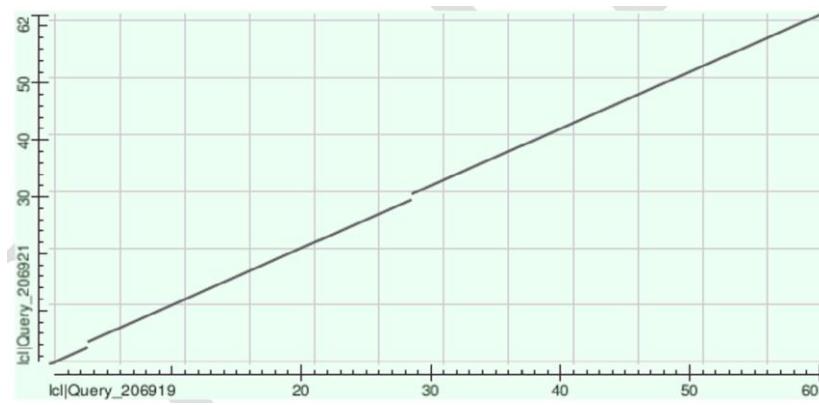
Sbjct 1 LKCHNTQLPFIYKTCPEGKNLCFKATLKKFPLKFPVKRGCADNCPKNSALLKYVCCSTDK 60

Query 59 CN 60

CN

Sbjct 61 CN 62

The dot plot for the sequences is depicted in the figure shown herein.

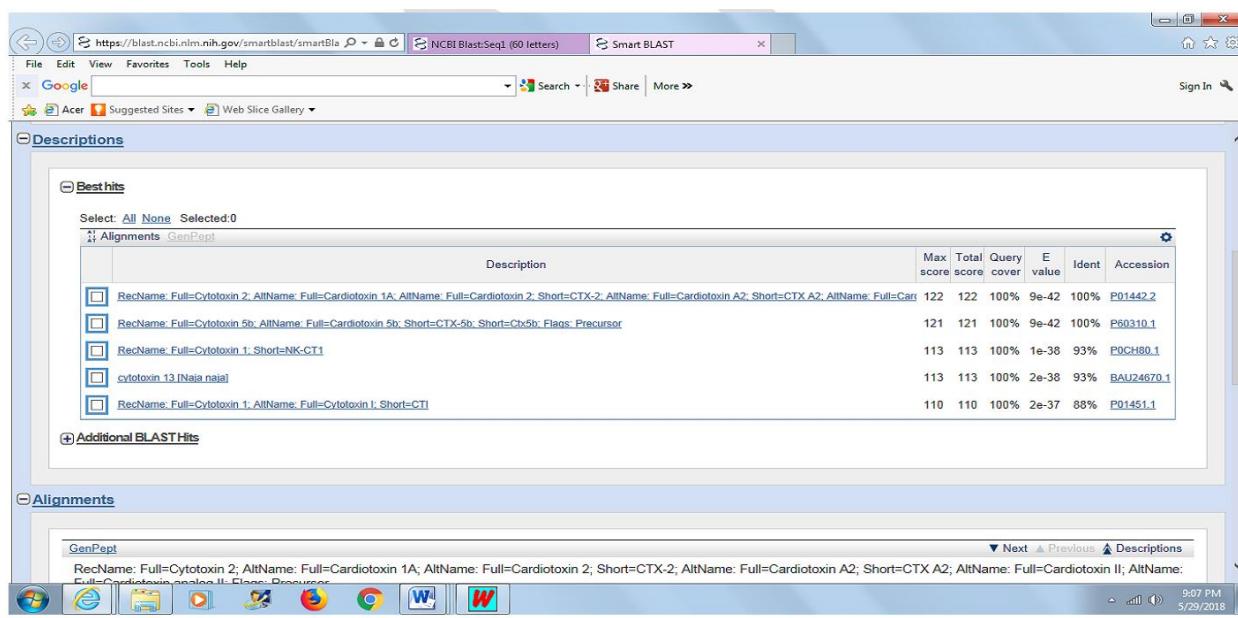
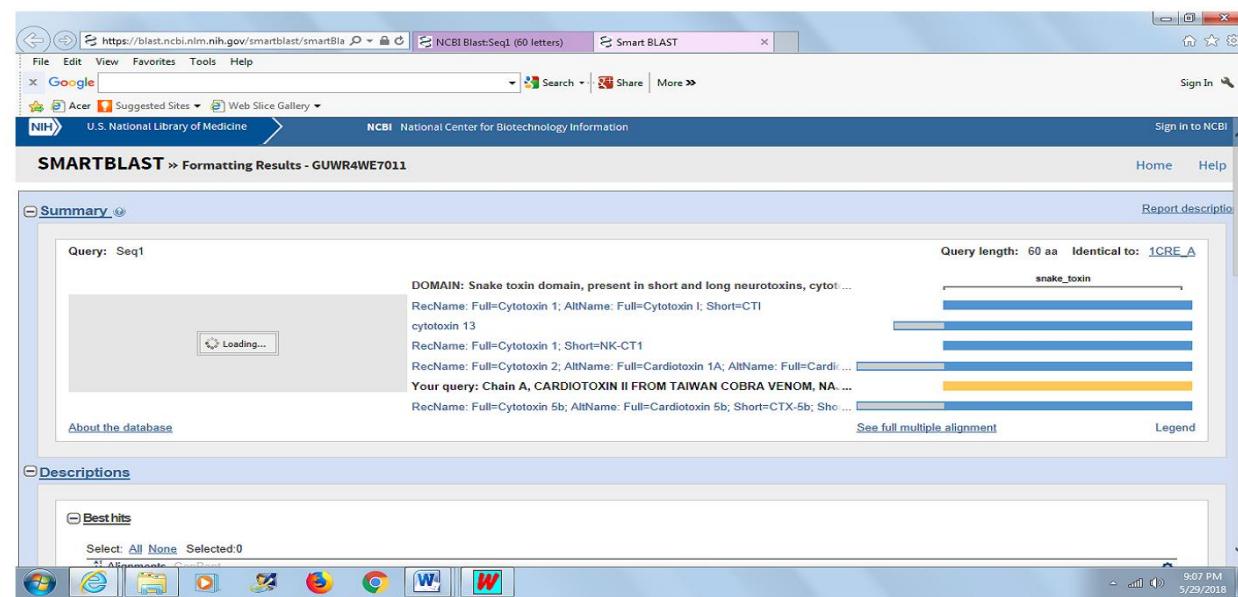


The sequence alignment can also be performed by using SmartBLAST and the results are as

KARPAGAM ACADEMY OF HIGHER EDUCATION

**CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A**

shown below herein.



Results:

Prepared by: Dr. Selvakumar, S. Assistant Professor, Department of Biotechnology, KAHE.

Page 25

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Exp No: 7. Multiple sequence alignment using ClustalW/Clustal Omega

Aim: To examine multiple sequences alignment of protein sequences using Clustal Omega

Procedure: Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between three or more sequences. This tool can align up to 4000 sequences or a maximum file size of 4 MB.

The following 11 sequences were first retrieved from PDB database and the data were stored in ‘Fasta’ format.

>1CRF Naja Naja Atra
LKCNKLVPLFYKTCPAGKNLCYKMFMSNLTVVKRGCIDVCPKNSALVKYVCCNTDRCN

>1CRE Naja Naja Atra
LKCNKLVPLFYKTCPAGKNLCYKMFMSNLTVVKRGCIDVCPKNSALVKYVCCNTDRCN

>1CHV Naja Naja Atra
LKCNKLVPLFYKTCPAGKNLCYKMFMSNKMVVKRGCIDVCPKSSLLVKYVCCNTDRCN

>P60301.1
MKTLLLTLVVVTIVCLDLGYTLKCNKLVPLFYKTCPAGKNLCYKMFMVATPKVPVKRGCIDVCPKSSLLVKYVCCNTDRCN

>P80245.2
MKTLLLTLVVVTIVCLDLGYTLKCNQLIPIASKTCPAGKNLCYKMFMSDLTIPVKRGCIDVCPKNSLLVKYVCCNTDRCN

>P60304.1
MKTLLLTLVVVTIVCLDLGYTLKCNKLPIASKTCPAGKNLCYKMFMSDLTIPVKRGCIDVCPKNSLLVKYVCCNTDRCN

>1CXO
LKCNQLIPPFWKTCPKGKNCYKMTMRAAPMVPVKRGCIDVCPKSSLLIKYMCCNTDKCN

>1CXN
LKCNQLIPPFWKTCPKGKNCYKMTMRAAPMVPVKRGCIDVCPKSSLLIKYMCCNTDKCN

>P01442.2
MKTLLLTLVVVTIVCLDLGYTLKCNKLVPLFYKTCPAGKNLCYKMFMSNLTVVKRGCIDVCPKNSALVKYVCCNTDRCN

>1CVO
LKCHNTQLPFIYKTCPEGKNLCFKATLKKFPLKFPVKRGCADNCPKNSALLKYVCCSTDKCN

>P62375.1
MKTLLLTMVVVTIVCLDLGYTLKCHNTQLPFIYKTCPEGKNLCFKATLKKFPLKFPVKRGCADNCPKNSALLKYVCCSTDKCN

All the 11 sequences were pasted one-by-one in the ‘sequence space’ provided in the Clustal

CLASS : III B. Sc., BT
COURSE NAME : BIOINFORMATICS PRACTICAL
BATCH : 2017 – 2020
COURSE CODE : 17BTU512A

Omega program. After setting appropriate alignment algorithm, matrix and output format, the program was executed and the outcomes are shown below herein'

CLUSTAL O(1.2.4) multiple sequence alignment

```
1CVO -----  
LKCHNTQLPFIYKTCPEGKNLCFKATLKKFPLKF  
PVKRG 39 P62375.1  
MKTLLLTMVVVTIVCLDLGYTLKCHNTQLPFIYKTCPEGKNLCFKATLKKFPLKF PVKRG 60  
1CXO ----- LKCNQ-LIPPFWKTCPKGKNLCYKMTMRA-  
APMVPVKRG 37  
1CXN ----- LKCNQ-LIPPFWKTCPKGKNLCYKMTMRA-  
APMVPVKRG 37  
P60301.1 MKTLLLTLVVVTIVCLDLGYTLKCNK-LVPLFYKTCPAGKNLCYKMFVA-  
TPKVPVKRG 58  
P80245.2 MKTLLLTLVVVTIVCLDLGYTLKCNQ-LIPPFYKTCAAGKNLCYKMFVA-  
APKVPVKRG 58  
P60304.1 MKTLLLTLVVVTIVCLDLGYTLKCNK-LIPIASKTCPAGKNLCYKMFMS-  
DLTIPVKRG 58  
1CRF ----- LKCNK-LVPLFYKTCPAGKNLCYKMFVMS-  
NLTVPVKRG 37  
1CRE ----- LKCNK-LVPLFYKTCPAGKNLCYKMFVMS-  
NLTVPVKRG 37  
P01442.2 MKTLLLTLVVVTIVCLDLGYTLKCNK-LVPLFYKTCPAGKNLCYKMFVMS-  
NLTVPVKRG 58  
1CHV ----- LKCNK-LVPLFYKTCPAGKNLCYKMFVMS-  
NKMVPVKRG 37  
***** : : : * *** * ***** : * :  
* *****  
1CVO CADNCPKNSALLKYVCCSTDKCN 62  
P62375.1 CADNCPKNSALLKYVCCSTDKCN 83  
1CXO CIDVCPKSSLLIKYMCCNTDKCN 60  
1CXN CIDVCPKSSLLIKYMCCNTDKCN 60  
P60301.1 CIDVCPKSSLLVKYVCCNTDRCN 81  
P80245.2 CIDVCPKSSLLVKYVCCNTDRCN 81  
P60304.1 CIDVCPKNSLLVKYVCCNTDRCN 81  
1CRF CIDVCPKNSALVKYVCCNTDRCN 60  
1CER CIDVCPKNSALVKYVCCNTDRCN 60  
P01442.2 CIDVCPKNSALVKYVCCNTDRCN 81  
1CHV CIDVCPKSSLLVKYVCCNTDRCN 60  
* * *** . * * : * : * * . * * : **
```

Results: