



National Center for Biotechnology

Information

By,

**Kavisa Ghosh,**

**V M.Sc.Biotechnology(Int.)**

# The National Center for Biotechnology Information



*Created in 1988 as a part of the  
National Library of Medicine at NIH*

- Establish public databases
- Research in computational biology
- Develop software tools for sequence analysis

# NCBI HOME PAGE

Address <http://www.ncbi.nlm.nih.gov/>

Search Site Rating News Links Customize Links Free Hotmail

softonic NCBI SEARCH WEB GADGETS Games Play! MTV Free TV Login zynga Gar

NCBI Resources How To

**NCBI** National Center for Biotechnology Information

**Resources**

- NCBI Home**
- All Resources (A-Z)
- Data & Software
- DNA & RNA
- Domains & Structures
- Genes & Expression
- Genetics & Medicine
- Genomes & Maps
- Homology
- Literature
- Proteins
- Sequence Analysis
- Small Molecules
- Taxonomy

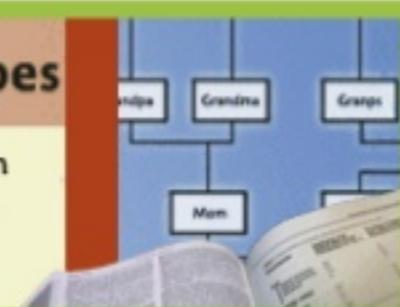
**Welcome to NCBI**

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[More about the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [RSS](#)

**Genotypes and Phenotypes**

Data from Genome Wide Association studies that link genes and diseases. See study variables, protocols, and analysis.



II 1 2 3 4

**How To...**

- Determine conserved synteny between the genomes of two organisms

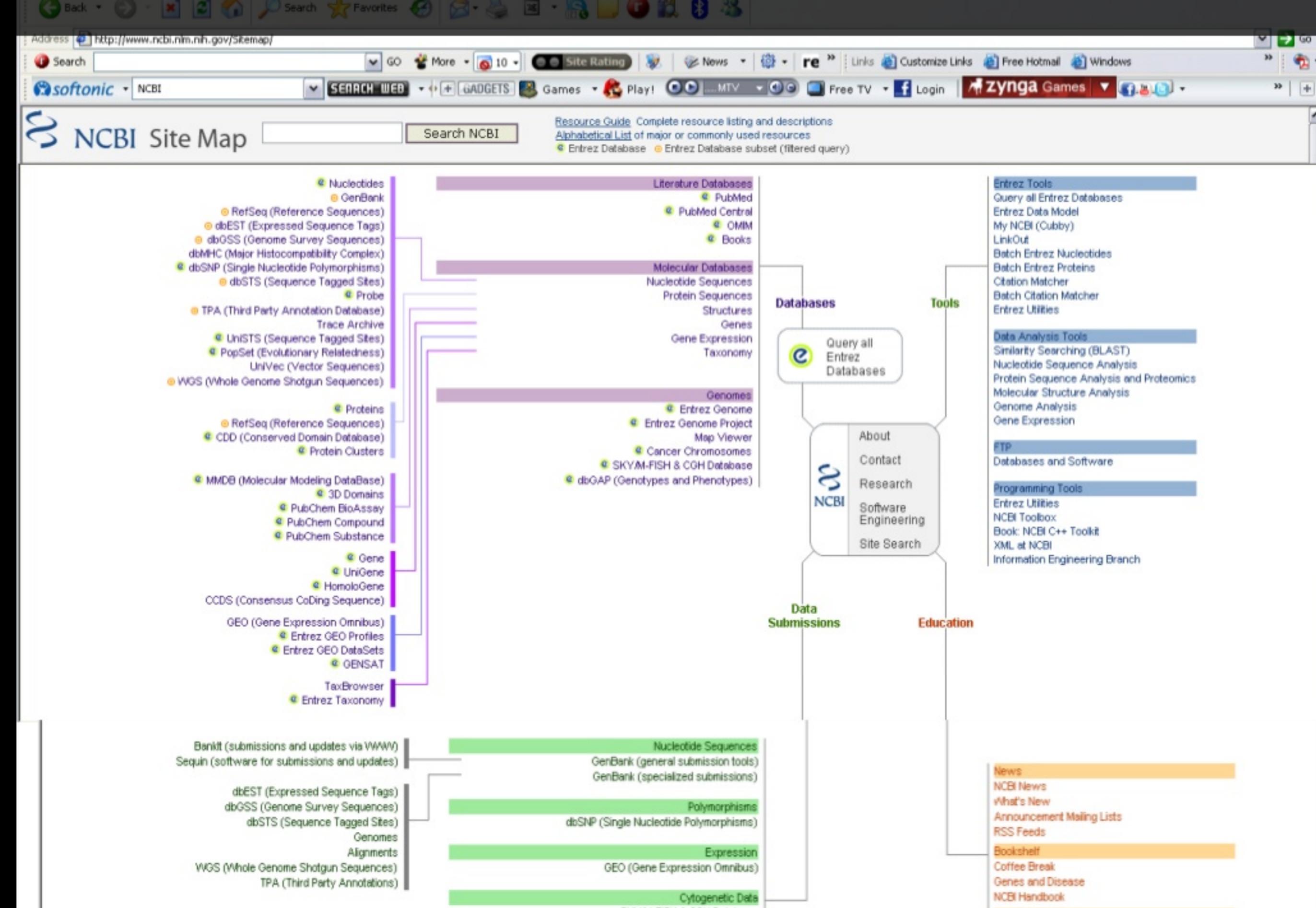
**Popular Resources**

- BLAST
- Bookshelf
- Gene
- Genome
- Nucleotide
- OMIM
- Protein
- PubChem
- PubMed
- PubMed Central
- SNP

**NCBI News**

Selected Structures, Taxonomy on Wikipedia 23 Jul 2010

The June NCBI News is available on the Bookshelf.





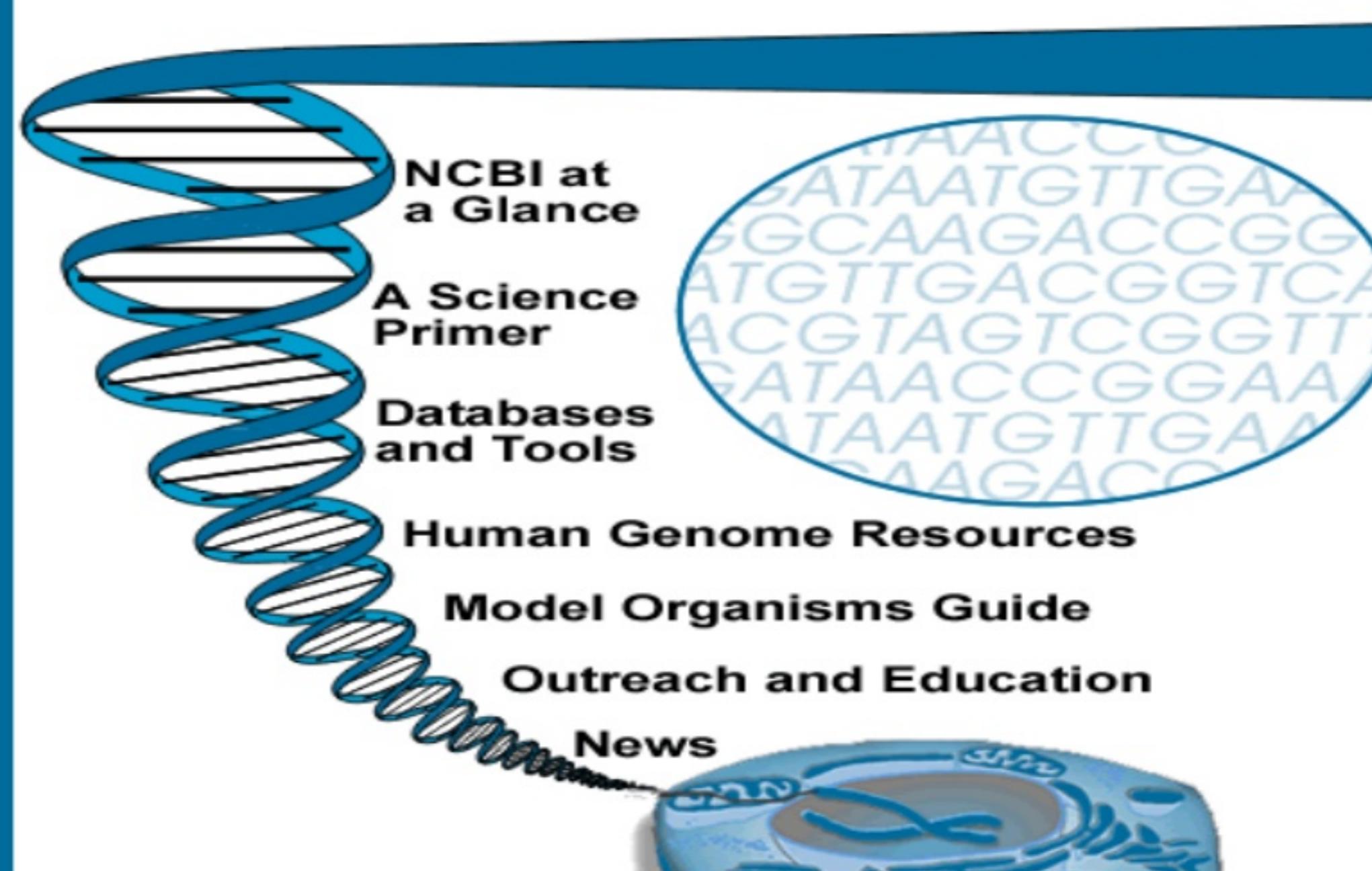
# ACGTACGATCGTAGCTAGC TACGTACGCTAGCTAGCTAGC TGCCTACGCTAGCTAGCTAGC

# About NCBI

[About NCBI](#)  
[Site Map](#)

National Center for Biotechnology Information

<a href="#">About NCBI</a>	<a href="#">NCBI at a Glance</a>	<a href="#">A Science Primer</a>	<a href="#">Databases and Tools</a>
<a href="#">Human Genome Resources</a>	<a href="#">Model Organisms Guide</a>	<a href="#">Outreach and Education</a>	<a href="#">News</a>



# Entrez: An Integrated Database Search and Retrieval System

The logo for Entrez, featuring the word "Entrez" in large orange letters with a blue shadow, and "search and retrieval system" in smaller blue text below it.

Entrez is a retrieval system for searching several linked databases. It provides access to:

- [PubMed](#): The biomedical literature (PubMed)
- [Nucleotide](#): sequence database (GenBank)
- [Protein](#): sequence database
- [Structure](#): three-dimensional macromolecular structures
- [Genome](#): complete genome assemblies
- [PopSet](#): population study data sets
- [OMIM](#): Online Mendelian Inheritance in Man
- [Taxonomy](#): organisms in GenBank
- [Books](#): BookShelf online books
- [3D Domains](#): domains from Entrez Structure
- [UniSTS](#): markers and mapping data
- [SNP](#): single nucleotide polymorphisms
- [CDD](#): conserved domains
- [Journals](#): journals in Entrez
- [UniGene](#): gene-oriented clusters of transcript sequences
- [PMC](#): full-text digital archive of life sciences journal literature





HOME | SEARCH | SITE MAP

PubMed

Entrez

Human Genome

GenBank

Map Viewer

BLAST

## Search across databases

all[filter]

GO

CLEAR

Help

14764283



**PubMed:** biomedical literature citations and abstracts

255107



**PubMed Central:** free, full text journal articles

44433



**Books:** online books

15831



**OMIM:** Online Mendelian Inheritance in Man

none



**Site Search:** NCBI web and FTP sites

37871567



**Nucleotide:** sequence database (GenBank)

4609013



**Protein:** sequence database

3343



**Genome:** whole genome sequences

24167



**Structure:** three-dimensional macromolecular structures

212265



**Taxonomy:** organisms in GenBank

10699891



**SNP:** single nucleotide polymorphism

805506



**Gene:** gene-centered information

none



**HomoloGene:** Eukaryotic homology groups

978830



**UniGene:** gene-oriented clusters of transcript sequences

18039



**CDD:** conserved protein domain database

99783



**3D Domains:** domains from Entrez Structure

257454



**UniSTS:** markers and mapping data

21947



**PopSet:** population study data sets

none



**GEO:** expression and molecular abundance profiles

none



**GEO DataSets:** experimental sets of GEO data

19940



**Journals:** detailed information about the journals indexed in

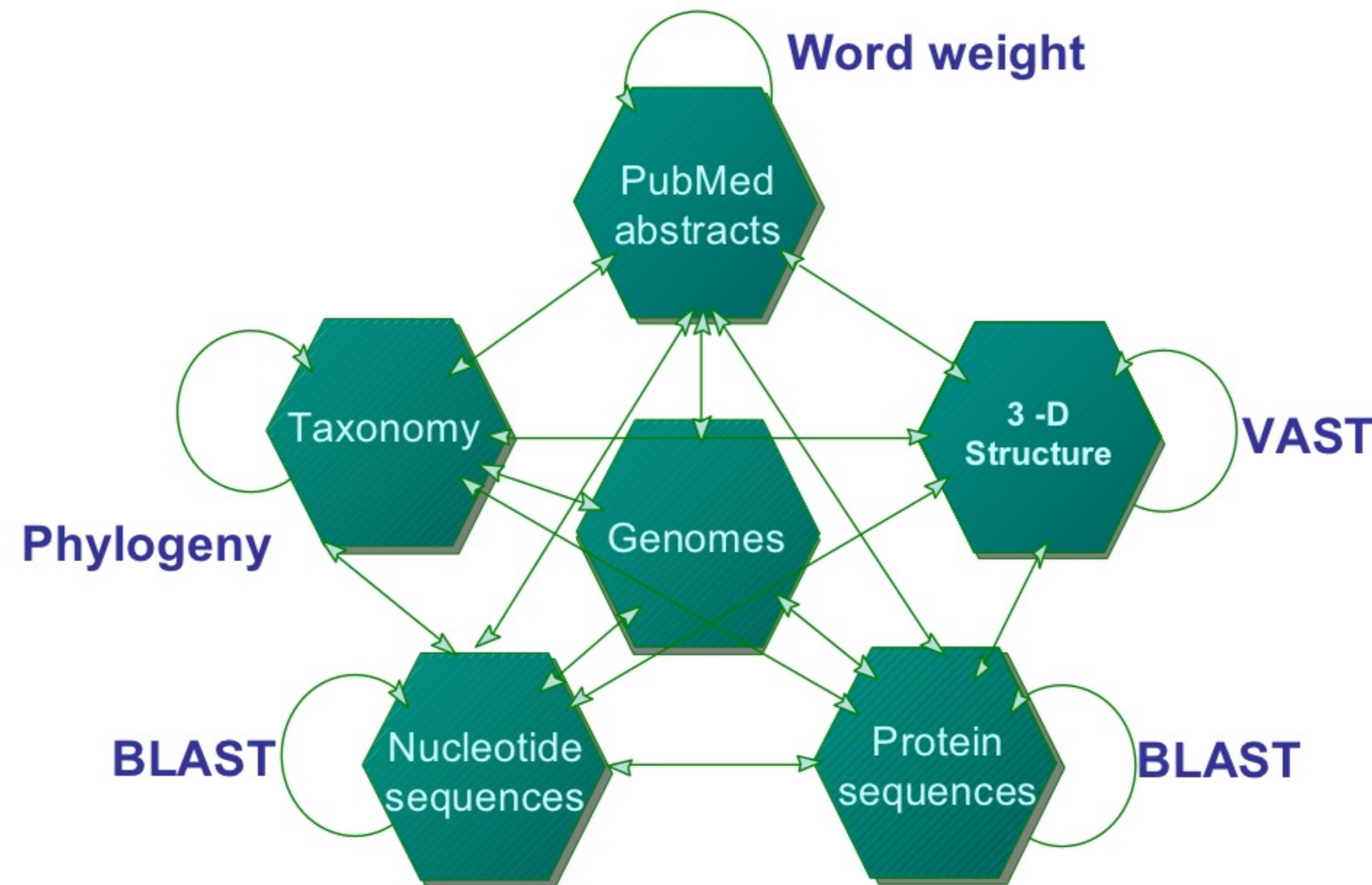
23031



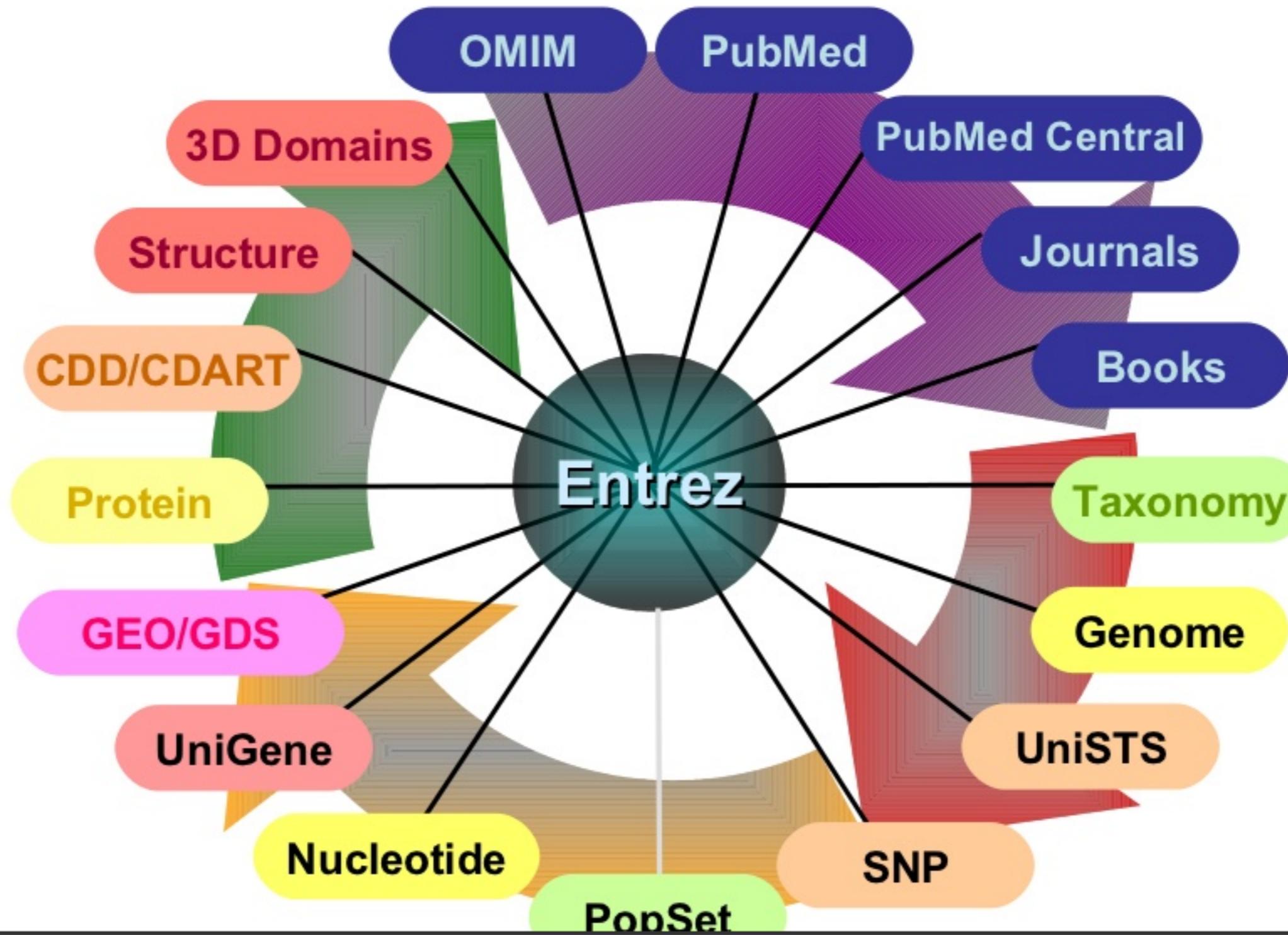
**MeSH:** detailed information about NLM's controlled



# Entrez: Database Integration



# The (ever expanding) Entrez System





# Entrez Databases

- All Molecular Database entries are organized by organism (*Taxonomy Database*).
- Each record is assigned a UID.
  - A “unique integer identifier” for internal tracking
- Each record is indexed by data fields.
  - [author], [title], [organism], and many others
- Each record is given a Document Summary.
  - a summary of the record’s content (DocSum)
- Each record is manually or computationally assigned [links](#) to biologically related UIDs in and across databases.

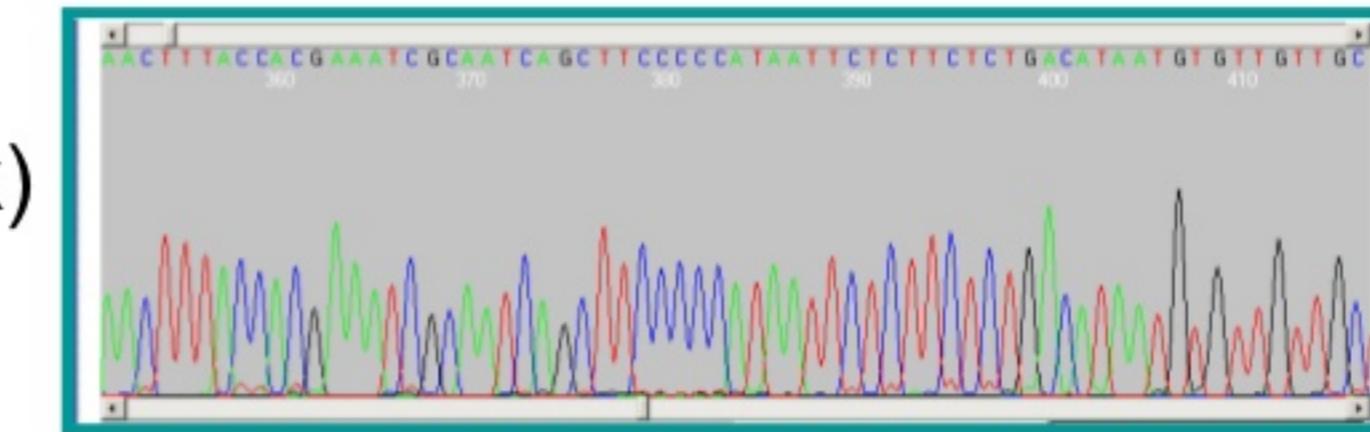
# Literature Databases

- PubMed
- Books
- PubMed Central
- Journals
- On-Line Mendelian Inheritance in Man  
(OMIM)



# Molecular Sequence Databases

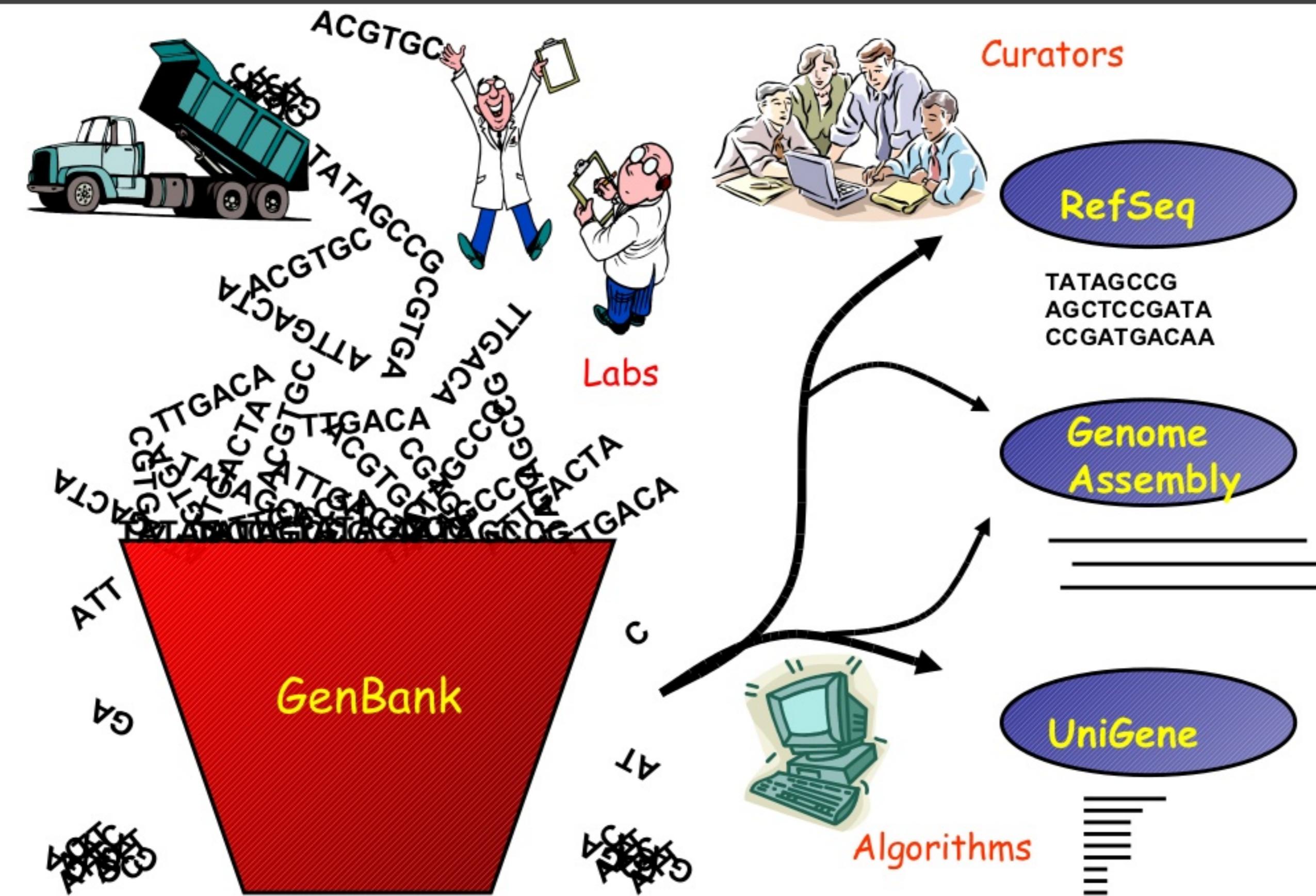
- Sequence Databases
  - Nucleotide (GenBank)
    - Taxonomy
    - PopSet
  - Protein
- Marker Databases
  - Single Nucleotide Polymorphisms (SNP's, dbSNP)
  - Sequence Tagged Sites (STS's, dbSTS)
  - Expressed Sequence Tags (EST's, dbEST)
    - UniGene



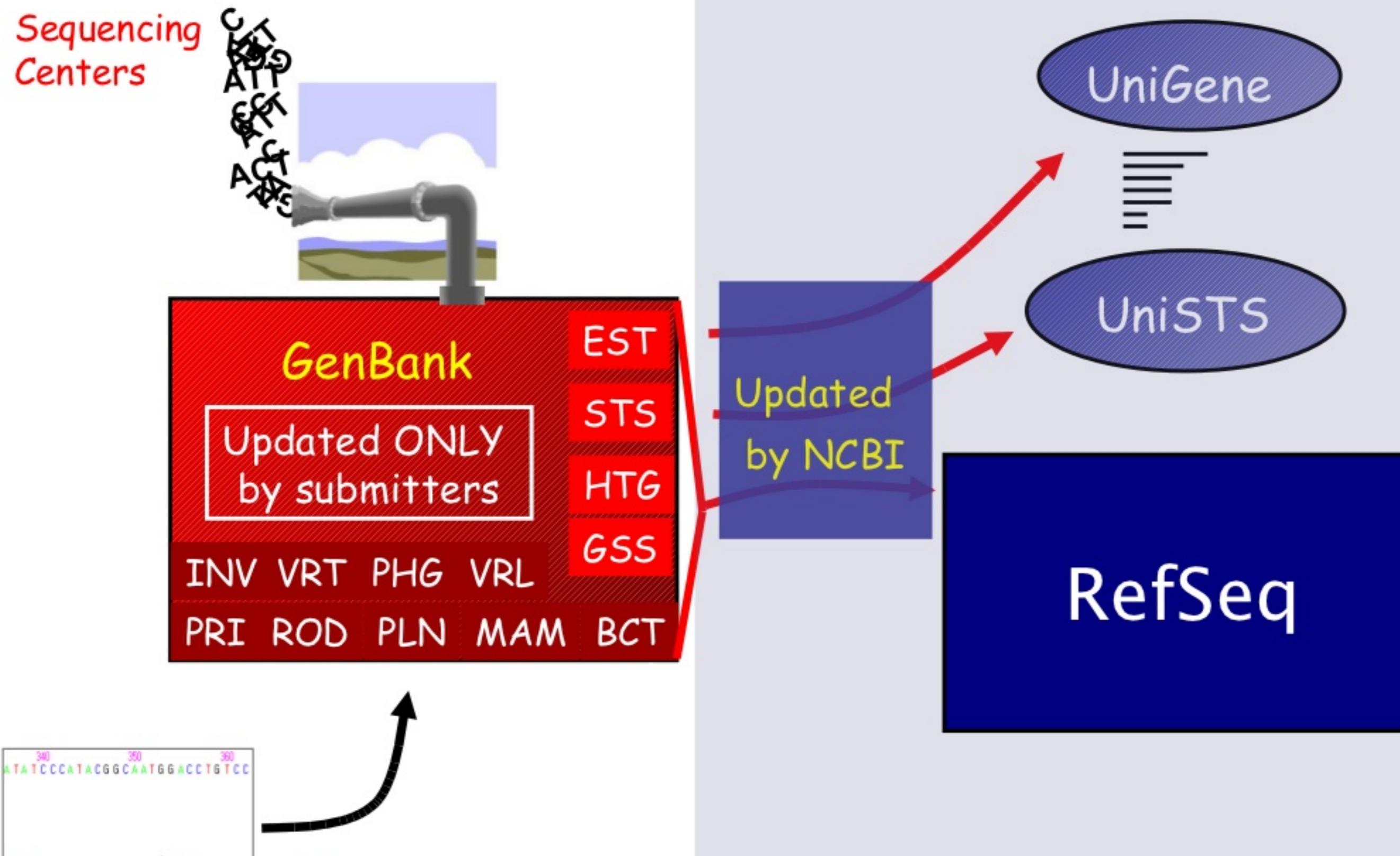
# Molecular Databases

- Primary Databases
  - Original submissions by experimentalists
  - Database staff organize but don't add additional information
    - **Example: GenBank**
- Derivative Databases
  - Human curated
    - compilation and correction of data
    - **Example: SWISS-PROT, NCBI RefSeq mRNA**
  - Computationally Derived
    - **Example: UniGene**
  - Combinations
    - **Example: NCBI Genome Assembly**





# Derivative Databases



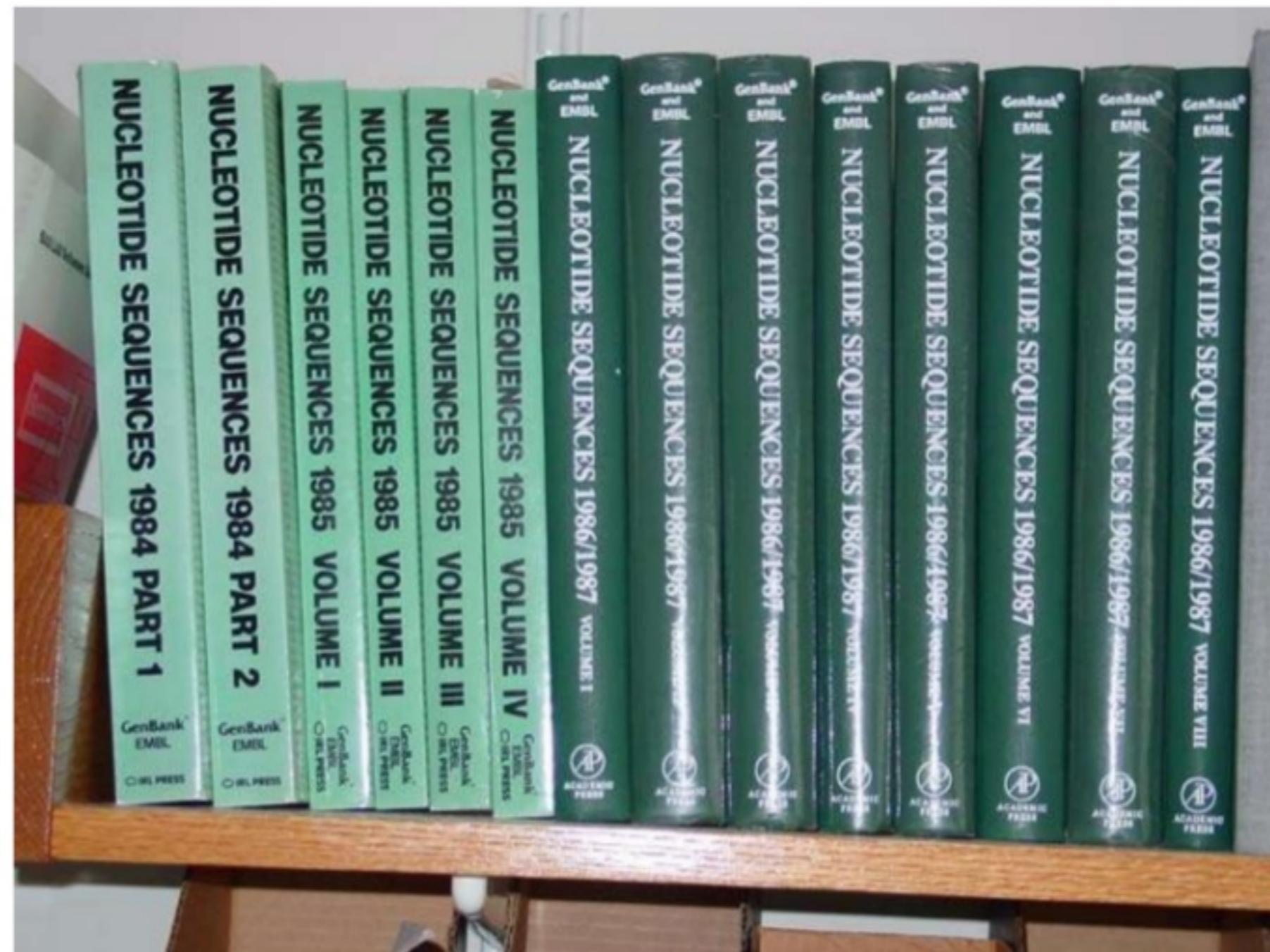
# What is

## GenBank?

- Nucleotide only sequence database
- Archival in nature
  - Historical
  - Reflective of submitter point of view (subjective)
  - Redundant
- GenBank Data
  - Direct submissions (traditional records)
  - Batch submissions (EST, GSS, STS)
  - ftp accounts (genome data)
- Three collaborating databases
  - GenBank
  - DNA Database of Japan (DDBJ)
  - European Molecular Biology Laboratory (EMBL) Database



# The Old Way



# GenBank: NCBI's Primary Sequence Database

Release 136

June 2003

25,592,865

18,197,119 (June 2002)

32,528,249,295

22,616,937,182 (June 2002)

110,000 +

Records

Nucleotides

Species

- full release every two months
- incremental and cumulative updates daily
- available only through internet

<ftp://ftp.ncbi.nih.gov/genbank/>

<ftp://genbank.sdsc.edu/pub>

<ftp://bio-mirror.net/biomirror/genbank/>

## GenBank Continued...

- 106,533,156,756 bases in 108,431,692 sequence records in the traditional GenBank divisions
- 148,165,117,763 bases in 48,443,067 sequence records in the WGS division as of August 2009.



# GenBank Divisions

<b>PRI</b> (28)	Primate
<b>ROD</b> (15)	Rodent
<b>PLN</b> (20)	Plant and Fungal
<b>BCT</b> (18)	Bacterial/Archeal
<b>INV</b> (7)	Invertebrate
<b>VRT</b> (7)	Other Vertebrate
<b>VRL</b> (4)	Viral
<b>MAM</b> (2)	Mammalian
<b>PHG</b> (1)	Phage
<b>SYN</b> (1)	Synthetic
<b>ENV</b> (4)	Envir. samples
<b>UNA</b> (1)	Unannotated

<b>EST</b> (570)	Expressed Sequence Tag
<b>GSS</b> (197)	Genome Survey Sequence
<b>HTG</b> (88)	High Throughput Genomic
<b>PAT</b> (27)	Patent
<b>STS</b> (9)	Sequence Tagged Site

## “Organismal” (Traditional)

- Organized by taxonomy (sort of)
- Direct submissions (Sequin/Bankit)
- Accurate (~1 error per 10,000 bp)
- **Well characterized**

## “Functional” (Bulk)

- Organized by sequence type
- Batch submissions (ftp/email)
- Less accurate



# GenBank Functional (Bulk) Divisions

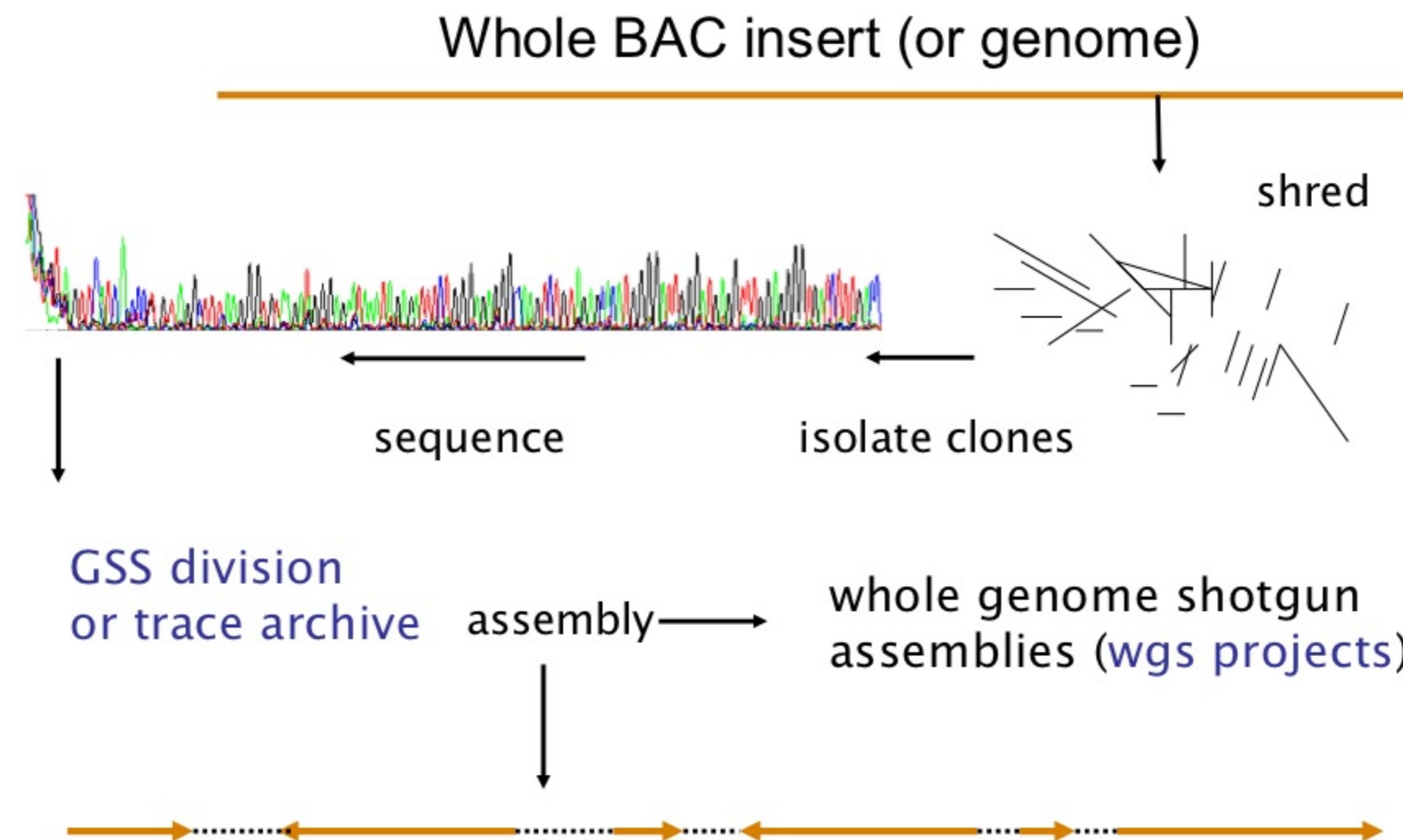


- **Expressed Sequence Tag**
  - 1st pass single read cDNA
- **Genome Survey Sequence**
  - 1st pass single read gDNA
- **High Throughput Genomic**
  - incomplete sequences of genomic clones
- **Sequence Tagged Site**
  - PCR-based mapping reagents

**Whole Genome Shotgun**



# GSS, HTG, WGS



# Whole Genome Shotgun Projects

- 685 projects
  - Bacteria (320)
  - Environmental sequences (14)
  - Archaea (8)
  - Eukaryotes (140), including:
    - Chicken, Rat, Mouse, Dog (2), Chimpanzee, Human
    - Pufferfish (2)
    - Honeybee, Anopheles, Fruit Flies (3), Silkworm
    - Nematode (2)
    - Yeasts (8), Aspergillus (2)
    - Rice (2)



# Whole Genome Shotgun (WGS) Projects

23: [AACC00000000](#)

Homo sapiens chromosome 7, whole genome shotgun sequencing project  
gi|50364594|gb|AACC00000000.2|AACC02000000[50364594]

Links

24: [AACQ00000000](#)

Candida albicans SC5314, whole genome shotgun sequencing project  
gi|46445633|gb|AACQ00000000.1|AACQ01000000[46445633]

wgs master [properties]

Links

25: [AABT00000000](#)

Aspergillus terreus ATCC 20542, whole genome shotgun sequencing project  
gi|27262064|gb|AABT00000000.1|AABT01000000[27262064]

**ftp://ftp.ncbi.nih.gov/genbank/wgs/**

26: [NZ\\_AAEZ00000000](#)

Pseudomonas syringae pv. phaseolicola 1448A, unfinished sequence, whole genome shotgun sequencing project  
gi|50591998|ref|NZ\_AAEZ00000000.1|NZ\_AAEZ01000000[50591998]

Links

27: [NZ\\_AAFA00000000](#)

Streptococcus suis 89/1591, unfinished sequence, whole genome shotgun sequencing project  
gi|50591969|ref|NZ\_AAFA00000000.1|NZ\_AAFA01000000[50591969]

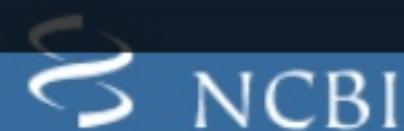
Links

28: [AAFA00000000](#)

Streptococcus suis 89/1591, whole genome shotgun sequencing project  
gi|50557642|gb|AAFA00000000.1|AAFA01000000[50557642]

Links





ENTREZ

# Genome Project

connection

information  
discovery

Well

All Databases

PubMed

Nucleotide

Protein

Genome

Structure

PMC

1

Search

Genome Project

for

Go

Clear

Limits

Preview/Index

History

Clipboard

Details

About Entrez

Entrez Genome  
ProjectHome  
Overview

Help

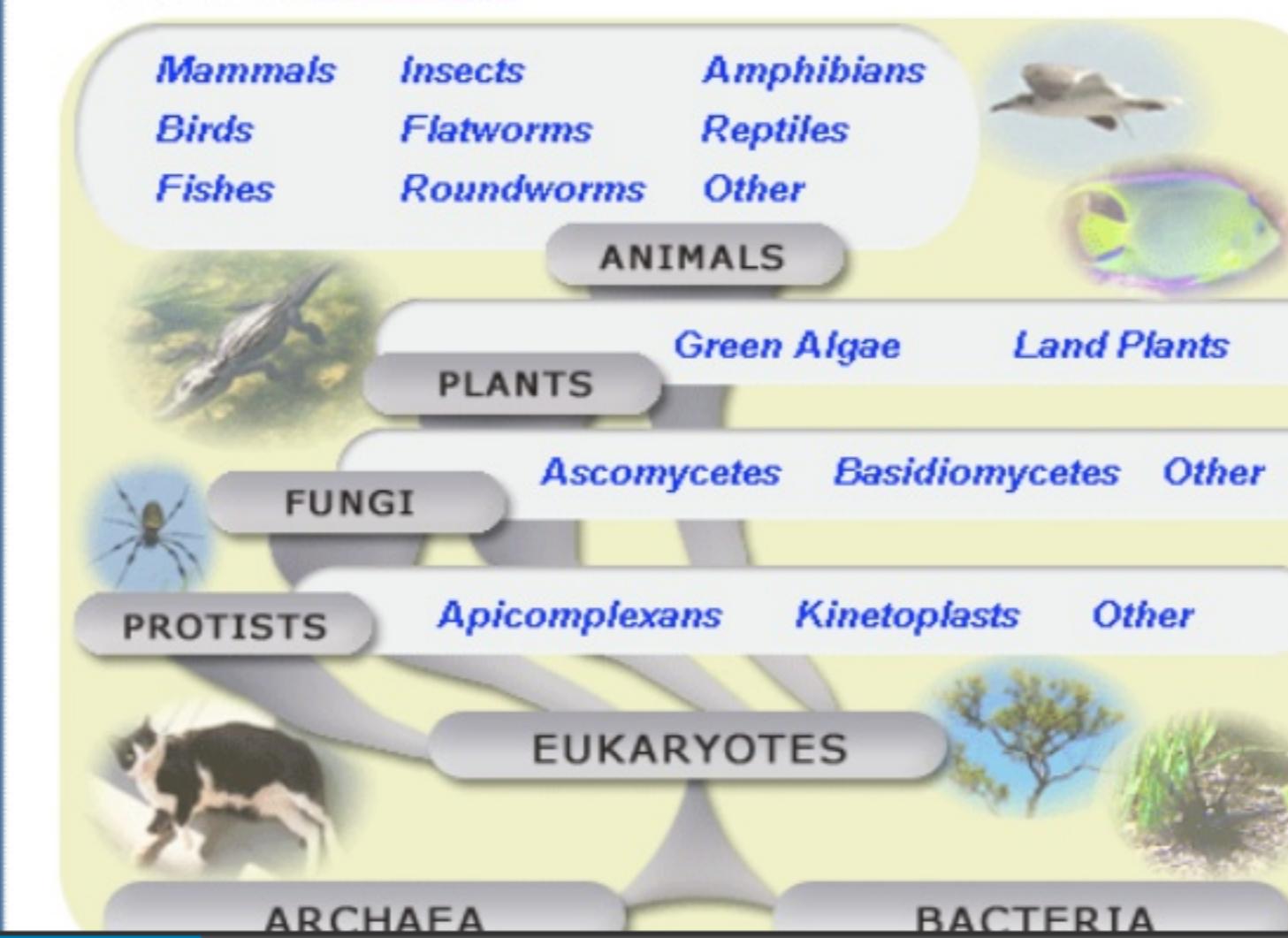
Statistics

Sequencing  
CentersSubmitting  
Project  
Submissions  
Project InstructionsGeneral Genome  
Submissions  
Feature Tables  
Bacterial Genome  
Submissions  
Whole Genome  
Shotgun  
Sequences

Related

Welcome to the NCBI Entrez Genome Project database.

This searchable database is a collection of complete and incomplete large-scale sequencing, assembly, annotation, and mapping projects for cellular organisms. The database is organized into organism-specific overviews that function as portals from which all projects in the database pertaining to that organism can be browsed and retrieved. [Read more...](#)

**NCBI Resources**[Entrez Gene](#)gene-related  
information[Entrez Genome](#)sequence and map data  
from whole genomes[Eukaryotic Projects](#)eukaryotic-specific  
genome projects[Genomic Biology](#)organism-specific  
links[Prokaryotic Projects](#)prokaryotic-specific  
genome projects[Organellar Genomes](#)organellar reference sequences  
and tools[Plant Genomes](#)major plant  
genome projects[RefSeq](#)the reference  
sequence project[Viral Genomes](#)viral reference sequences  
and tools[WGS Sequences](#)

# What is UniGene?

A gene-oriented view of sequence entries

- Megablast based automated sequence clustering
- Now informed by genome hits **New!**
- Nonredundant set of gene oriented clusters
- Each cluster a unique gene
- Information on tissue types and map locations
- Includes well-characterized genes and novel ESTs
- Useful for gene discovery and selection of mapping reagents



**UniGene**  
ORGANIZED VIEW OF THE TRANSCRIPTOME

PubMed Nucleotide Protein Genome Structure Popset Taxonomy

Search UniGene ▾ Go

Limits Preview/Index History Clipboard Details

UniGene Dr.12379 *Danio rerio* il17rd

Interleukin 17 receptor D (il17rd)

**SELECTED PROTEIN SIMILARITIES**

*Comparison of sequences in UniGene with proteins supported by a complete genome. The alignments can suggest function of a gene.*

<i>H. sapiens</i>	<a href="#">pir:T42695</a> - T42695 hypothetical protein DKFZp434N1928.1 - human	55.56 % / 553 aa (see <a href="#">ProtEST</a> )
<i>M. musculus</i>	<a href="#">ref:NP_032385.1</a> - interleukin 17 receptor [Mus musculus]	26.71 % / 269 aa (see <a href="#">ProtEST</a> )

**GENE EXPRESSION**

*Tissues and development stages from this gene's sequences survey gene expression. Links to other NCBI expression resources.*

cDNA sources: kidney, olfactory rosettes, ovary, testis, gastrula, larval, adult

Restricted Expression: gastrula [[Show more like this](#)]

Expression Profile: View expression levels using UniGene's EST ProfileViewer

ZFIN: Gene Expression provided by the Zebrafish Information Network

**MAPPING POSITION**

*Genomic location specified by transcript mapping, radiation hybrid mapping, genetic mapping or cytogenetic mapping.*

**Links** ▾ **Links** ▾

- ▶ Gene
- ▶ HomoloGene
- ▶ Nucleotide
- ▶ OMIM
- ▶ Protein
- ▶ Pubmed
- ▶ Unigene\_homologous
- ▶ ZFIN

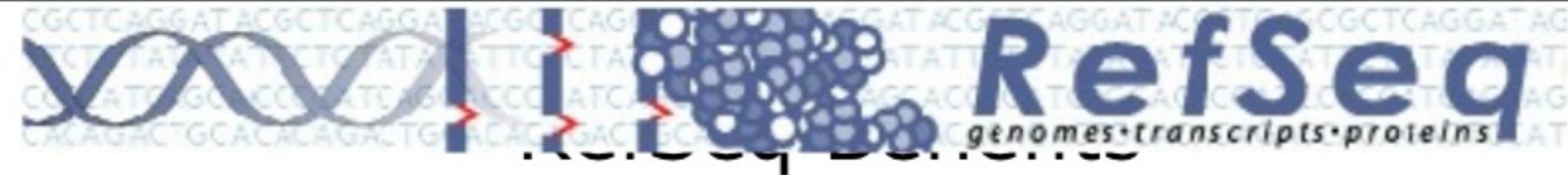


Species	UniGene Entries
Chordata	
Mammalia	
<i>Bos taurus</i>	23,572
<i>Canis familiaris</i>	3,905
<i>Homo sapiens</i>	105,651
<i>Mus musculus</i>	79,846
<i>Rattus norvegicus</i>	39,874
<i>Sus scrofa</i>	20,426
Aves	
<i>Gallus gallus</i>	12,220
Amphibia	
<i>Xenopus laevis</i>	22,129
<i>Xenopus tropicalis</i>	17,102
Actinopterygii	
<i>Danio rerio</i>	17,951
<i>Oncorhynchus mykiss</i>	13,458
<i>Oryzias latipes</i>	8,127
<i>Salmo salar</i>	1,033
Asciidae	
<i>Ciona intestinalis</i>	13,600
Echinodermata	
Echinoidea	
<i>Strongylocentrotus purpuratus</i>	2,592
Arthropoda	
Insecta	
<i>Anopheles gambiae</i>	14,085
<i>Apis mellifera</i>	4,756
<i>Bombyx mori</i>	2,008
<i>Drosophila melanogaster</i>	12,322
Nematoda	
Chromadorea	
<i>Caenorhabditis elegans</i>	15,655
Platyhelminthes	

# UniGene

Embryophyta	
Bryopsida	
<i>Physcomitrella patens</i>	6,946
Coniferopsida	
<i>Pinus taeda</i>	8,669
Eudicotyledons	
<i>Arabidopsis thaliana</i>	19,666
<i>Glycine max</i>	12,567
<i>Helianthus annuus</i>	1,897
<i>Lactuca sativa</i>	9,629
<i>Lycopersicon esculentum</i>	3,658
<i>Medicago truncatula</i>	5,398
<i>Populus tremula x Populus tremuloides</i>	2,711
<i>Solanum tuberosum</i>	5,476
<i>Vitis vinifera</i>	11,832
Liliopsida	
<i>Hordeum vulgare</i>	11,980
<i>Oryza sativa</i>	25,089
<i>Saccharum officinarum</i>	4,771
<i>Sorghum bicolor</i>	7,075
<i>Triticum aestivum</i>	24,543
<i>Zea mays</i>	12,941
Chlorophyta	
Chlorophyceae	
<i>Chlamydomonas reinhardtii</i>	5,874
Mycetozoa	
Dictyosteliida	
<i>Dictyostelium discoideum</i>	3,853
Apicomplexa	
Coccidia	
<i>Toxoplasma gondii</i>	6,001





- non-redundant; best representative
- updates to reflect current sequence data and biology
- distinct, stable accession series

# RefSeq: NCBI's Derivative Sequence Database

- **Curated transcripts and proteins**
  - reviewed
  - human, mouse, rat, fruit fly, zebrafish, arabidopsis
- **Model transcripts and proteins**
- **Assembled Genomic Regions (contigs)**
  - human genome
  - mouse genome
- **Chromosome records**
  - Human genome
  - microbial
  - organelle

`srcdb_refseq[Properties]`

`ftp://ftp.ncbi.nih.gov/refseq/release/`



# RefSeq Benefits

- non-redundancy
- explicitly linked nucleotide and protein sequences
- updates to reflect current sequence data and biology
- data validation
- format consistency
- distinct accession series
- stewardship by NCBI staff and collaborators



# RefSeq Accession Numbers

**Curated mRNA**

**Curated Protein**

**Curated non-coding RNA**

**Predicted mRNA**

**Predicted Protein**

**Predicted non-coding RNA**

**Reference Genomic Sequence**

**Microbial replicons, organelle**

**genome**

**Contig**

**WGS Supercontig**



# Third Party Annotation (TPA) Database

- Annotations of [existing](#) GenBank sequences
- Allows for community annotation of genomes
- Direct submissions
  - BankIt
  - Sequin

**tpa [Properties]**



# Other NCBI Databases

- **dbSNP:** nucleotide polymorphism
- **Geo:** Gene Expression Omnibus  
microarray and other expression data
- **Gene:** gene records  
Unifies LocusLink and Microbial Genomes
- **Structure:** imported structures (PDB)  
Cn3D viewer, NCBI curation
- **CDD:** conserved domain database  
Protein families (COGs and KOGs)  
Single domains (PFAM, SMART, CD)

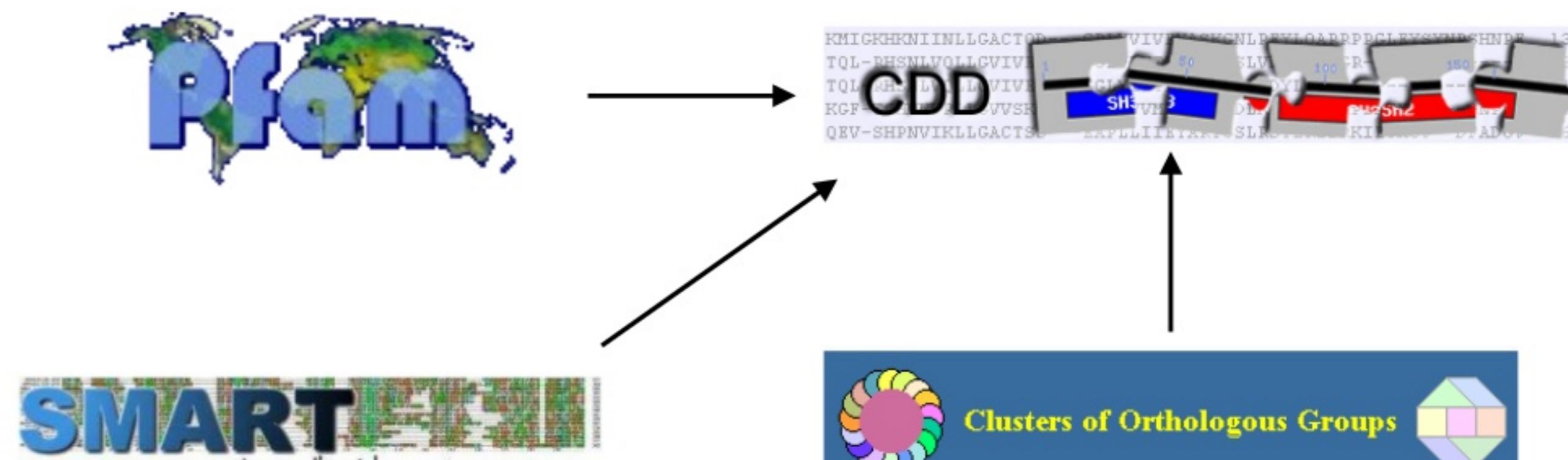
# NCBI Protein Databases

- GenPept GenBank, EMBL, DDBJ CDS translations
- RefSeq mRNA based (NP\_) and genome based (XP\_)
- Swiss-Prot curated high quality protein reviews
- PIR protein information resource Georgetown University
- PRF protein resource foundation
- PDB Protein Databank sequences from structures

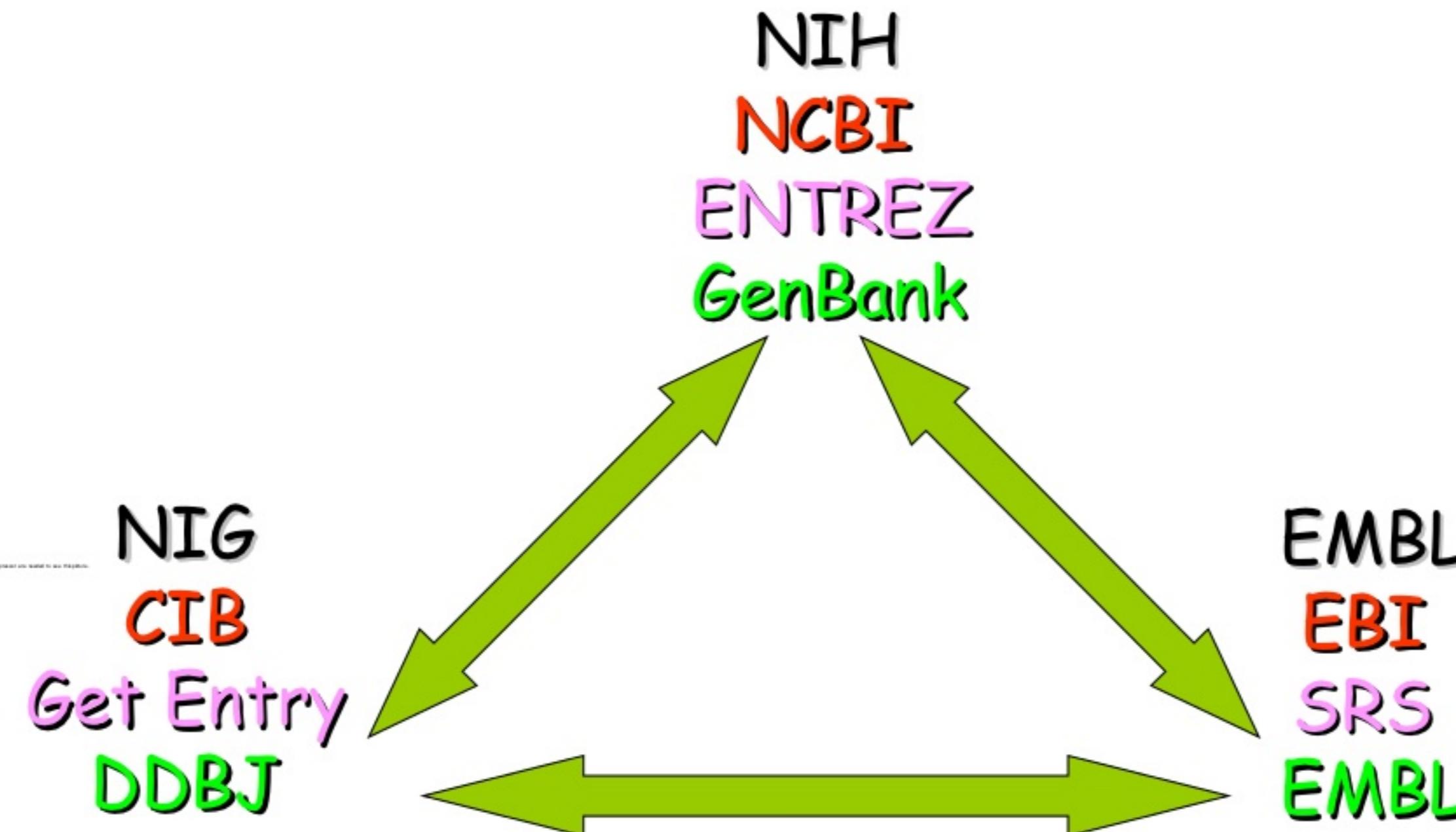




# NCBI Structures and Domains



# The International Nucleotide Sequence Database Collaboration



## Sequence formats

ASN.1  
DNAStrider  
EMBL  
Fitch  
GCG  
[GenBank](#)/GB  
IG/Stanford  
MSF  
NBRF  
Olsen  
PAUP/NEXUS  
Pearson/[Fasta](#)  
Phylip  
PIR/CODATA  
Plain/Raw  
Pretty  
Zuker

**NOTE:**

- FASTA is a popular sequence format



**NCBI Nucleotide**

PubMed Nucleotide Protein Genome Structure PopSet Taxonomy

Search **Nucleotide** for  Go Clear

Limits Preview/Index History Clipboard

Display **GenBank** Save Text Add to Clipboard Get Subsequence

1: AF467571. *Hydrobia acuta* ac...[gi:22416476]

LOCUS AF467571 638 bp DNA linear INV 22-AUG-2002

DEFINITION *Hydrobia acuta* acuta isolate 1479 cytochrome c oxidase subunit I (COI) gene, partial cds; mitochondrial gene for mitochondrial product.

ACCESSION AF467571

VERSION AF467571.1 GI:22416476

KEYWORDS .

SOURCE *Hydrobia acuta* acuta.

ORGANISM Mitochondrion *Hydrobia acuta* acuta  
Eukaryota; Metazoa; Mollusca; Gastropoda; Caenogastropoda;  
Mesogastropoda; Rissocoidea; Hydrobiidae; Hydrobia.

REFERENCE 1 (bases 1 to 638)

AUTHORS Wilke,T. and Pfenniger,M.

TITLE Separating historic events from recurrent processes in cryptic species: phylogeography of mud snails (*Hydrobia* spp.)

JOURNAL Mol. Ecol. 11 (8), 1439-1451 (2002)

PUBMED [12144664](#)

REFERENCE 2 (bases 1 to 638)

AUTHORS Wilke,T.

TITLE Direct Submission

JOURNAL Submitted (11-JAN-2002) Department of Microbiology and Tropical Medicine, The George Washington University, 2300 Eye Street, Washington, DC 20037, USA

FEATURES Location/Qualifiers

source 1..638  
 /organism="Hydrobia acuta acuta"  
 /organelle="mitochondrion"  
 /isolate="1479"  
 /sub\_species="acuta"

**GenBank format**



□ 1: AF467571. *Hydrobia acuta* ac...[gi:22416476]

>gi|22416476|gb|AF467571.1| Hydrobia acuta acuta isolate 1479 cytochrome c oxidase subunit I (COI) gene, parti  
ATTTTATTGGTATGTGGCTGGGTTAGTAGGTACAGCACTAACGTTAATTCGTGCTGAACCTAGGTC  
AGCCTGGTGCCTTGGGTGATGATCAGCTTATAACGTAATTGTTACTGCTCATGCCTTGTATAAT  
TTTTTTCTTGTAAATGCCTATAATAATTGGTGGCTTGGAAATTGATTAGTGCCTTAATACTTGGTGCT  
CCAGATATAAGCTTTCCTCGGCTTAATAACATAAGTTCTGACTTTACCTCCTGCTTGCTATTATTAC  
TTTCTTCGGCAGCTGTAGAGAGAGGAGCAGGGACAGGATGAACCGTGTATCCCCCATTATCTAGTAACAT  
TGCTCACGCCGGGGGGCTGTAGATTAGCTATTTCTCTCCACTTAGCGGGTGGTCTTCTATTCTT  
GGGGCTGTAAATTTATTACAACATCATTAAATAACGGTGACGAGGAATGCAGTTGAGCGGGCTTCCGT  
TGTCGTATGATCTGTAAAAATTACTGCCATTCTATTACTATCTTACCTGTCTAGCTGGTGCTAT  
TACTATGCTTTAACGGATCGAAATTAAATACTGCATTTCGACCCAGCAGGAGGTGGAGACCCTATT  
TTATACCA

Revised: July 5, 2002.

[Disclaimer](#) | [Write to the Help Desk](#)  
NCBI | NLM | NIH

## Fasta format

# FASTA Format

```
>gi|603218|gb|U18238.1|MSU18238 Medicago sativa glucose-6-phosphate dehyd  
CCACCAAGATATAATTAAAGTAGATCAGAGTAGAAGAAGATGGGAACAAAATGAATGGCATGTAGAAAGAAGA  
GATACCCATACCTTACTTCATTCCTTACCAACACCTTACTTCAGACTCCCCACTCTCTTATTCCTTAC
```

## FASTA Definition Line

**>gi|603218|gb|U18238.1|MSU18238**

```
CCTTCAGTGTACATCCGTTGCAAATGATCAAAACTTGTGATGAATAAAATCTGATTTGGTGGAT  
GTTGAGAAACCGTAGGGATCTAGAATCTAGAAGAACTCAGTACGTATTGATCACTATTTAAGGAACAA  
TAAAGAACCCACAATCGTATTGATCACTATTTAAGGAACAAATCTGCTCTGTGGAACCAACATTGAAATGGL  
ATGTTAGTACTTCGTTTGCAAAAGATAGTATTAGAGAGAGATATCATTGAAATGGAATTATCCG  
AGATATCATTCCAAACCACTGAGACATTGAGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
AGATATCATTCCAAACCACTGAGACATTGAGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
CCTGAGCACATTGAGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
TTCTTGGACAATATGAAATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
AACTACTATTCTGGGGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
CTAAATTCTAGGAAGGCCATTGAGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
AGCAAGGGAGAAACGAGTACATTGAGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
GCAACCTGGACTGGAAATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
ATAACCATTCCAGAGGCTTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
GCAGAGACGAATTAAAGCTTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
GAAGCCGGTTCTTACAAATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
TATGTTCAAACACCCGGATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
AGGATTAGGATTATCAGCTTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC  
TCATTTGGCTCTATAATTTGGACAAATATGAAATTTGGACAAACCCGACTTTGC
```

gi number

Locus Name

## Database Identifiers

gb GenBank

emb EMBL

dbj DDBJ

sp SWISS-PROT

pdb Protein Databank

pir PIR

prf PRF

ref RefSeq

## Accession number

# Data Analysis Tools



# BLAST

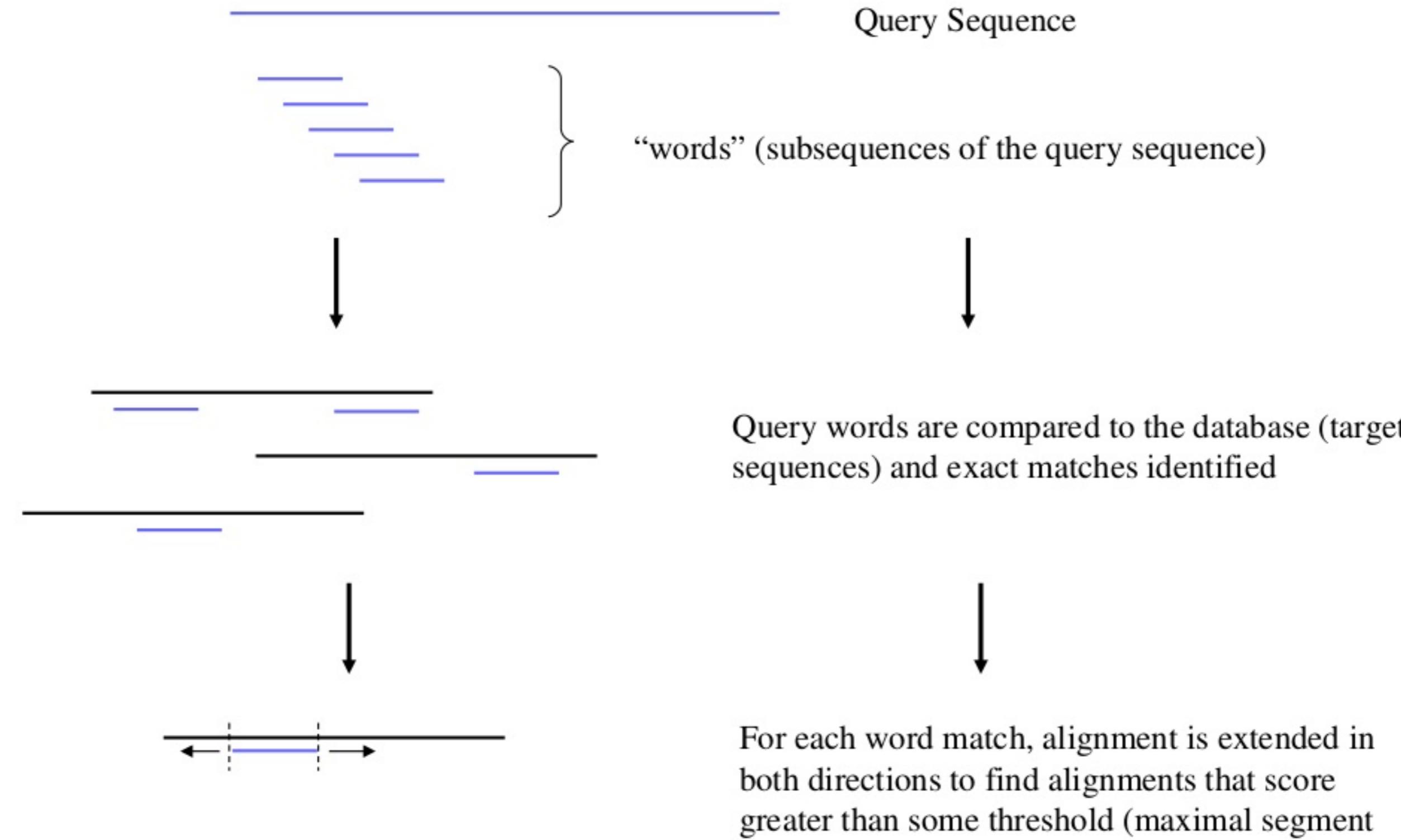
Basic Local Alignment Search Tool



# Basic Local Alignment Search Tool

- Widely used similarity search tool
- Heuristic approach based on Smith Waterman algorithm
- Finds best local alignments
- Provides statistical significance
- All combinations (DNA/Protein) query and database.
  - DNA vs DNA
  - DNA translation vs Protein
  - Protein vs Protein
  - Protein vs DNA translation
  - DNA translation vs DNA translation
- www, standalone, and network clients

# How BLAST works - pictoral



<http://blast.ncbi.nlm.nih.gov/Blast.cgi>

## Basic BLAST

---

Choose a BLAST program to run.

[nucleotide blast](#)

Search a **nucleotide** database using a **nucleotide** query  
*Algorithms:* blastn, megablast, discontiguous megablast

[protein blast](#)

Search **protein** database using a **protein** query  
*Algorithms:* blastp, psi-blast, phi-blast

[blastx](#)

Search **protein** database using a **translated nucleotide** query

[tblastn](#)

Search **translated nucleotide** database using a **protein** query

[tblastx](#)

Search **translated nucleotide** database using a **translated nucleotide** query



# BLASTing a sequence at NCBI – programs

**BLAST** Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

New Designing or Testing PCR Primers? Try your search in Primer-BLAST. [Go](#)

**BLAST Assembled Genomes**

Choose a species genome to search, or [list all genomic BLAST databases](#).

<input type="checkbox"/> <a href="#">Human</a>	<input type="checkbox"/> <a href="#">Oryza sativa</a>	<input type="checkbox"/> <a href="#">Gallus gallus</a>
<input type="checkbox"/> <a href="#">Mouse</a>	<input type="checkbox"/> <a href="#">Bos taurus</a>	<input type="checkbox"/> <a href="#">Pan troglodytes</a>
<input type="checkbox"/> <a href="#">Rat</a>	<input type="checkbox"/> <a href="#">Danio rerio</a>	<input type="checkbox"/> <a href="#">Microbes</a>
<input type="checkbox"/> <a href="#">Arabidopsis thaliana</a>	<input type="checkbox"/> <a href="#">Drosophila melanogaster</a>	<input type="checkbox"/> <a href="#">Apis mellifera</a>

**Basic BLAST**

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a nucleotide database using a nucleotide query <i>Algorithms:</i> blastn, megablast, discontiguous megablast
<a href="#">protein blast</a>	Search protein database using a protein query <i>Algorithms:</i> blastp, psi-blast, phi-blast
<a href="#">blastx</a>	Search protein database using a translated nucleotide query
<a href="#">tblastn</a>	Search translated nucleotide database using a protein query
<a href="#">tblastx</a>	Search translated nucleotide database using a translated nucleotide query

**News**

[Align Sequences with BLAST](#)

A new Bi2seq functionality has been added to the standard BLAST pages that allows you to align a query against a set of subject sequences.  
2008-09-04 12:56:52

[More BLAST news...](#)

**Tip of the Day**

[How to Search Custom Databases in Web-Blast Using Entrez Queries.](#)

A powerful feature of the BLAST Web interface is the ability to limit BLAST searches to a subset of any database using a standard Entrez query. Skillful use of Entrez queries allows the equivalent of on-the-fly construction of databases of exact composition.

[More tips...](#)

# BLASTing a sequence at NCBI – enter accession

BLAST Basic Local Alignment Search Tool My NCBI [Sign In] [Register]

NCBI/ BLAST/ blastp suite

blastn blastp blastx tblastn tblastx

Enter Query Sequence

Enter accession number, gi, or FASTA sequence  Clear

Query subrange  From   
To

Or, upload file  Browse... Job Title

Enter a descriptive title for your BLAST search

Blast 2 sequences

Choose Search Set

Database Non-redundant protein sequences (nr)

Organism Optional Enter organism name or id—completions will be suggested  
Enter organism common name, binomial, or taxid. Only 20 top taxa will be shown.

Entrez Query Optional Enter an Entrez query to limit search

Program Selection

Algorithm  blastp (protein-protein BLAST)

## BLASTing a sequence at NCBI – enter sequence

# BLASTing a sequence at NCBI - job status

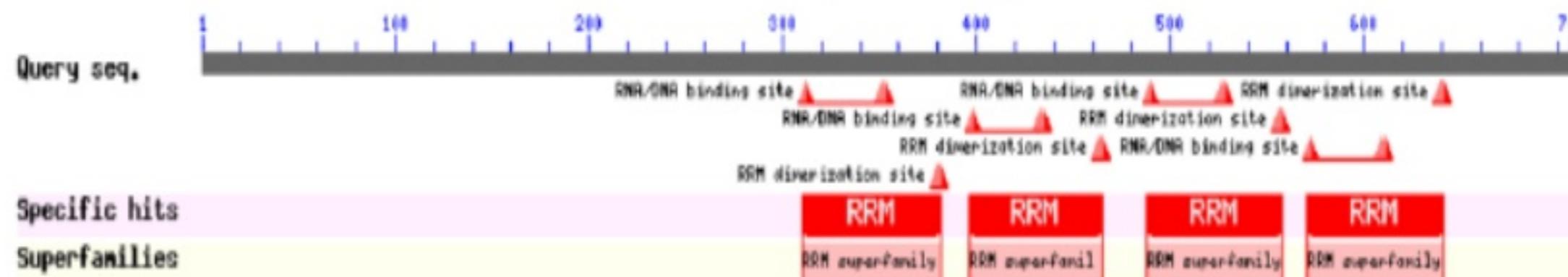
BLAST Basic Local Alignment Search Tool My NCBI [Sign In] [Register]

Home Recent Results Saved Strategies Help

NCBI BLAST/blastp suite/Formatting Results - SVTzb9YT011 [Formatting options]

Job Title: gi|128843|sp|P09405.2|NUCL\_MOUSE RecName...

Putative conserved domains have been detected, click on the image below for detailed results.



Request ID	SVTzb9YT011
Status	Searching
Submitted at	Sat Feb 7 14:58:35 2009
Current time	Sat Feb 7 14:58:48 2009
Time since submission	00:00:13

This page will be automatically updated in 6 seconds

# BLASTing a sequence at NCBI – blast summary

**BLAST** Basic Local Alignment Search Tool My NCBI [Sign In] [Register]

Home Recent Results Saved Strategies Help

NCBI BLAST! blastp suite! Formatting Results - SVTZB9YT011

Edit and Resubmit Save Search Strategies ►Formatting options ►Download

gi|128843|sp|P09405.2|NUCL\_MOUSE RecName:...

Query ID Icl|22157 Database Name nr  
Description gi|128843|sp|P09405.2|NUCL\_MOUSE RecName:  
Full=Nudeolin; AltName: Full=Protein C23 Description All non-redundant GenBank CDS  
translations+PDB+SwissProt+PIR+PRF excluding  
environmental samples from WGS projects  
Molecule type amino acid Program BLASTP 2.2.19+ ►Citation  
Query Length 707  
Other reports: ►Search Summary [Taxonomy reports] [Distance tree of results] [Related Structures]

▼ Graphic Summary

▼ Show Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.

Query seq. 188 288 388 488 588 688 787  
RNA/DNA binding site RNA/DNA binding site RRM dimerization site  
RNA/DNA binding site RRM dimerization site RRM dimerization site RNA/DNA binding site  
RRM dimerization site

Specific hits Superfamilies RRM RRM RRM RRM  
RRM superfamily RRM superfamily RRM superfamily RRM superfamily

Distribution of 172 Blast Hits on the Query Sequence ⓘ  
Mouse over to see the define, click to show alignments

Color key for alignment scores  
<40 40-50 50-80 80-200 >=200

# BLASTing a sequence at NCBI – used parameters

Other reports: [▼Search Summary](#) [Taxonomy reports](#) [Distance tree of results](#) [Related Structures](#)

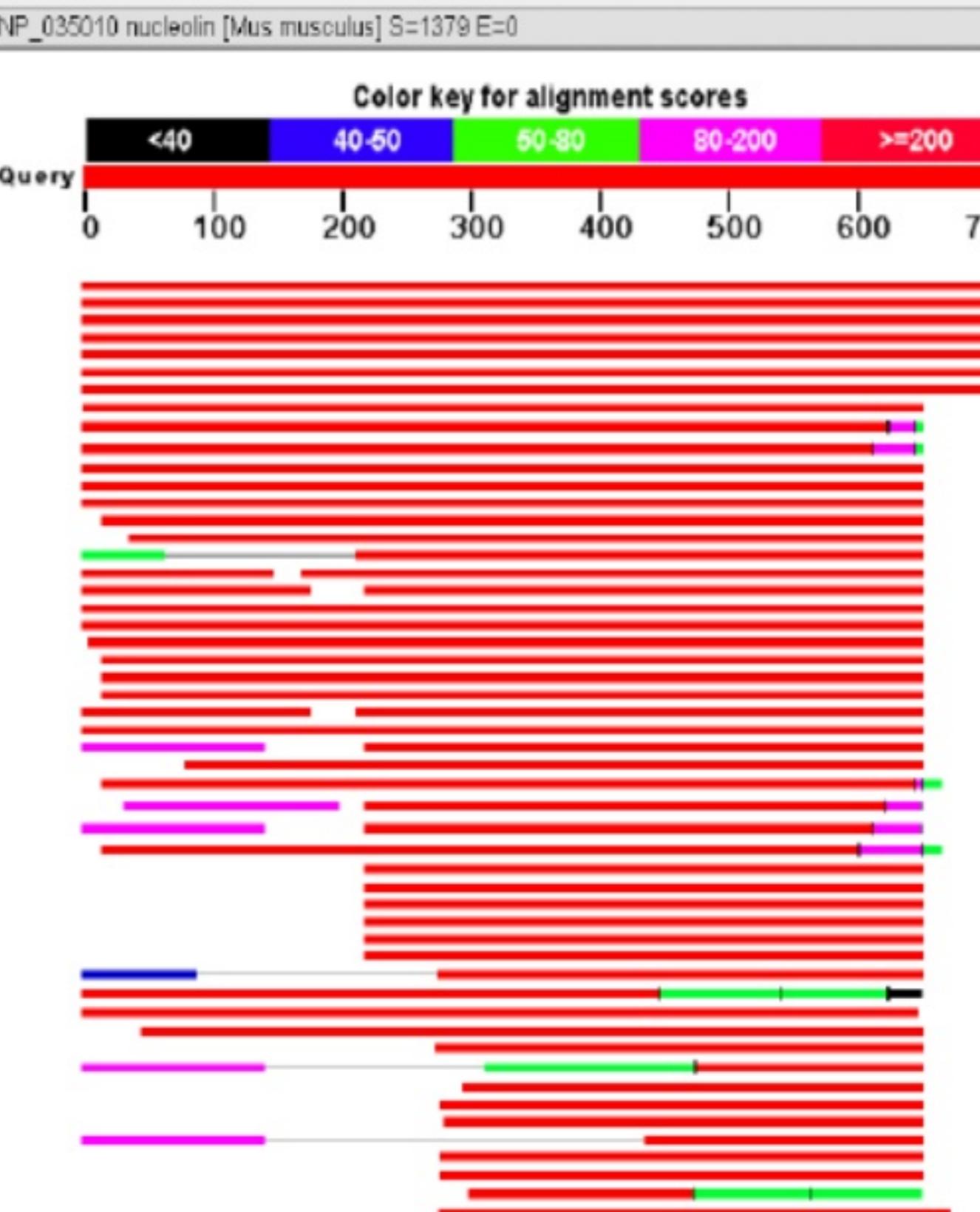
Search Parameters	
Program	blastp
Word size	3
Expect value	10
Hitlist size	100
Gapcosts	11,1
Matrix	BLOSUM62
Threshold	t1
Composition-based stats	2
Filter string	F
Genetic Code	1
Window Size	40

Database	
Posted date	Feb 6, 2009 5:53 PM
Number of letters	2,699,408,701
Number of sequences	7,831,890
Entrez query	none

Karlin-Altschul statistics		
Params	Ungapped	Gapped
Lambda	0.302661	0.267
K	0.127079	0.041
H	0.344587	0.14

Results Statistics	
Length adjustment	143
Effective length of query	564
Effective length of database	1579448431

## BLASTing a sequence at NCBI – graphical display

Distribution of 172 Blast Hits on the Query Sequence [?](#)

# BLASTing a sequence at NCBI – hit list

**▼ Descriptions**

Sequences producing significant alignments:		Score (Bits)	E Value	
<a href="#">ref NP_035010.3 </a>	nucleolin [Mus musculus] > <a href="#">sp P09405.2 NUCL_M...</a>	<a href="#">1379</a>	<a href="#">0.0</a>	
<a href="#">dbj BAE36484.1 </a>	unnamed protein product [Mus musculus]	<a href="#">1378</a>	<a href="#">0.0</a>	
<a href="#">dbj BAE38940.1 </a>	unnamed protein product [Mus musculus]	<a href="#">1378</a>	<a href="#">0.0</a>	
<a href="#">dbj BAE40448.1 </a>	unnamed protein product [Mus musculus] > <a href="#">dbj B...</a>	<a href="#">1375</a>	<a href="#">0.0</a>	
<a href="#">dbj BAC26311.1 </a>	unnamed protein product [Mus musculus]	<a href="#">1373</a>	<a href="#">0.0</a>	
<a href="#">gb AAH05460.1 </a>	Nucleolin [Mus musculus]	<a href="#">1371</a>	<a href="#">0.0</a>	
<a href="#">dbj BAC27474.1 </a>	unnamed protein product [Mus musculus]	<a href="#">1363</a>	<a href="#">0.0</a>	
<a href="#">gb EDL40224.1 </a>	nucleolin, isoform CRA_e [Mus musculus]	<a href="#">1009</a>	<a href="#">0.0</a>	
<a href="#">gb EDL40223.1 </a>	nucleolin, isoform CRA_d [Mus musculus]	<a href="#">966</a>	<a href="#">0.0</a>	
<a href="#">gb EDL40222.1 </a>	nucleolin, isoform CRA_c [Mus musculus]	<a href="#">942</a>	<a href="#">0.0</a>	
<a href="#">sp P13383.3 NUCL RAT</a>	RecName: Full=Nucleolin; AltName: Full=...	<a href="#">941</a>	<a href="#">0.0</a>	
<a href="#">ref NP_036881.2 </a>	nucleolin [Rattus norvegicus] > <a href="#">gb AAH85751.1 ...</a>	<a href="#">941</a>	<a href="#">0.0</a>	
<a href="#">sp P08199.2 NUCL MESAU</a>	RecName: Full=Nucleolin; AltName: Full...	<a href="#">919</a>	<a href="#">0.0</a>	
<a href="#">gb EDL75577.1 </a>	nucleolin, isoform CRA_b [Rattus norvegicus]	<a href="#">912</a>	<a href="#">0.0</a>	
<a href="#">gb AAA36966.1 </a>	nucleolin, C23	<a href="#">893</a>	<a href="#">0.0</a>	
<a href="#">gb EDL40220.1 </a>	nucleolin, isoform CRA_a [Mus musculus]	<a href="#">797</a>	<a href="#">0.0</a>	
<a href="#">dbj BAC34476.1 </a>	unnamed protein product [Mus musculus]	<a href="#">796</a>	<a href="#">0.0</a>	
<a href="#">gb EDL40221.1 </a>	nucleolin, isoform CRA_b [Mus musculus]	<a href="#">786</a>	<a href="#">0.0</a>	
<a href="#">gb AAD56625.1 AF151373_1</a>	nucleolin-related protein NRP [Rattu...	<a href="#">781</a>	<a href="#">0.0</a>	
<a href="#">sp Q4R4J7.3 NUCL MACFA</a>	RecName: Full=Nucleolin > <a href="#">dbj BABD0345....</a>	<a href="#">768</a>	<a href="#">0.0</a>	
<a href="#">ref XP_001116949.1 </a>	PREDICTED: similar to nucleolin [Macaca m...	<a href="#">762</a>	<a href="#">0.0</a>	
<a href="#">ref XP_861643.1 </a>	PREDICTED: similar to nucleolin-related prot...	<a href="#">761</a>	<a href="#">0.0</a>	
<a href="#">ref XP_861613.1 </a>	PREDICTED: similar to nucleolin-related prot...	<a href="#">761</a>	<a href="#">0.0</a>	
<a href="#">ref XP_850477.1 </a>	PREDICTED: similar to nucleolin-related prot...	<a href="#">761</a>	<a href="#">0.0</a>	
<a href="#">gb EDL75581.1 </a>	nucleolin, isoform CRA_e [Rattus norvegicus]	<a href="#">756</a>	<a href="#">0.0</a>	
<a href="#">ref XP_516145.2 </a>	PREDICTED: hypothetical protein [Pan troglod...	<a href="#">755</a>	<a href="#">0.0</a>	
<a href="#">gb EDL75579.1 </a>	nucleolin, isoform CRA_d [Rattus norvegicus] >...	<a href="#">749</a>	<a href="#">0.0</a>	

How often would this hit have occurred by chance?

Rule of thumb:  
E-value < 0.0001

# BLASTing a sequence at NCBI

>gb|AAP62554.1| G nucleolin [Oncorhynchus mykiss]  
Length=255

[GENE ID: 100135911](#) [LOC100135911](#) | nucleolin [Oncorhynchus mykiss]

Score = 239 bits (610), Expect = 7e-61, Method: Compositional matrix adjust.  
Identities = 133/260 (51%), Positives = 182/260 (70%), Gaps = 11/260 (4%)

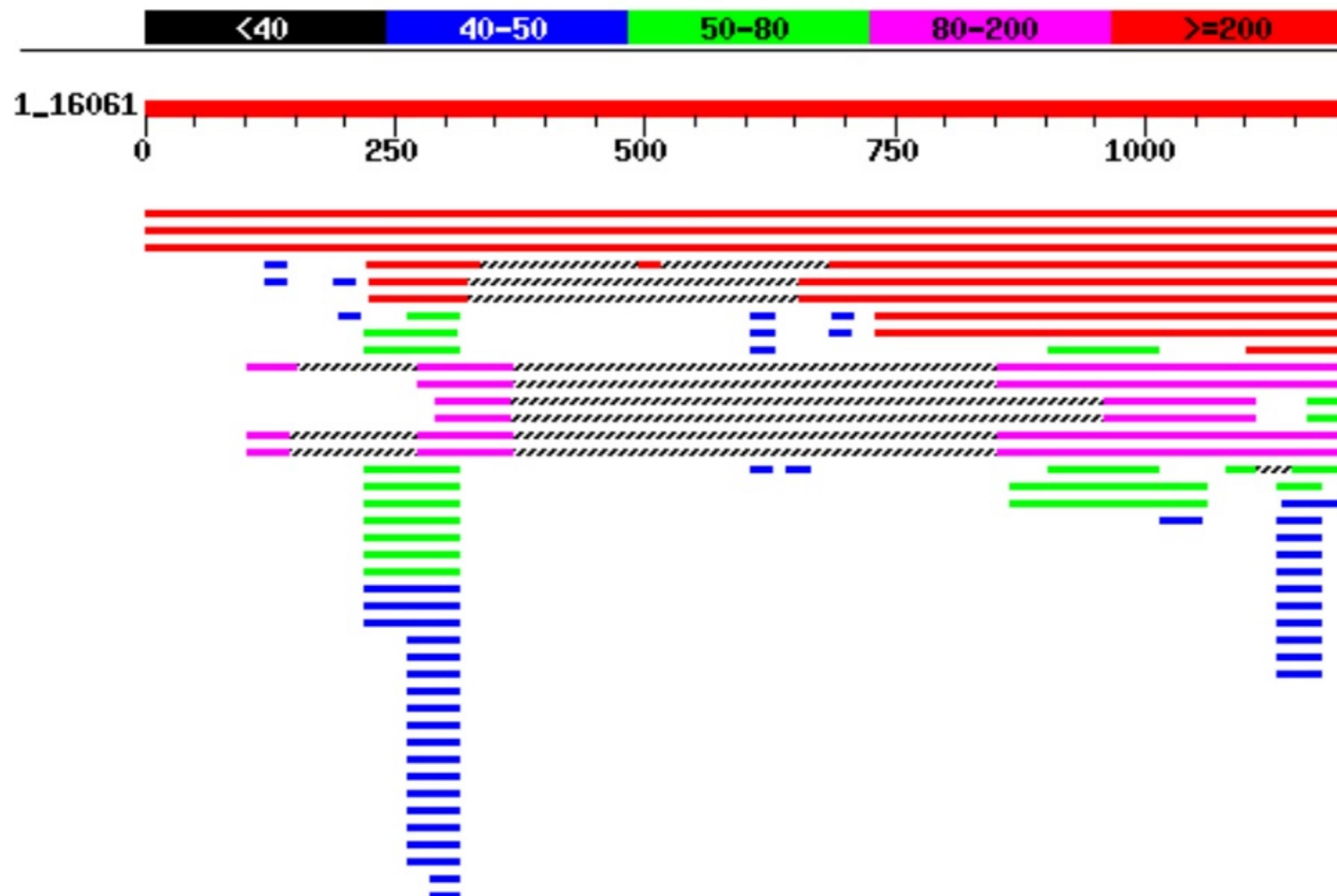
Query 283	KRBMTKQKBAPEAKRQKVEGSSEPTPFNLFIGNINENKSVNELKFAISELFAKNDLAWD	342
	K++ +KE D AKK K SE F LFIGNIN NK +E+K A++ F+K +L V D	
subjct 2	KRKADNKKETPPAKKAK---SESDDTFCLFIGNINENKDFDEIKEALAAFFSKKNLEVQD	58
Query 343	VRTGTINRKPGYVDFESAEDLEKALBLTGLKVGNEBIRLERPKGR---DSKEVRAARTLL	398
	VR G ++KFGYV++ SAED++ A+BL G K G E+K++K + + + KK R ARTL	
subjct 59	VRLGASKKPGYVFASAEADMQTAMELNGKKCMGQELKMDKARSKGNSQEEKKDRDARTLF	118
Query 399	AENLSPNITEDDELKEVFEAMBIKL-VSQDGK8RGIAYIEFKS8ADAERKNNLEEKQGAEIF	457
	EKL F+ TED+LKEVF +A+BIRT QDG ++GIAYI FK+EA A+K I E OGAA+	
subjct 119	VENLPSSATEDDLKEVFAANAVEBIRIPGQDGGSNRGIAYIAFPKTBAKADKMLTEAQGADVO	178
Query 458	GRSVSILYYTGERGQRQERTGKTSTWSGB8SKTLVLSNLSSATKETLEEVEEKFIFKVPQ	517
	GRSt + YTG K Q+ R + + ESKTL+++NLSSAT++L+ FE A ItVPO	
subjct 179	GRSIMVDYTPGIKSQEGGRP--PAQAAAESKTLIVNNLSSATEDSLQSAFESEGAVSIRVPQ	236
Query 518	NPHGKPKGYAFIEFASFEDA 537	
	N +G+PKG+A++EF S E A	
subjct 237	N-NGRPKGPAPVEFESAEKA 255	

Score = 99.8 bits (247), Expect = 8e-19, Method: Compositional matrix adjust.  
Identities = 76/242 (31%), Positives = 118/242 (48%), Gaps = 29/242 (11%)

Query 396	TLLAKNLSPNITEDDELKEYFE-----DAMEIRLVSDGK8RGIAYIEFKSEADAEN	447
	L NL+ N DE+KE + ++RL G SK Y+EF S D +	
subjct 26	CLPIGNLN8NKDPDEIKEALAAFFSKKNLEVQDVRL---GASKKPGYVFASAEADMQTA	81
Query 448	LEEKQGAEIDGRVSILYYTGERGQRQERTGKTSTWSGB8SKTLVLSNLSSATKETLEEVF	507
	+E G + G+ + + KG QE +++TL + NL +SAT++ L+EVF	
subjct 82	ME-LNGKKCMGQELKMDKARSKGNSQEEKKDR----DARTLFVKNLDFSATEDDLKEVF	135
Query 508	EKATFIKVPOINPHGKPKGYAFIEFASFEDAKEALNSCNKMEIEGR7IRLQLQGSNSR---	564
	A I++P G +G A+I F + A + L +++GR+I ++ G S+	
subjct 136	ANAVEIRIPGQDGGSNRGIAYIAFPKTEAMADKMLTEAQGADVQGREIMVDYTGIKSQKGG	195
Query 565	-----SQDEKTLFVKGLSEDTTTEETLKESFEGSYRARIVTDRETGS8KCFGBYDFNSEE	618
	+ EKTL V LE TE++L+ +FEG+V R+ + G KGF FV+F S E	
subjct 196	RPPAQAAAESKTLIVNNLSSATEDSLQSAFESEGAVSIRV--DQINGRDKGPFAFVEFESAE	253



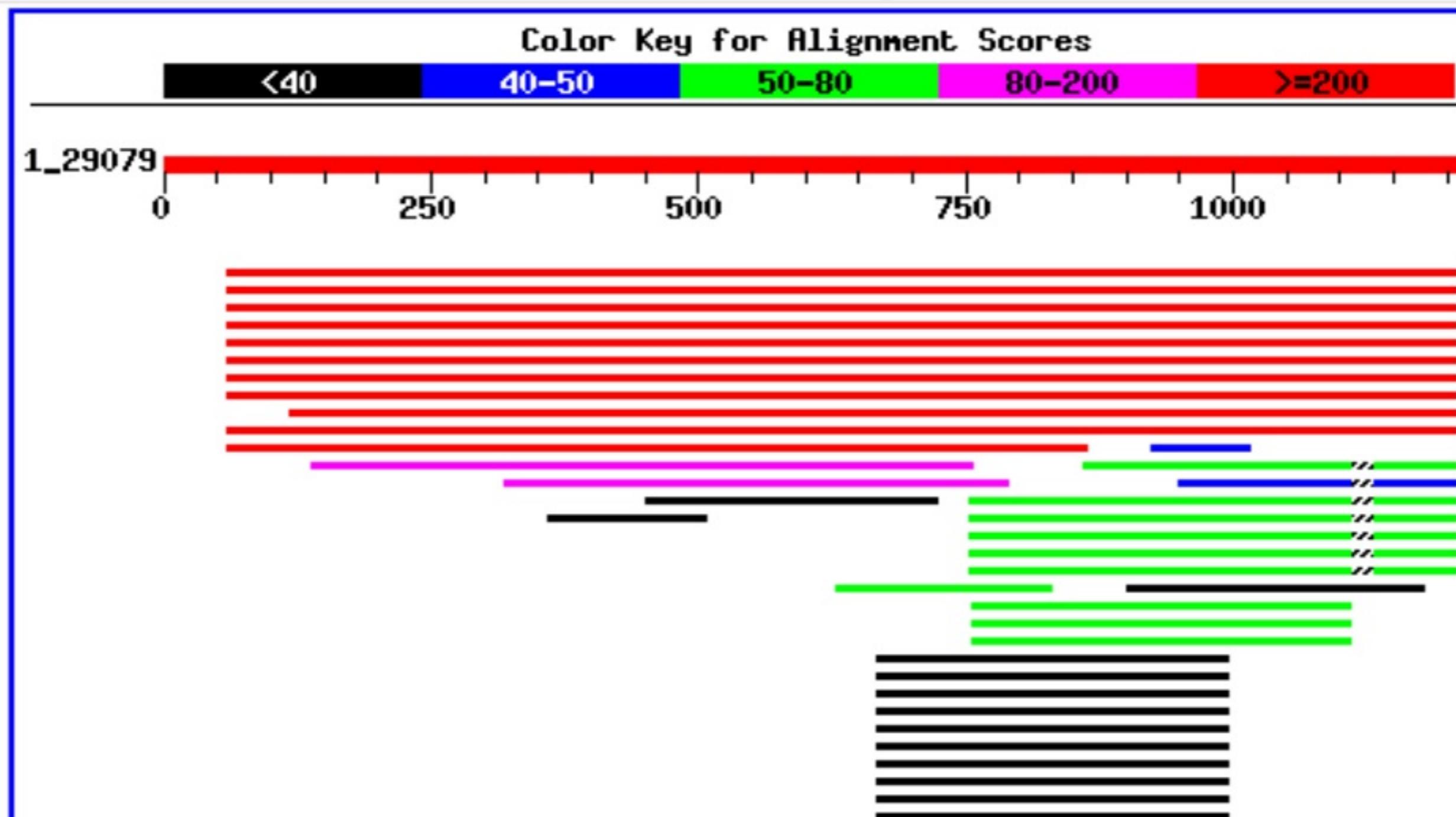
## Result Page 1 of BLASTn



# Result Page of BLASTx

## Distribution of 59 Blast Hits on the Query Sequence

Mouse-over to show defline and scores. Click to show alignments



# Result Page of BLASTx

			Score	E	
			(bits)	Value	
<b>Sequences producing significant alignments:</b>					
<a href="#">gi 51095067 gb EAL24310.1 </a>	IMP (inosine monophosphate) dehy...	<a href="#">411</a>	e-125	G	
<a href="#">gi 57999523 emb CAI45968.1 </a>	hypothetical protein [Homo sapi...	<a href="#">411</a>	e-125	G	
<a href="#">gi 51095068 gb EAL24311.1 </a>	IMP (inosine monophosphate) dehy...	<a href="#">411</a>	e-125	G	
<a href="#">gi 25014074 sp P20839 IMD1_HUMAN </a>	Inosine-5'-monophosphate d...	<a href="#">411</a>	e-125	G	
<a href="#">gi 106722 pir  A35566 </a>	IMP dehydrogenase (EC 1.1.1.205) I - ...	<a href="#">411</a>	e-125	S	
<a href="#">gi 13543973 gb AAH06124.1 </a>	IMP (inosine monophosphate) dehy...	<a href="#">407</a>	e-124	G S	
<a href="#">gi 4504689 ref NP_000875.1 </a>	IMP (inosine monophosphate) deh...	<a href="#">406</a>	e-123	G	
<a href="#">gi 307067 gb AAA36114.1 </a>	IMP dehydrogenase type 1 (EC 1.1.1...	<a href="#">407</a>	e-123	G	
<a href="#">gi 47077068 dbj BAD18464.1 </a>	unnamed protein product [Homo s...	<a href="#">397</a>	e-120	G	
<a href="#">gi 16549223 dbj BAB70780.1 </a>	unnamed protein product [Homo s...	<a href="#">357</a>	e-108	G	
<a href="#">gi 51467033 ref XP_496992.1 </a>	PREDICTED: similar to inosine ...	<a href="#">246</a>	8e-65	G	
<a href="#">gi 51467471 ref XP_497019.1 </a>	PREDICTED: similar to Inosine-...	<a href="#">197</a>	3e-50	G	
<a href="#">gi 44979607 gb AA550155.1 </a>	IMP dehydrogenase 2 [Homo sapiens]	<a href="#">147</a>	3e-35	G	
<a href="#">gi 45708411 gb AAH03053.1 </a>	GMPR2 protein [Homo sapiens]	<a href="#">76</a>	3e-20	G	
<a href="#">View All</a>					



>gi|51095067|gb|EAL24310.1| G IMP (inosine monophosphate) dehydrogenase 1 [Homo sapiens]  
 gi|34328930|ref|NP\_000874.2| G inosine monophosphate dehydrogenase 1 isoform a [Homo sapiens]  
 Length = 599

Score = 411 bits (1057), Expect(2) = e-125  
 Identities = 207/356 (58%), Positives = 271/356 (76%), Gaps = 4/356 (1%)  
 Frame = +1

Query: 61 DGLSVQELMDSKIRGGLAYNDFLILPGLVDFASSEVSLQTKLTRNITLNIPLVSSPMDTV 240  
 DGL+ Q+L S GL YNDFLILPG +DF + EV L + LTR ITL PL+SSPMDTV  
 Sbjct: 101 DGLTAQQQLFASA--DGLTYNDFLILPGFIDFIADEVDLTSALTRKITLKPLISSPMDTV 158

Query: 241 TESEMATFMALLDGIGFIHHNCTPEDQADMVRVKNYENGFINNPIVISPTTVGEAKSM 420  
 TE++MA MAL+ GIGFIHHNCTPE QA+ VR+VK +E GFI +P+V+SP+ TVG+  
 Sbjct: 159 TEADMAIAMALMGGIGFIHHNCTPEFQANEVRKVKKFEQGFITDPVVLSPSHTVGDVLEA 218

Query: 421 KEKYGFAGFPVTADGKRNAKLVGAITSRDIQFV--EDNSLLVQDVMTK-NPVTGAQGIT 588  
 K ++GF+G P+T G +KLVG +TSRDI F+ +D++ L+ +VMT V G+T  
 Sbjct: 219 KMRHGFSGIPITETGTMGSKLVGVTSRDIDFLAEKDHTLLSEVMTPRIELVVAPAGVT 278

Query: 589 LSEGNEILKKIKKGRLVVDEKGNLVSMLSRTDLMKNQKYPLASKSANTKQLLWGASIGT 768  
 L E NEIL++ KKG+L +V++ LV++++RTDL KN+ YPLASK + KQLL GA++GT  
 Sbjct: 279 LKEANEILQRSKKGKLPIVNDCDELVAIIAARTDLKKNRDYPLASKDSQ-KQLLCGAAVGT 337

Query: 769 MDADKERLRLVKAGLDVVILDSSQGNSIFQLNMKWIETFPDLEIIAGNVVTKEQAAN 948  
 + DK RL LL +AG+DV++LDSSQGNS++Q+ M+ +IK+ +P L++I GNVVT QA N  
 Sbjct: 338 REDDKYRLDLLTQAGVDVIVLDSSQGNSVYQIAMVHYIKQKYPHLQVIGGNVVTAAQAKN 397

Query: 949 LIAAGADGLRIGMGTGSICITQKVMACGRPQGTAVYNVCEFANQFGVPCMADGGVQ 1116  
 LI AG DGLR+GMG GSICITQ+VMACGRPQGTAVY V E+A +FGVP +ADGG+Q  
 Sbjct: 1209 LHDAGVDPGLVCMGCCSICITQVMAACGPQGTAVYKVAEYARPGCVPHDGGQ 452

Score = 46.2 bits (108), Expect(2) = e-125  
Identities = 21/32 (65%), Positives = 28/32 (87%)  
Frame = +2

Query: 1115 KNIGHITKALALGSSTVMMGGMLAGTTEAPG 1210  
+ +GH++ KALALG+STVMMG +LA TTE+PG

Sbjct: 453 QTVGHVV-KALALGASTVMMGSLLAATTEAPG 483

Fig 5.14: Result Page of Pairwise Alignment



Address <http://www.ncbi.nlm.nih.gov/projects/gorf/>

Search GO More 10 Site Rating

softonic SEARCH WEB GADGETS Games

# ORF Finder (Open Reading Frame Finder)

PubMed Entrez BLAST OMIM Taxonomy Structure

NCBI

Tools for data mining

GenBank sequence submission support and software

FTP site download data and software

The ORF Finder (Open Reading Frame Finder) is a graphical analysis tool which finds all open reading frames of a selectable minimum size in a user's sequence or in a sequence already in the database. This tool identifies all open reading frames using the standard or alternative genetic codes. The deduced amino acid sequence can be saved in various formats and searched against the sequence database using the WWW BLAST server. The ORF Finder should be helpful in preparing complete and accurate sequence submissions. It is also compatible with the GenBank sequence submission

Enter GI or ACCESSION

OrFFind

Clear

or sequence in FASTA format


FROM:

TO:

Genetic codes

1 Standard



Comments and suggestions to:

[info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)

Search Taxonomy Taxon ID Taxon Name Taxon Type



Enter GI or ACCESSION

OrfFind

Clear

or sequence in FASTA format


FROM:  TO:

Genetic codes

1 Standard



---

Comments and suggestions to:

[info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)

ORIGIN TRANSLATE PROTEIN TBLASTN TBLASTX BLASTN BLASTX



**Enter GI or ACCESSION****or sequence in FASTA format**

```
GC
CACTTCTGAGTTGGGGCAGCGGGTTCTAGCTCAG
CTCATGCTGAGAATGTAAGAACTACAAACAAAAAT
TT
CTATTAAAATTAAAGTTTGTGTCTTAAAAA
AAAAAA
```

**FROM:****TO:**Genetic codes

1 Standard



---

*Comments and suggestions to:*[info@ncbi.nlm.nih.gov](mailto:info@ncbi.nlm.nih.gov)Credit to: Tatjana Tatusova and Roman Tatusov



# ORF Finder (Open Reading Frame Finder)

[PubMed](#)[Entrez](#)[BLAST](#)[OMIM](#)[Taxonomy](#)[Structure](#)

gi|31343070|ref|NM\_174000.2| Bos taurus calreticulin (CALR), mRNA

[View](#)[1 GenBank](#)[Redraw](#)[100](#)[SixFrames](#)

Frame	from	to	Length
+2	71..1324	1254	
-1	77..505	429	
+3	705..1013	309	
+3	1155..1361	207	
-2	613..795	183	
+2	1646..1771	126	
+3	240..362	123	
+1	1..123	123	
-2	868..984	117	



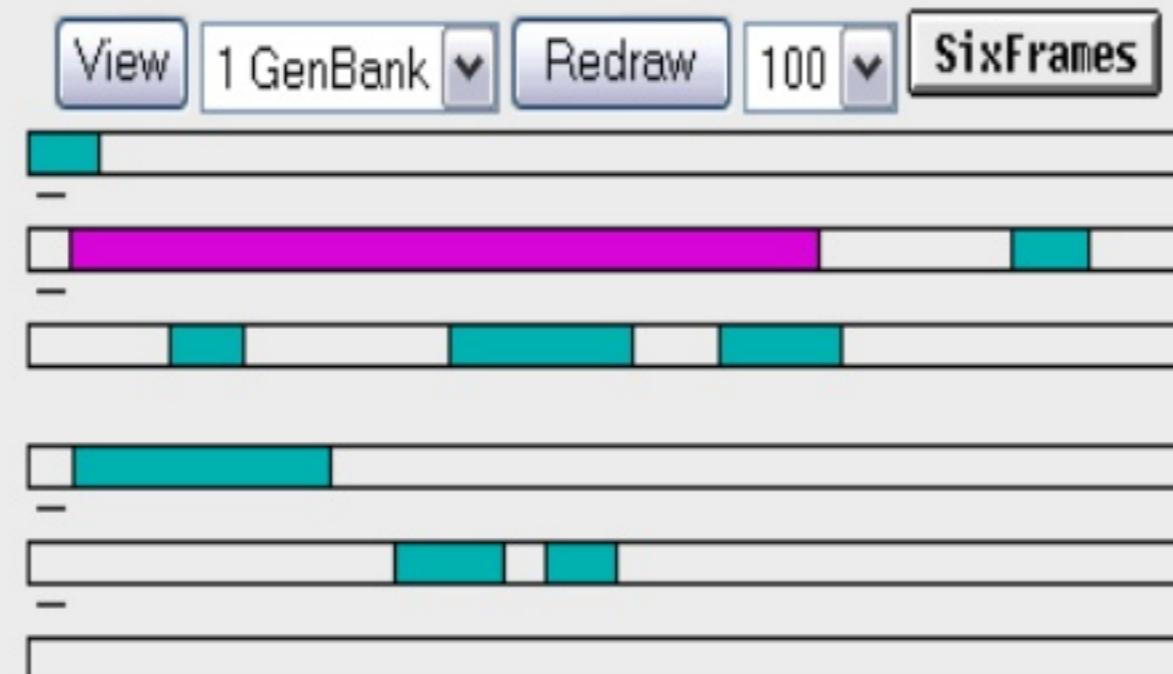


# ORF Finder (Open Reading Frame Finder)

[PubMed](#)[Entrez](#)[BLAST](#)[OMIM](#)[Taxonomy](#)[Structure](#)

gi|31343070|ref|NM\_174000.2| Bos taurus calreticulin (CALR), mRNA

Program: blastp    Database: nr    BLAST    with parameters    Cognitor



Frame	from to	Length
+2	71..1324	1254
-1	77.. 505	429
+3	705..1013	309
+3	1155..1361	207
-2	613.. 795	183
+2	1646..1771	126
+3	240.. 362	123
+1	1.. 123	123
-2	868.. 984	117

Length: 417 aa



71 **atgtctgttacccgtggcgctgtgtgttttcggcctggcc**  
**H L L P V P L L L G L L G L A**  
 116 **gcgcgtgatcccaccgtttacttaaggagcagtttctggacgga**  
**A A D P T V Y F K E Q F L D G**  
 161 **gacgggtggaccggagcgatggatcgaatccaaggcacaaaccggat**  
**D C W T E R W I E S K H K P D**  
 206 **tttggcaatttgttctcagttccggcaagtttatggtgaccas**  
**F G K F V L S S G K F Y G D Q**  
 251 **gagaaagataaaaggcctgcagactagccaggatgccccggttctac**  
**E K D K G L Q T S Q D A R F Y**  
 296 **gtctgtcgccagatggagcccttcagcaacaagggtcagacg**  
**A L S A R F E P F S N K G Q T**  
 341 **ctgggtggtgcagttcacagtgaacacacgaggcagaacatcgactgt**  
**L V V Q F T V K H E Q N I D C**  
 386 **ggggggggctatgtgaagctgtttccagctgggttggatcagaca**  
**C C C Y V K L F P A C L D Q T**  
 431 **gac**atg**cacggagactccgaatacacaata**atg**tttggccggac**  
**D H H G D S E Y N I H F G P D**  
 476 **atctgtggccctggcaccataaaggatcatgtcatcttcaactac**  
**I C G P G T K K V H V I F N Y**  
 521 **aaggggcaagaatgtgtatcaacaaggatatccgctgcaaggac**  
**K G K N V L I N K D I R C K D**  
 566 **gatgaattcacccacgtacacgtgttgcgggctaaataat**  
**D E F T H L Y T L I V R P N N**  
 611 **acctatgaggtgaagattgacaaacagccagggtggagtcaggctc**  
**T Y E V K I D N S Q V E S G S**  
 656 **ttggaggacgattgggattttttgccaccaagaagataaaggat**  
**L E D D W F L P P R K K I K D**  
 701 **cctgatggcgttaaggctgtaaagactgggacgtcgcccaagatc**  
**P D A A R P E D W D D R A K I**  
 746 **gatgaccctacagactccaaaggctgttgcggattgggacaaaggctgag**  
**D D P T D S K P E D W D K P E**  
 791 **cacattctgaccctgtatgttaagaaacctgaggactgggatgaa**  
**H I P D P D A K K P E D W D E**  
 836 **gag**atg**gacggagagtggaaaccacactgttattcagaacccagag**  
**E H D G E W E P P V I Q N P E**  
 881 **tacaaggggaaatggaaaccacggcagatcgacaaacccagagttac**  
**Y K G E W K P R Q I D N P E Y**  
 926 **aagggcatttggatccacccagagattgacaaacctgttgttcc**  
**K G I W I H P E I D N P E Y S**  
 971 **cctgacagcaacatctatgttgcataaaaacttcgtgttcttaggc**  
**P D S N I Y A Y E N F A V L G**  
 1016 **ttggatcttggcaggtcaagtctggcaccatcttgcataacttc**  
**L D L W Q V K S G T I F D N F**  
 1061 **ctcatcaccaacgtgaagcgtatgttgcggatggcaacgg**  
**L I T N D E A Y A E E F G N E**  
 1106 **acgtgggtgttacaaaggcagcagaaaagcaa**atg**aaggacaaag**  
**T W G V T K A A E K Q H K D K**  
 1151 **caggatgaagagcagaggctacatgaggaggaggagaagaaa**  
**Q D E E Q R L H E E E E K K**  
 1196 **ggcaaggaggaggaagagaggcagacaaagatgtgacgaagacaag**  
**G K E E E A D K D D D E D K**  
 1241 **gatgaggatgaggaggatgaagatgagaaggaagaggaggagggaa**  
**D E D E F D F D F K E F E F E**



Address <http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>

Search GO More 10 Site Rating News GADGETS Games Play! re

softonic SEARCH WEB GADGETS Games Play! MTV

## BLAST Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

▶ NCBI/ BLAST/ Format Request

Query Protein Sequence (417 letters)

Database nr

Job title Protein Sequence (417 letters)

Request ID 5P5G6A2X01S   Show results in a new window

Format

Show Alignment as HTML  Advanced View  Use old BLAST report format [Reset to defaults](#)

Alignment View Pairwise [?](#)

Display  Graphical Overview  Linkout  Sequence Retrieval  NCBItaxID [?](#)

Masking Character: Lower Case Color: Grey [?](#)

Limit results Descriptions: 100 Graphical overview: 100 Alignments: 100 [?](#)

Organism Type common name, binomial, taxid, or group name. Only 20 top taxa will be shown.  
Enter organism name or id--completions will be suggested [?](#)

Entrez query: [?](#)

Expect Min: Expect Max: [?](#)

Format for  PSI-BLAST with inclusion threshold: [?](#)

**BLAST** Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

► NCBI/ BLAST/ blastp suite/ Formatting Results - 5P5G6A2X01S

Edit and Resubmit Save Search Strategies ►Formatting options ►Download

Protein Sequence (417 letters)

**Query ID:** Icl|31691  
**Description:** None  
**Molecule type:** amino acid  
**Query Length:** 417

**Database Name nr**  
**Description:** All non-redundant GenBank CDS translations from WGS projects  
**Program:** BLASTP 2.2.24+ ►Citation

Other reports: ►Search Summary [Taxonomy reports] [Distance tree of results] [Multiple alignment]

Graphic Summary

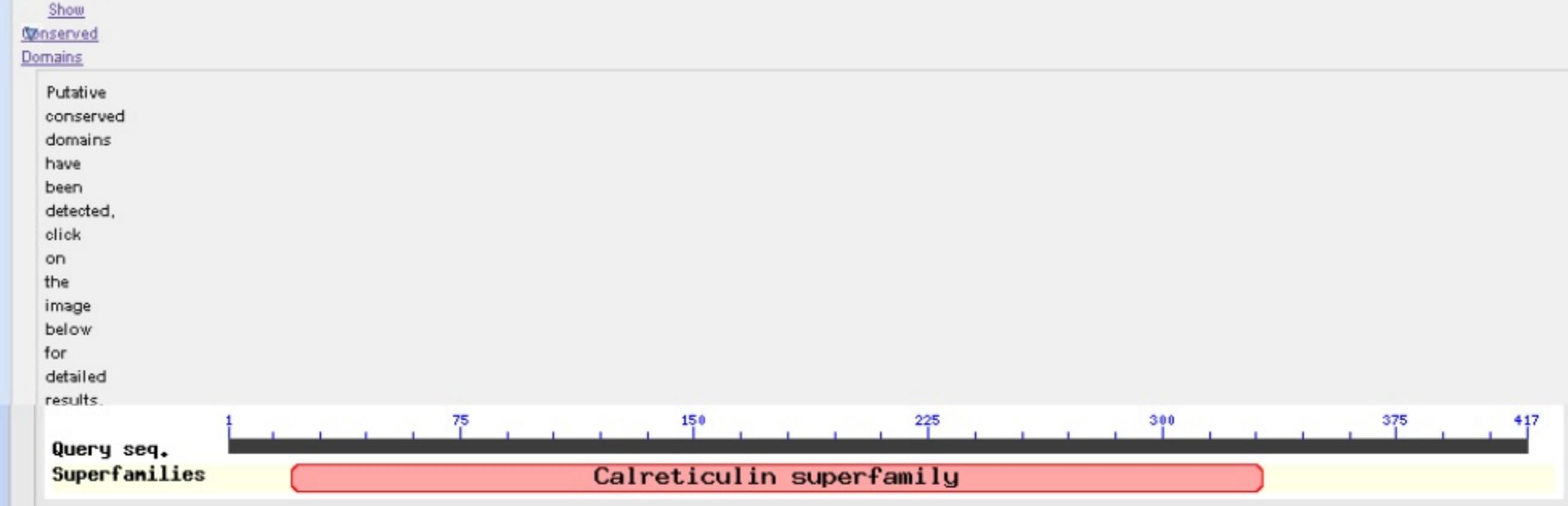
Show  
Conserved  
Domains

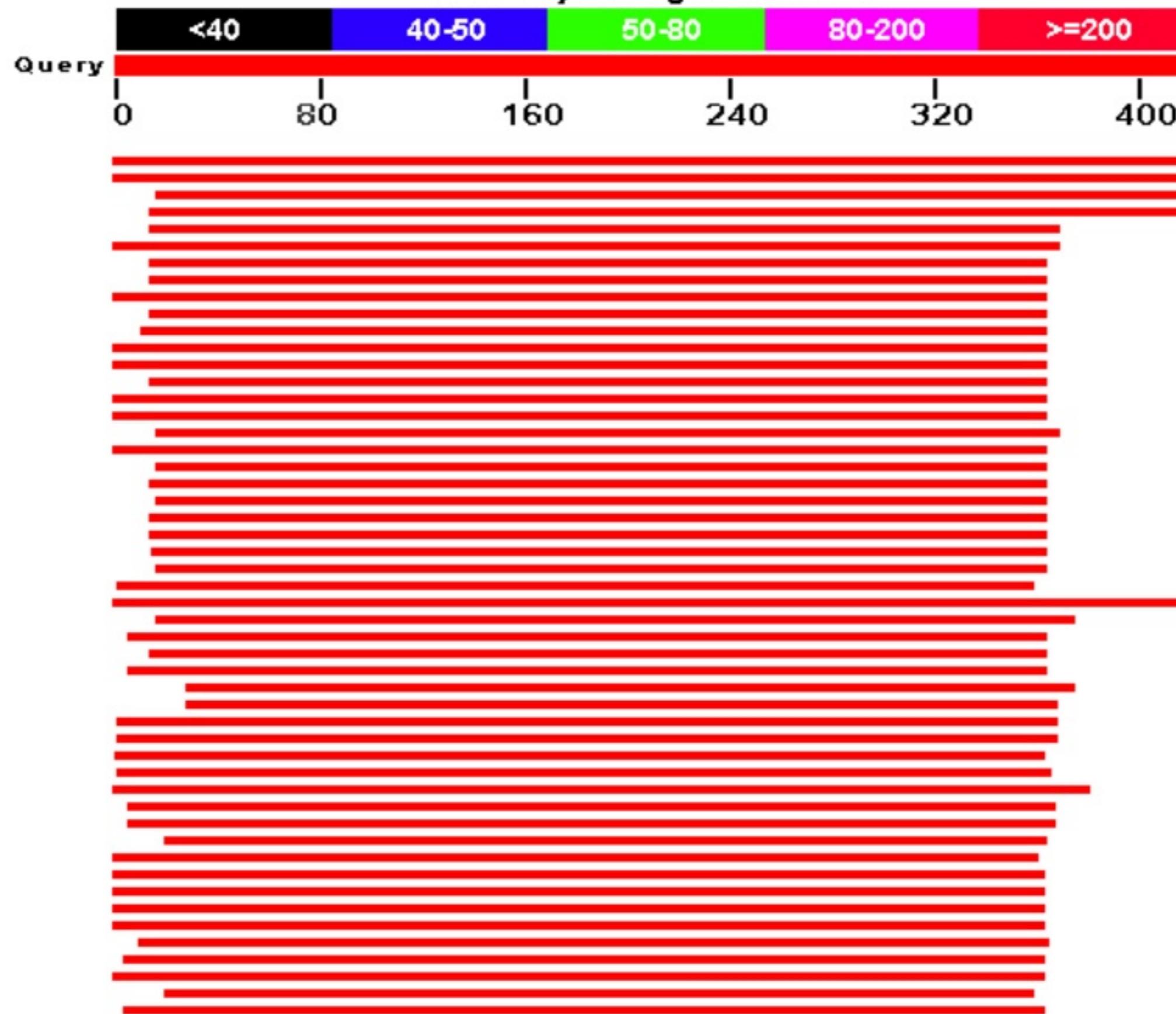
Putative  
conserved  
domains  
have  
been  
detected,  
click  
on  
the  
image  
below  
for  
detailed  
results.

Query seq.  
Superfamilies

Calreticulin superfamily

Distribution  
of  
100  
Blast  
Hits  
on



**Color key for alignment scores**

## ▼ Descriptions

Sequences producing significant alignments:	Score (Bits)	E Value	
<a href="#">ref NP_776425.1  calreticulin precursor [Bos taurus] &gt;sp P521...</a>	<a href="#">643</a>	0.0	
<a href="#">gb AAI40583.1  CALR protein [Bos taurus] &gt;gb DAA28020.1  calr...</a>	<a href="#">640</a>	0.0	
<a href="#">gb AAB30209.1  calreticulin [cattle, brain, Peptide, 400 aa]</a>	<a href="#">607</a>	0.0	
<a href="#">ref NP_001167504.1  calreticulin [Sus scrofa] &gt;gb ADD52600.1 ...</a>	<a href="#">719</a>	0.0	
<a href="#">ref NP_001075704.1  calreticulin precursor [Oryctolagus cunic...</a>	<a href="#">704</a>	0.0	
<a href="#">dbj BAE22855.1  unnamed protein product [Mus musculus]</a>	<a href="#">701</a>	0.0	
<a href="#">ref XP_853393.1  PREDICTED: similar to calreticulin isoform 2...</a>	<a href="#">698</a>	0.0	
<a href="#">ref XP_002921056.1  PREDICTED: calreticulin-like [Ailuropoda ...]</a>	<a href="#">697</a>	0.0	
<a href="#">ref NP_071794.1  calreticulin precursor [Rattus norvegicus] &gt;...</a>	<a href="#">695</a>	0.0	
<a href="#">ref XP_001504932.1  PREDICTED: similar to calreticulin [Equus...</a>	<a href="#">694</a>	0.0	
<a href="#">dbj BAD96780.1  calreticulin precursor variant [Homo sapiens]</a>	<a href="#">692</a>	0.0	
<a href="#">ref NP_031517.1  calreticulin precursor [Mus musculus] &gt;sp P1...</a>	<a href="#">692</a>	0.0	
<a href="#">dbj BAE35687.1  unnamed protein product [Mus musculus]</a>	<a href="#">691</a>	0.0	
<a href="#">ref XP_867320.1  PREDICTED: similar to calreticulin isoform 5...</a>	<a href="#">690</a>	0.0	
<a href="#">sp Q8K3H7.1 CALR CRIGR RecName: Full=Calreticulin; AltName: F...</a>	<a href="#">690</a>	0.0	
<a href="#">ref XP_001110217.1  PREDICTED: calreticulin isoform 2 [Macaca...</a>	<a href="#">688</a>	0.0	
<a href="#">gb AAB20095.1  calreticulin [rabbits, skeletal muscle, Peptid...</a>	<a href="#">688</a>	0.0	
<a href="#">ref XP_002761834.1  PREDICTED: calreticulin-like [Callithrix ...]</a>	<a href="#">687</a>	0.0	
<a href="#">ref NP_004334.1  calreticulin precursor [Homo sapiens] &gt;ref X...</a>	<a href="#">682</a>	0.0	
<a href="#">ref XP_867302.1  PREDICTED: similar to calreticulin isoform 3...</a>	<a href="#">678</a>	0.0	
<a href="#">gb EAW84330.1  calreticulin, isoform CRA_a [Homo sapiens]</a>	<a href="#">673</a>	0.0	
<a href="#">ref XP_533899.2  PREDICTED: similar to calreticulin isoform 1...</a>	<a href="#">669</a>	0.0	
<a href="#">ref XP_867310.1  PREDICTED: similar to calreticulin isoform 4...</a>	<a href="#">650</a>	0.0	
<a href="#">ref XP_001377711.1  PREDICTED: similar to precursor (AA -17 t...</a>	<a href="#">645</a>	0.0	
<a href="#">ref XP_001514084.1  PREDICTED: similar to calreticulin [Ornit...</a>	<a href="#">645</a>	0.0	
<a href="#">gb AAA49610.1  calreticulin [Gallus gallus]</a>	<a href="#">645</a>	0.0	



# Literature Links

PubMed  
OMIM



# NM\_000249: PubMed

NCBI **PubMed** National Library of Medicine NLM

Search  for

About Entrez  Preview/Index History Clipboard Details

Text Version Display Abstract Show: 20 Sort Send to Text

1: **Science**. 1994 Mar 18;263(5153):1625-9.

Comment in:

- [Science. 1994 Mar 18;263\(5153\):1559-60.](#)

**Mutation of a mutL homolog in hereditary colorectal cancer**

**Papadopoulos N, Nicolaides NC, Wei YF, Ruben SM, Carter KC, Rosen CA, Haseltine WA, Fleischmann RD, Fraser CM, Adams MD, et al.**

Johns Hopkins Oncology Center, Baltimore, MD 21231.

Some cases of hereditary nonpolyposis colorectal cancer (HNPCC) are due to alterations in a mutS-related mismatch repair gene. A search of a large database of expressed sequence tags derived from random complementary DNA clones revealed three additional human mismatch repair genes, all related to the bacterial mutL gene. One of these genes (hMLH1) resides on chromosome 3p21, within 1 centimorgan of markers previously linked to cancer susceptibility in HNPCC kindreds. Mutations of hMLH1 that would disrupt the gene product were identified in such kindreds, demonstrating that this gene is responsible for the disease. These results suggest that defects in any of several mismatch repair genes can cause HNPCC.

Related Links

- Nucleotide
- Protein
- OMIM
- Cited in PMC
- Cited in Books
- Books
- LinkOut

# Books Link

1: Science. 1994 Mar 18;263(5153):1625-9.

Comment in:

- [Science. 1994 Mar 18;263\(5153\):1625-9.](#)

[Mutation of a](#)

**Papadopoulos N,  
CA, Haseltine V**

Johns Hopkins C

Some cases of [he](#)  
alterations in a [mu](#)  
[expressed sequen](#)  
revealed the [add](#)  
[mutL gene](#) one c  
centimorg [of m](#)  
kindreds. [Statistica](#)  
in such [kindreds](#),  
results suggest the  
[HNPCC](#).

The screenshot shows a section of the NCBI Handbook. At the top right, there's a sidebar with a yellow arrow pointing up to the title '18 items in The NCBI Handbook'. Below this are links to various books like 'Primer of Medicine (US)', 'C. elegans II.', and 'C. elegans I.'. The main content area has a blue header with the NCBI logo and the title 'The NCBI Handbook'. Below the header are navigation links: 'Short Contents | Full Contents' and 'Other books @ NCBI'. The main text starts with 'The NCBI Handbook → Part 3. Querying and Linking the Data'. It includes creation and update dates ('Created: October 9, 2002' and 'Updated: August 13, 2003'). A large section is dedicated to '21. UniGene: A Unified View of the Transcriptome' by Joan U. Pontius, Lukas Wagner, and Gregory D. Schuler. This section includes a summary and detailed text about UniGene's purpose and development. A yellow arrow points down from the sidebar to the '21. UniGene' link in the main text.

18 items in The NCBI Handbook

NCBI Books

National Library of Medicine (US),  
National Institutes of Health, US Gov.

C. elegans II.

Riddle, David L.; Blumenthal, Thomas; Meyer, Barbara J.; Priess, James R., editors.  
C. elegans I. Cold Spring Harbor Laboratory Press, c1997.

The NCBI Handbook

Short Contents | Full Contents

Other books @ NCBI

**Navigation**

[About this book](#)

[Part 3. Querying and Linking the Data](#)

→ [21. UniGene: A Unified View of the Transcriptome](#)

[Expressed Sequence Tags \(ESTs\)](#)

[Sequence Clusters](#)

[UniGene Cluster Browser](#)

[Protein Similarity Analysis](#)

[Digital Differential Display \(DDD\)](#)

[HomoloGene](#)

[References](#)

[Glossary](#)

**The NCBI Handbook → Part 3. Querying and Linking the Data**

Created: October 9, 2002

Updated: August 13, 2003

**21. UniGene: A Unified View of the Transcriptome**

by Joan U. Pontius, Lukas Wagner, and Gregory D. Schuler

**Summary**

The task of assembling an inventory of all genes of *Homo sapiens* and other organisms began more than a decade ago with large-scale survey sequencing of transcribed sequences. The resulting Expressed Sequence Tags (ESTs) were a gold mine of novel gene sequences that provided an infrastructure for additional large-scale projects, such as gene maps, expression systems, and full-length cDNA projects. In addition, untold numbers of targeted gene-hunting projects have benefited from the availability of these sequences and the physical clone reagents. However, the high level of redundancy found among transcribed sequences, not to mention a variety of common experimental artifacts, made it difficult for many people to make effective use of the data. This problem was the motivation for the development of [UniGene](#), a largely automated analytical system for producing an organized view of the transcriptome. In this chapter, we discuss the properties of the input sequences, the process by which they are analyzed in

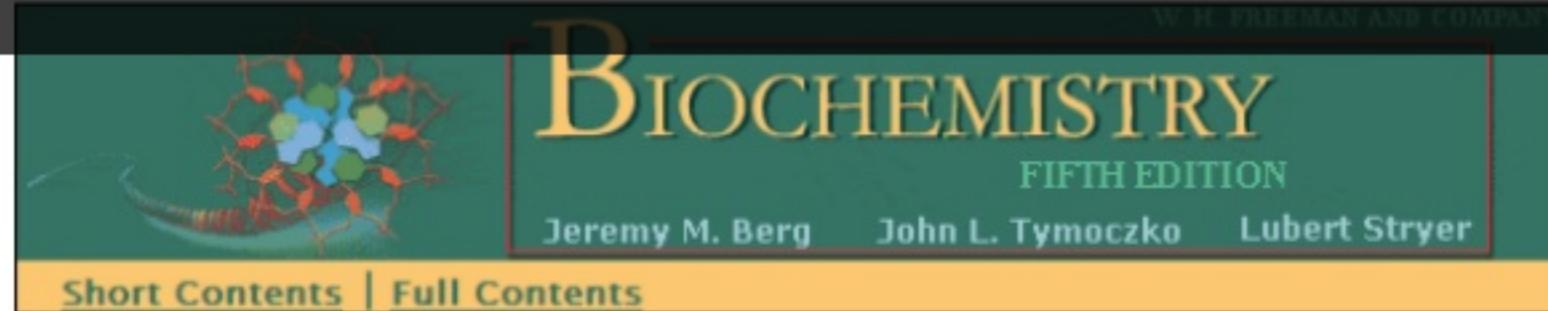
James E. M. Frei, Emil

Search

This book  All books  PubMed

GO

< 74 of 80 >



Biochemistry → III. Synthesizing the Molecules of Life → 31. The Control of Gene Expression → 31.3. Transcriptional Activation and Repression Are Mediated by Protein-Protein Interactions



**Figure 31.29. Structure of a Bromodomain.** This four-helix-bundle domain binds peptides containing acetyllysine. An acetylated



# OMIM: Human Disease Genes

**OMIM**  
Online Mendelian Inheritance in Man

Johns Hopkins University

PubMed Nucleotide Protein Genome Structure PMC Taxonomy OMIM

Search OMIM for Go Clear

Limits Preview/Index History Clipboard Details

Display Detailed Show: 20 Send to Text

**\*120436** COLON CANCER, FAMILIAL NONPOLYPOSIS, TYPE 2 Links

*Alternative titles; symbols*

FCC2  
COCA2  
COLORECTAL CANCER, HEREDITARY NONPOLYPOSIS, TYPE 2; HNPCC2  
MutL, E. COLI, HOMOLOG OF, 1, INCLUDED; MLH1, INCLUDED

Gene map locus [3p21.3](#)

**TEXT**

Using RFLPs and microsatellite markers for linkage analysis in 3 hereditary nonpolyposis colon cancer families, [Lindblom et al. \(1993\)](#) demonstrated linkage to 3p23-p21. Tumor DNA from 1 tumor in each family was included in the study to look for rearrangements related to tumor development. None of the colon tumors showed loss of heterozygosity (LOH) for any of the informative markers used on 20 different chromosomes. However, after they had detected linkage to 3p, [Lindblom et al. \(1993\)](#) observed a gain of bands for several dinucleotide markers located on 3p. A gain of bands was observed with markers on many chromosomes.

After human homologs of the mutS gene of bacteria and yeast were found to have mutations responsible for hereditary nonpolyposis colorectal cancer ([120435](#)), [Papadopoulos et al. \(1994\)](#) searched for other human mismatch repair (MMR) genes. A survey of a large database of expressed sequenced tags (ESTs) derived from random cDNA clones revealed 3 additional human MMR genes, all related to the bacterial mutL gene. [Papadopoulos et al. \(1994\)](#) mapped one of these genes (MLH1) to 3p21.3 by fluorescence in situ hybridization. The other 2 genes had a slightly greater similarity to the yeast mutL homolog PMS1 ([600258](#)) and were therefore denoted PMS1 and PMS2, respectively. The mapping of MLH1 to 3p21 was of interest because

# Taxonomy Link



# Taxonomy Link

 NCBI 

Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search for  as   lock

Go Clear

Display  levels using filter:

Nucleotide  Protein  Structure  Genome  Popset  SNP  
 3D Domains  Domains  GEO Datasets  GEO Expressions  UniGene  UniSTS  
 PubMed Central  Gene  MapView  LinkOut  BLAST  TRACE

**Lineage (full):** root; cellular organisms; Eukaryota; Fungi/Metazoa group; Metazoa; Eumetazoa; Bilateria; Coelomata; Deuterostomia; Chordata; Craniata; Vertebrata; Gnathostomata; Teleostomi; Euteleostomi; Sarcopterygii; Tetrapoda; Amniota; Mammalia; Theria; Eutheria; Primates; Catarrhini; Hominidae; Homo/Pan/Gorilla group; Homo

---

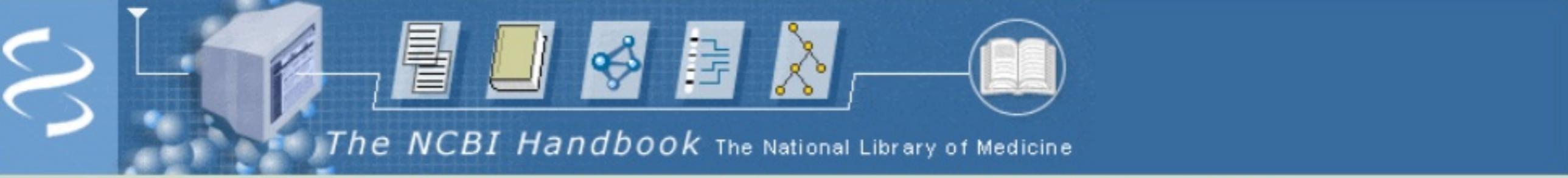
◦ **Homo sapiens** (human) 7,356,169 208,191 4,800 25 8,008 4,145,589 118,517

*Click on organism name to get more information.*

▪ **Homo sapiens neanderthalensis** 6



# For More Information...



The NCBI Handbook The National Library of Medicine

[Short Contents](#) | [Full Contents](#)      [Other books @ NCBI](#)

## The NCBI Handbook

Part 1. [The Databases](#)

 [1. GenBank: The Nucleotide Sequence Database](#)  
Ilene Mizrachi.  
Created: October 9, 2002, Updated: July 27, 2004

 [2. PubMed: The Bibliographic Database](#)  
Kathi Canese, Jennifer Jentsch, and Carol Myers.  
Created: October 9, 2002, Updated: August 13, 2003

 [3. Macromolecular Structure Databases](#)  
Eric Sayers and Steve Bryant.  
Created: October 9, 2002, Updated: August 13, 2003

 [4. The Taxonomy Project](#)  
Scott Federhen.  
Created: October 9, 2002, Updated: August 13, 2003

 [5. The Single Nucleotide Polymorphism Database \(dbSNP\) of Nucleotide Sequence Variation](#)  
Adrienne Kitts and Stephen Sherry.  
Created: October 09, 2002, Updated: June 16, 2006

 [6. The Gene Expression Omnibus \(GEO\): A Gene Expression and Hybridization Repository](#)  
Peter Ermolaev and Alon Shklar

Search

This book    All books  
 PubMed

# Thank You!

