



Dept. of Computer Science and Engineering  
University of Rajshahi  
www.ru.ac.bd

Dr. Shamim Ahmad

# Databases

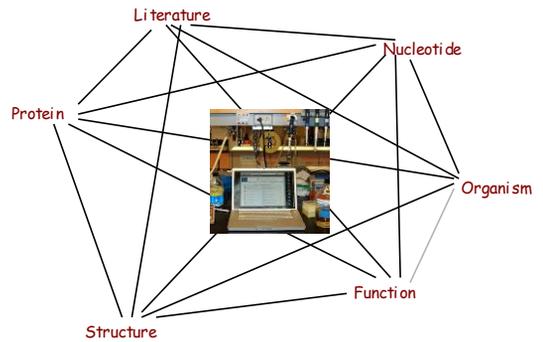


Organizing information in the post-genomic era  
The rise of bioinformatics

What is a database?

Which databases are important for molecular cell biology research?

How is information processed in databases?



Biological databases use different organizing principles

Hyperlinks connect records in different databases

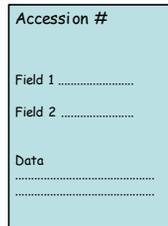
## Databases are organized collections of information

Information is stored in records

Databases assign each record a unique **accession number** using their own numbering system

**Fields** are used to cross-reference the data. Records can be searched by fields.

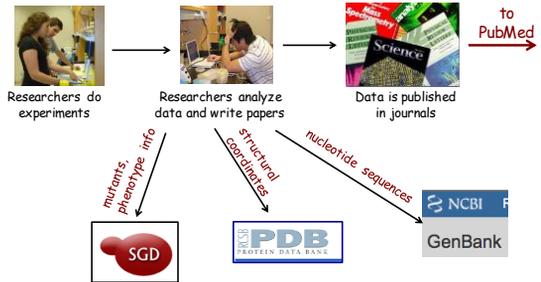
**Data** is entered in the record using a defined format



**Bioinformaticians** work with computer scientists to set up the database structure

**Curators** review and link records within and between databases

The information in databases ultimately derives from experimental data



**Curators** will process the submissions and link entries in different databases

## Database Examples in Bioinformatics

- **Primary (archival)**
  - GenBank/EMBL/DBJ (seqs)
  - PDB (protein structures)
  - Medline (literature)
  - IMEx databases (protein interactions)
- **Secondary (curated)**
  - RefSeq (seqs)
  - UniProt-SwissProt (seqs)
  - Taxon (taxonomy)
  - PROSITE (binding sites)
  - OMIM (genetics literature/reviews)
  - IMEx databases (protein interactions)

7

## Sequence Databases

### ➤ DNA

- ❖ NCBI: GenBank -> RefSeq
- ❖ EBI: EMBL

### ➤ Protein

- ❖ NCBI: GenPept
- ❖ EBI: UniProt: TrEMBL -> UniProt: Swiss-Prot

TrEMBL = "translated EMBL"

**National Center for Biotechnology Information** [www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)

**European Bioinformatics Institute** [www.ebi.ac.uk](http://www.ebi.ac.uk)

8

## GenBank

- Direct submissions from individual laboratories
- Bulk submissions from large-scale sequencing centers.

### Direct submissions to GenBank

- **BankIT**, which is a Web-based form
- **Sequin** stand-alone submission program

### Upon receipt of a sequence submission

- GenBank staff examines the originality of the data
- Assigns an accession number to the sequence
- Performs quality assurance checks

## GenBank

- GenBank has become an important database for research in biological fields
- It has grown in recent years at an exponential rate by doubling roughly every 18 months
- Entries are retrievable by Entrez
- Downloadable by FTP

### Bulk submissions

- Expressed Sequence Tag (EST)
- Sequence-tagged site (STS)
- Genome Survey Sequence (GSS)
- High-Throughput Genome Sequence (HTGS)



Growth in GenBank base pairs, 1982 to 2007, on a semi-log scale

## GenBank

An **expressed sequence tag (EST)** is a short sub-sequence of a cDNA sequence

A **sequence-tagged site (or STS)** is a short (200 to 500 base pair) DNA sequence that has a single occurrence in the genome

**Genome Survey Sequences (GSS)** are nucleotide sequences similar to EST's that the only difference is that most of them are genomic in origin, rather than mRNA

## Complementary DNA

- **Complementary DNA (cDNA)** is DNA synthesized from a single stranded RNA (e.g., messenger RNA (mRNA))
- DNA is derived from mRNA, so it contains only **exons**, with no introns.
- cDNA is also produced naturally by retroviruses (such as HIV-1, HIV-2, simian immunodeficiency virus, etc)

## Complementary DNA

- cDNA is often used to clone eukaryotic genes in prokaryotes.
- When scientists want to express a specific protein in a cell that does not normally express that protein, they will transfer the cDNA that codes for the protein to the recipient cell.
- cDNA is also produced naturally by retroviruses (such as HIV-1, HIV-2, simian immunodeficiency virus, etc)

Organism	base pairs
<i>Homo sapiens</i>	16,310,774,187
<i>Mus musculus</i>	9,974,977,889
<i>Rattus norvegicus</i>	6,521,253,272
<i>Bos taurus</i>	5,386,258,455
<i>Zea mays</i>	5,062,731,057
<i>Sus scrofa</i>	4,887,861,860
<i>Danio rerio</i>	3,120,857,462
<i>Strongylocentrotus purpuratus</i>	1,435,236,534
<i>Macaca mulatta</i>	1,256,203,101
<i>Oryza sativa Japonica Group</i>	1,255,686,573
<i>Nicotiana tabacum</i>	1,197,357,811
<i>Xenopus (Silurana) tropicalis</i>	1,249,938,611
<i>Drosophila melanogaster</i>	1,119,965,220
<i>Pan troglodytes</i>	1,008,323,292
<i>Arabidopsis thaliana</i>	1,144,226,616
<i>Canis lupus familiaris</i>	951,238,343
<i>Vitis vinifera</i>	999,010,073
<i>Gallus gallus</i>	899,631,338
<i>Glycine max</i>	906,638,854
<i>Triticum aestivum</i>	898,689,329



## UniProt

- UniProtKB
  - Swiss-Prot  
(reviewed, manually annotated entries)  
1.93 B Amino acid, 1914.3
  - TrEMBL  
(unreviewed, automatically annotated entries)  
14 M Amino Acid, 1914.3

## UniProt

- UniProtKB (UniProt Knowledgebase)
  - Information extracted from scientific literature and [biocurator](#)-evaluated computational analysis
  - To provide all known relevant information about a particular protein

## UniProtKB (Annotation)

- Protein and gene names
- Function
- Enzyme-specific
- Subcellular location
- Protein-protein interactions
- Pattern of expression
- Ion-, substrate- and cofactor-binding sites
- Protein variant forms produced by natural genetic variation, RNA editing,

## UniProt: Swiss-Prot – An example of curated, reviewed annotation

- Incorporates:
  - ✓ Function of the protein
  - ✓ Subcellular localization of protein
  - ✓ Post-translational modification
  - ✓ Domains and sites
  - ✓ Secondary structure
  - ✓ Quaternary structure
  - ✓ Similarities to other proteins
  - ✓ Diseases associated with deficiencies in the protein
  - ✓ Sequence conflicts, variants, etc.

24

## Swiss-Prot

- Reliable protein sequences associated
  - With a high level of annotation
  - Description of the function of a protein
  - Its domain structure
  - Post-translational modifications
  - Variants
  - Minimal level of redundancy
  - High level of integration with other databases.

## TrEMBL (Translated EMBL Nucleotide Sequence Data Library)

was created to provide automated annotations for those proteins not in Swiss-Prot

## UniProt Archive (UniParc)

- It is a **comprehensive** and **non-redundant** database,
- **Proteins may exist in several different source databases**, and in **multiple copies** in the same database
- Information extracted from **scientific literature** and [biocurator](#)-evaluated **computational analysis**
- Each sequence is given a stable and unique identifier (UPI)

## UniProt Archive (UniParc)

- INSDC **EMBL-Bank/DDBJ/GenBank** nucleotide sequence databases
- [Ensembl](#)
- **European Patent Office (EPO)**
- Explain Post-Translational Modification **FlyBase**: the primary repository of genetic and molecular data for the insect family Drosophilidae (FlyBase)
- **H-Invitational Database (H-Inv)**
- **International Protein Index (IPI)**
- **Japan Patent Office (JPO)**

## UniRef

The UniProt Reference Clusters (UniRef) consist of three databases

- Protein sequences from UniProtKB
- Selected UniParc records

## The Protein Data Bank(PDB)

- Crystallographic database for the three-dimensional structural data of large biological molecules,
  - Proteins
  - Nucleic acids.
- The data, typically obtained by
  - X-ray Crystallography
  - NMR spectroscopy,
  - Cryo-electron microscopy

## The Protein Data Bank(PDB)

Experimental Method	Proteins	Nucleic Acids	Protein/Nucleic complexes			Total
			Acid	Other		
<a href="#">X-ray diffraction</a>	106595	1820	5471	4	113890	
<a href="#">NMR</a>	10296	1190	241	8	11735	
<a href="#">Electron microscopy</a>	1021	30	367	0	1418	
Hybrid	99	3	2	1	105	
Other	181	4	6	13	204	
<b>Total:</b>	<b>118192</b>	<b>3047</b>	<b>6087</b>	<b>26</b>	<b>127352</b>	

NCBI maintains both primary and derivative databases  
We'll look at three of them



PubMed is the premier literature database in the world

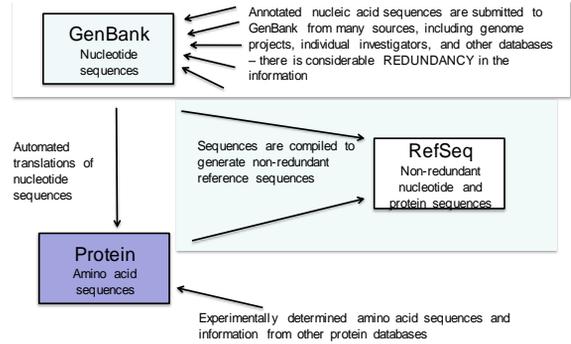
Questions for today:

What is a database?

Which databases are important for molecular cell biology research?

How is information processed in databases?

Curators are responsible for data flow between the NCBI databases



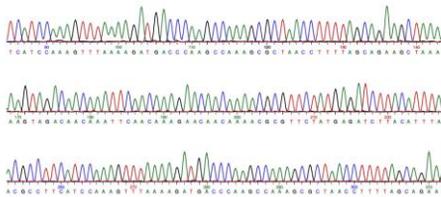
Most records in the Protein database have been derived by automated translation of nucleotide sequences



On a larger scale: Genome projects have produced the reference sequences in nucleotide databases

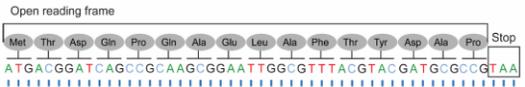
(robots and computers do much of the work)

1. Pieces of chromosomal DNA are sequenced, each ~1000 bp long



*S. cerevisiae* genome is ~12 Mbp - how many reads would be necessary to cover each base pair in the genome once?

3. Chromosomal sequences are analyzed for the presence of potential transcripts (open reading frames: ORFs)



ORFs are characterized by an under-representation of stop codons

ORF-finding computer algorithms look for sequences that

- begin with a methionine
- methionine is separated from a stop codon in the same reading frame by a large number of amino acids (often 100, equiv. to 300bp)

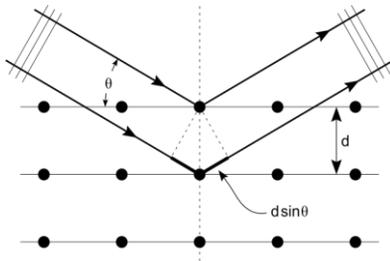
GenBank NM\_##### records are predicted ORFs

4. Protein sequences are computationally predicted from ORF sequences

GenBank NP\_##### records

## X-ray crystallography

For determining the atomic and molecular structure of a crystal, in which the crystalline atoms cause a beam of incident X-rays to diffract into many specific directions



## Nuclear magnetic resonance spectroscopy of proteins (NMR)

- It is a field of structural biology
- NMR spectroscopy is used
- To obtain information about the structure and dynamics of
  - proteins, and also nucleic acids,

Entrez  
Global Query Cross-Database  
Search System  
[CDS](#)  
Coding sequence;

### The locus name:

The first three characters usually designated the organism

Accession Number [ACCN]

### Sequence Length

Sequence Length [SLEN] 2500:2600[SLEN]

### Molecule Type

The type of molecule that was sequenced. In this example, the molecule type is [DNA](#)

Properties [PROP]

biomol\_genomic, biomol\_mRNA {PROP}

### GenBank Division

gbdiv\_pri, gbdiv\_est, Properties [PROP]

human[ORGN] NOT gbdiv\_est[PROP]

- 1) To retrieve records within a range of lengths, use the colon as the range operator, e.g., 2500:2600[SLEN].
- 2) To retrieve all sequences shorter than a certain number, use 2 as the lower bound, e.g., 2:100[SLEN].
- 3) To retrieve all sequences longer than a certain number, use a series of 9's as the upper bound, e.g., 325000:99999999[SLEN].

### Molecule Type

The type of molecule that was sequenced. In this example, the molecule type is [DNA](#)

- genomic DNA
- genomic RNA
- precursor RNA
- mRNA (cDNA)
- ribosomal RNA
- transfer RNA,
- small nuclear RNA
- small cytoplasmic RNA.

### Modification Date

Modification Date [MDAT]

To retrieve records modified between two dates, 1999/07/25:1999/07/31[MDAT].

### DEFINITION

Brief description of sequence

Coding region (CDS),

Word [TITL]

### ACCESSION

The unique identifier for a sequence record

### VERSION

If there is any change to the sequence data (even a single base), the version number will be increased, e.g., U12345.1 ? U12345.2,

### KEYWORDS

Word or phrase describing the sequence

### SOURCE

Free-format information including an abbreviated form of the organism name, sometimes followed by a molecule type.

### Organism

The formal scientific name for the source organism

### REFERENCE

Publications by the authors

### AUTHORS

List of authors in the order in which they appear in the cited article

### TITLE

Title of the published work or tentative title of an unpublished work.

< symbol indicates partial on the 5' end. Example: <1..206

\_> symbol indicates partial on the 3' end.  
 Example: **4821..>5028**

The GenBank database is divided into 18 divisions:

1. PRI - primate sequences
2. ROD - rodent sequences
3. MAM - other mammalian sequences
4. VRT - other vertebrate sequences
5. INV - invertebrate sequences
6. **PLN - plant, fungal, and algal sequences**
7. BCT - bacterial sequences
8. **VRL - viral sequences**
9. PHG - bacteriophage sequences
10. SYN - synthetic sequences
11. UNA - unannotated sequences
12. **EST - EST sequences (expressed sequence tags)**
13. PAT - patent sequences
14. **STS - STS sequences (sequence tagged sites)**
15. **GSS - GSS sequences (genome survey sequences)**
16. HTG - HTG sequences (high-throughput genomic sequences)
17. HTC - unfinished high-throughput cDNA sequencing
18. ENV - environmental sampling sequences

## FASTA format

- It is a text-based format for representing
  - Nucleotide sequences
  - Peptide sequences
- Nucleotides or amino acids are represented using single-letter codes
- The format also allows for sequence names and comments to precede the sequences.

- The simplicity of FASTA format
- Makes it easy to manipulate and parse sequences using text-processing tools and scripting languages
  - R programming language
  - Python
  - Ruby,
  - Perl

- A sequence in FASTA format is represented as a series of lines
- Each of which should be no longer than 120 characters and usually do not exceed 80 characters.
- The first line in a FASTA file starts either with a ">" (greater-than) symbol or, less frequently, a ";" (semicolon) and was taken as a comment.

```

:LCBO - Prolactin precursor - Bovine
: a sample sequence in FASTA format
MDSK655QK6SRLLLLLVVSNLLLCQ6VVVSTPVCPNGP6NCQVSLRDLFDRAVMVSHYIHDLS5
EWNNEFDKRYAKGK6FTTMAINS6HTSSLPFPEDKEQAQQTHHEVLM5LTL6LRSWINDPLYHL
VTEVR6MKG6APDAILSRATEIEENKRLLE6MEMIF6QVTP6AKETEPV6VWS6LP5LQTKD6E6
ARYSAFYNLHLCLRRD55KIDTYLKLNNCRIIYNNNN6*

>MCHU - Calmodulin - Human, rabbit, bovine, rat, and chicken
ADQLTEEQIAEFKEAFSLFDK6D6GTITTKEL6TVMRSL6QNPT6AELQDMINEVDAD6NGTID
FPEFLTMMARKMKD6TDEEEI6REA6RVF6DKD6NGYISA6ELRH6VMTNL6EKL6TDE6VDEMI6REA
DID6D6QVNYEEFVQMM6TAK6*

>gi|5524211|gb|AAD44166.1| cytochrome b [Elephas maximus maximus]
LQLYTHIGRNIYGSYLYSETWNTGIMLLITMATAFM6YVLPW6QMSF6W6ATVITNLFSAIPYI6TNLV
EWI6W6G6F6SV6KATLNRFFAFHILPFTMV6AL6GVHL6TFL6HET6SN6NPL6L6TSD6DKIP6PHY6TI6KD6FL6
LLTLILL6LLL6ALL6SP6ML6G6DP6NH6P6AD6PL6N6T6L6IK6PEWY6FL6FAY6IL6RS6VP6NK6L6G6VL6AL6FL6SIV6IL
6LMP6FL6HT6SK6RS6MML6RPL6SQ6L6FW6TL6TMD6LL6TL6TWI6G6SQ6V6EY6P6YTI6IG6MAS6YL6YF6SII6LA6FL6PI6A6G
6IENY

```

```

>SEQUENCE_1
MTEITAAMVKELRES T6AG6MMDCKNAL6SETN6DFDKA6VQL6REK6L6KA6KK6ADR6LAA6E6G
LV5VKV5DDFTIA6MRP5SYLSY6EDLDM6TF6VENEY6KALVA6ELEKENE6ERRR6LK6DPNK6PEHK
IPQF6ASRK6QL6SDA6IL6KE6AE6EK6EEL6KA6Q6KP6EKI6W6DN6IIP6GKM6NS6FI6AD6NS6QL6DS6KL6TL
M6QFY6VMDD6KKT6VE6QVIA6E6KE6KEF6G6KIKI6VE6FICF6EV6E6LE6KKT6EDFA6EVA6A6QL
>SEQUENCE_2
SATVSEINSETDFVAKNDQFIALTKD6T6TAHI6QS6NS6L6QS6VEEL6HS6TIN6VK6FEY6LKS6QI
ATI6ENL6VRR6FA6TL6KA6GAN6GV6VNGY6IHT6N6RV6GVVIA6AC6D6SA6EVA6SK6SR6DL6LR6QIC6MH

```

Nucleic Acid Code	Meaning	Mnemonic
A	A	<a href="#">Adenine</a>
C	C	<a href="#">Cytosine</a>
G	G	<a href="#">Guanine</a>
T	T	<a href="#">Thymine</a>
U	U	<a href="#">Uracil</a>
R	A or G	<a href="#">puRine</a>
Y	C, T or U	<a href="#">pYrimidines</a>
K	G, T or U	bases which are <a href="#">Ketones</a>
M	A or C	bases with <a href="#">aMino groups</a>
S	C or G	<b>Strong</b> interaction
W	A, T or U	<b>Weak</b> interaction
B	not A (i.e. C, G, T or U)	<b>B</b> comes after A
D	not C (i.e. A, G, T or U)	<b>D</b> comes after C
H	not G (i.e., A, C, T or U)	<b>H</b> comes after G
V	neither T nor U (i.e. A, C or G)	<b>V</b> comes after U
N	A C G T U	<b>Nucleic acid</b>
-	gap of indeterminate length	



This screenshot shows the NCBI Nucleotide search results for the query "(Coding region (CDS) AND 28002800[LEN])". The search results are displayed in a table with columns for Species, Nucleotide, and other details. The top results include:

- 1. 2574 bp linear DNA** from *Nicotiana benthamiana* calcium ATPase (NCA1) gene promoter region and partial cds. Accession: GQ591211.1. Organism: Nicotiana benthamiana. Database: GenBank. FASTA: GenBank. Statistics: GenBank.
- 2. 2568 bp linear DNA** from *Escherichia coli* dta-pylE gene region. Accession: KP1211.1. Organism: Escherichia coli. Database: GenBank. FASTA: GenBank. Statistics: GenBank.

The interface includes a search bar, filters, and options to view details or download the results.

This screenshot shows the My NCBI dashboard. It features a search bar at the top with the query "(Coding region (CDS) AND 28002800[LEN])". Below the search bar, there are sections for "My Bibliography", "Recent Activity", and "Filters". The "Recent Activity" section shows a list of recent searches and their results. The "Filters" section allows users to refine their search results based on various criteria.

This screenshot shows the NCBI PubMed Advanced Search Builder interface. It provides a structured way to build complex search queries. The interface includes a search bar, a "Builder" section with fields for "Date - Entrez" and "Date - Entrez", and a "History" section for tracking previous searches. The search criteria are currently set to "Date - Entrez" from "2017/1/1" to "2017/1/1".

This screenshot shows the NCBI PubMed search results for the query "(2019/04/01[Date - Entrez] - 2000[Date - Entrez])". The search results are displayed in a table with columns for Article types, Format, Summary, Sort, and other details. The top results include:

- 1. Are changes in binding related to changes in catalytic use among Snp64?** by Olga I. Denisovskaya, Ramstedt M. Address: 2019 Apr 21; doi: 10.1101/3981434. [Epub ahead of print]. PMID: 30879564. Status: Article.
- 2. Anticoagulation by the anti-oxidized phospholipid antibody EDS.** by Mohammad M. Carter B, Kikta J, Martin C, Black A, Blum R, Risher HL, et al. Address: 2019 Apr 21; doi: 10.1101/3981434. [Epub ahead of print]. PMID: 30879565. Status: Article.

The interface includes a search bar, filters, and options to view details or download the results.



## RefSeq: NCBI Reference Sequence Database

- A comprehensive
- Integrated
- non-redundant
- well-annotated set of reference sequences including genomic, transcript, and protein.

## BioSystems

The NCBI BioSystems Database provides integrated access to biological systems and their component

- Genes
- Proteins
- Small molecules
- Literature describing those biosystems
- Other related data throughout Entrez.

## Gene

- Gene integrates information from a wide range of species. A record may include
- Nomenclature
- Reference Sequences (RefSeqs)
- Maps
- Pathways
- Variations
- Phenotypes
- Links to genome-, phenotype-

Search NCBI databases - NCBI X National Center for Biotechnology Information

U.S. National Library of Medicine  
National Center for Biotechnology Information

Search NCBI BRCA1 Search

NCBI Databases  
Results found in 23 databases for BRCA1  
Did you mean [breast](#)?

Literature	Genes	Genetics
Bookshelf 0	Gene 57	ClinVar 57
MeSH 1	GEO DataSets 1278	dbGAP 0
NLM Catalog 1	GEO Profiles 1988	dbSNP 0
PubMed 220	HomoloGene 0	dbVar 0
PubMed Central 203	PopSet 1	GTR 1
	UniGene 1	MedGen 2
		OMIM 18

https://www.ncbi.nlm.nih.gov/entrez/term=BRCA1

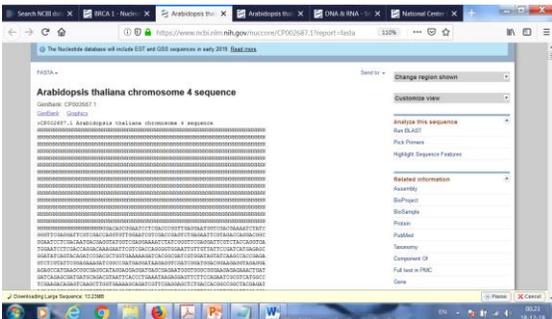
Nucleotide is selected



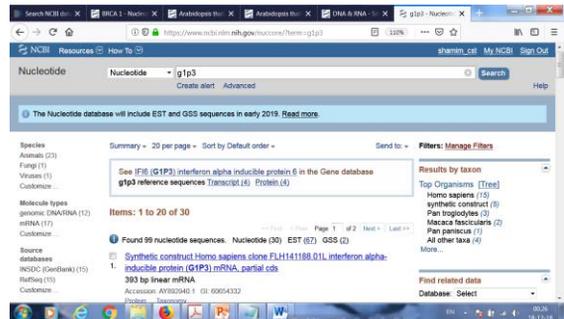
Genbank is selected



FASTA format is selected



Search for gip3 (GIP3, an interferon inducible gene)



Search NCBI | Home | Home - Gene - NCBI | Home - Biotechnology - NCBI

https://www.ncbi.nlm.nih.gov/uccore/terms.cgi?3

See **IFI6 (G1P3) interferon alpha inducible protein 6** in the Gene database  
**g1p3 reference sequences Transcript (4) Protein (4)**

Items: 1 to 20 of 30

Found 99 nucleotide sequences. Nucleotide (30) EST (67) GSS (2)

1. **Synthetic construct Homo sapiens clone FLH141188.011, interferon alpha-inducible protein (G1P3) mRNA, partial cds**  
 393 bp linear mRNA  
 Accession: AY892940.1 GI: 60654332  
[Protein](#) [Taxonomy](#)  
[GenBank](#) [FASTA](#) [Graphics](#)

Results by taxon

Top Organisms

- Homo sapiens (7)
- synthetic constru
- Macaca fascicular
- Pan paniscus (1)
- All other taxa (4)

Find related data  
 Database: Select  
 Find Items

Search NCBI | Home | Home - Gene - NCBI | Home - Biotechnology - NCBI

https://www.ncbi.nlm.nih.gov/guide/genes-expression/

NCBI National Center for Biotechnology Information

Genes & Expression

All Databases Downloads Submissions Tools How To

Quick Links

- BiProject (formerly Genome Project)
- Database of Genotypes and Phenotypes (dbGP)
- Gene
- Gene Expression Omnibus (GEO) Database
- Gene Expression Omnibus (GEO) Profiles
- Online Mendelian Inheritance in Man (OMIM)
- Protein
- Sequence Analysis

Databases

**BiProject (formerly Genome Project)**  
 A collection of genomics, functional genomics, and genetics studies and links to their resulting datasets. This resource describes project scope, materials, and objectives and provides a mechanism to retrieve datasets that are often difficult to find due to inconsistent annotation, multiple independent submissions, and the varied nature of diverse data types which are often stored in different databases.

**BiSystems**  
 Database that groups biomedical literature, small molecules, and sequence data in terms of biological relationships.

**ChIRP**  
 A resource to provide a public, tracked record of reported relationships between human variation and observed health status with supporting evidence. Related information

Search NCBI | Home | Home - Gene - NCBI | Home - Biotechnology - NCBI

https://www.ncbi.nlm.nih.gov/guide/dna-rna/

NCBI National Center for Biotechnology Information

DNA & RNA

All Databases Downloads Submissions Tools How To

Quick Links

- BiProject (formerly Genome Project)
- Database of Short Genetic Variations (dbSNP)
- GenBank
- Nucleotide Database
- PopSet
- RefSeqGene
- Reference Sequence (RefSeq)
- Sequence Read Archive (SRA)
- Trace Archive

Databases

**Assembly**  
 A database providing information on the structure of assembled genomes, assembly names and other meta-data, statistical reports, and links to genomic sequence data.

**BioCollections**  
 A curated set of metadata for culture collections, museums, herbaria and other natural history collections. The records display collection codes, information about the collection's home institutions, and links to relevant data at NCBI.

**BiProject (formerly Genome Project)**  
 A collection of genomics, functional genomics, and genetics studies and links to their resulting datasets. This resource describes project scope, materials, and objectives and provides a mechanism to retrieve datasets that are often difficult to find due to inconsistent annotation, multiple independent submissions, and the varied nature of diverse data types

Search NCBI | Home | Home - Gene - NCBI | Home - Biotechnology - NCBI

https://blast.ncbi.nlm.nih.gov/Blast.cgi

BLAST®

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance.

Web BLAST

Nucleotide BLAST  
 (nucleotide vs. nucleotide)

blastx  
 translated nucleotide vs. protein

tblastn  
 protein vs. translated nucleotide

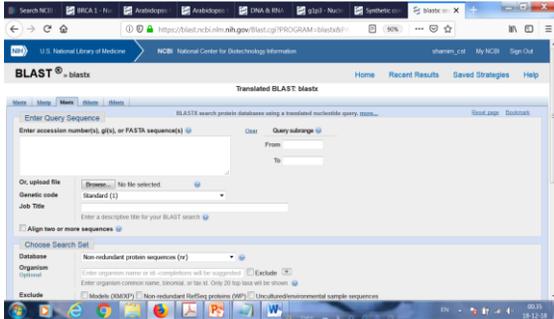
Protein BLAST  
 (protein vs. protein)

BLAST Genomes

Enter organism common name, accession name, or tax ID

Search

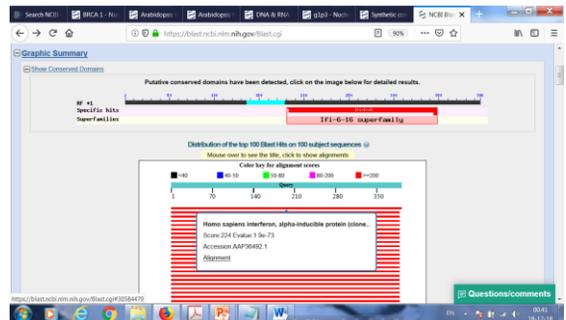
Select Blastx



Gip3 sequence found from FASTA

```

ATGCCGACAAAGCGGTA TCGC TTTTC TTG TGC TACCTGCTGC TCTTCACTTGCAG TG
GGTGGAGGCAG
GTAAAGAA AAGTGTCTCGG AGA GCTC GGACA GC GGC TCCGGG TTC TGGAA GGCCCTGA
CCTTCA TGGCCGT
CGGAGGAGACTCGCAG TCGCCGGGC TGCCCGCGC TGGGC TTCACCGCGCCGGCA TCG
CGCCAACTCG
GTGGCTGCTCGCTGATGAGC TGG TCTGCGA TCC TGA A TGGGGGCGCGGTGCCCGCCG
GGGGCTAGTGG
CCACGC TCGAGAGCCTCGGGCTGG TG GCA GCACG CTC GTCA TA GG TAA TAT TGGTGC
CTTGA TGGGCTA
CGCCACCCACAAGTA TCTCGATAGTGAGGA GGA TGA GGA GATTG
    
```





The screenshot shows the NCBI Nucleotide database homepage. At the top, there is a search bar with the text "Nucleotide" and a "Search" button. Below the search bar, there is a navigation menu with "Nucleotide" selected. A blue banner below the search bar contains the text: "The Nucleotide database will include EST and GSS sequences in early 2019. Read more." The main content area features a large graphic with DNA sequence letters (A, C, G, T) and the heading "Nucleotide". Below this graphic, a paragraph states: "The Nucleotide database is a collection of sequences from several sources, including GenBank, RefSeq, TPA and PDB. Genome, gene and transcript sequence data provide the foundation for biomedical research and discovery." At the bottom, there are three columns of links: "Using Nucleotide" (Quick Start Guide, FAQ, Help), "Nucleotide Tools" (Submit to GenBank, LinkOut, E-Utilities), and "Other Resources" (GenBank Home, RefSeq Home, Gene Home).

The screenshot shows the NCBI RefSeq database homepage. At the top, there is a search bar with the text "RefSeq" and a "Search" button. Below the search bar, there is a navigation menu with "RefSeq" selected. A large blue banner below the search bar features a DNA double helix graphic and the heading "RefSeq: NCBI Reference Sequence Database". Below the heading, a paragraph states: "A comprehensive, integrated, non-redundant, well-annotated set of reference sequences including genomic, transcript, and protein." The main content area is divided into three columns of links: "Using RefSeq" (About RefSeq, Human Reference Genome, Prokaryotic RefSeq Genomes, FAQ, NCBI Handbook, FactSheet), "RefSeq Access" (Human Genome Resources and Download, RefSeq FTP, RefSeq genomic FTP, New RefSeq genomic (last 30 days), New RefSeq transcripts (last 30 days), New RefSeq proteins (last 30 days)), and "RefSeq projects" (Consensus CDS (CCDS), RefSeq Functional Elements, RefSeqGene, Targeted Loci, Virus Variation).



