

## Lesson 7. Content

### 1. Biological databases.

## Availability of biological data

In the early '50 Frederick Sanger determined the first complete amino acid sequence (the two polypeptide chains of bovin insulin, A and B)

After 10 years the first attempts were made to create a database of protein sequences

Nowadays the number of known protein sequences is in the order of hundreds of millions,  $10^8$  ( $\approx$ 230 million in UniProt at October 2022)

The first protein 3D structures to be solved were hemoglobin by Max Perutz and myoglobin by John Kendrew, in 1958

Nowadays the number of known protein 3D structures in the Protein Data Bank (PDB) are  $\approx$ 200,000 (October 2022)

# Relevance and structure of biological databases

Biological data and therefore databases (DBs) which contain it constitute the backbone of bioinformatics research

There are many different ways to design databases, both in terms of the ways the information is stored and the ways it can be retrieved and analysed

Modern databases store data and contain complex technology to store it in a structured manner, to allow it to be efficiently INPUT, ACCESSED and ANALYSED

Databases require databases management system (DBMS), a software to control the DB

# Relevance and structure of biological databases

Databases contain entries, made of **data** (e.g. the sequence in a sequence DB) and **annotation** (e.g. organism, gene location and name, name and sequence of the encoded protein, etc.)

ID	HSIGHAF	standard; RNA; HUM; 1089 BP.	FH	Key	Location/Qualifiers
XX			FH		
AC	J00231;		FT	source	1..1089
XX			FT		/db xref="taxon:9606"
SV	J00231.1		FT		/organism="Homo sapiens"
XX			FT		/map="14q32.33"
DT	13-JUN-1985 (Rel. 06, Created)		FT	mRNA	<1..1089
DT	02-JUL-1999 (Rel. 60, Last updated, Version 7)		FT		/note="gamma3 mRNA"
XX			FT	gene	23..964
DE	Human Ig gamma3 heavy chain disease OMM protein mRNA.		FT	CDS	/gene="IGHG3"
XX			FT		23..964
KW	C-region; gamma heavy chain disease protein;		FT		/codon start=1
KW	gamma3 heavy chain disease protein; heavy chain disease; hinge exon;		FT		/db xref="SWISS-PROT:P01860"
KW	immunoglobulin gamma-chain; immunoglobulin heavy chain;		FT		/note="OMM protein (Ig gamma3) heavy chain"
KW	secreted immunoglobulin; V-region.		FT		/gene="IGHG3"
XX			FT		/protein id="AAA52805.1"
OS	Homo sapiens (human)		FT		/translation="MKXLMFFLLVAAAPRWVLSQVHLQESGFGGLKFPPELK
OC	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Mammalia; Eutheri		FT		TCPRCFEPKSCDTPPPCPRCFEPKSCDTPPPCPRCFEPKSCDTPPPCPKCP
OC	Primates; Catarrhini; Hominidae; Homo.		FT		SVFLFFPKPKDTLMISRTPEVTCVVVDVSHEDPKVQFKNYVDGVEVHNAKT
VV			FT		STFRVSVSLTVLHQDWLNGKEYKCKVSNKALPAPIEKTIISKAKGQFPXXXXX
			FT		EMTKNQVSLTCLVKGFYPSDIAVEWESNGQPENNYNTTPPMLDSDGSFFLY
			FT		RMQQGNIFSCVMHEALHNRVYTKSLSLSPGK*

Feature Table

DATA

```
50 Sequence 1089 BP; 240 A; 358 C; 271 G; 176 T; 44 other;
cctggacctc ctgtgcaaga acatgaacaa nctgtgggtc ttcttctcc tgggtggcagc
60
tccagatgg gtctgtctcc aggtgcacct gcaggagtgc ggcacaggac tggggaagcc
120
tccagagtc aaaacccccc ttggtgacac aactcacaca tgcacacggt gccacagacc
180
caaatcttgt gacacacctc ccccggtgcc acggtgccca gagcccaaat ctgtgtacac
240
acctccccc tgcacacggt gccacagacc caaatcttgt gacacacctc ccccggtgcc
300
nnngtgccca gacacctgaac tcttgggagg acggtcagtc ttctcttcc ccccaaaacc
360
caaggatacc cttatgattt cccggacccc tgaggtcaag tgcgtggtgg tggacgtgag
420
ccacgaagac cccnnngtcc agttcaagtg gtaogtgga ggcgtggagg tgcataatgc
480
caagacaaag ctgcgggagg agcagtagaa cagcacgttc cgtgtggtca ggcgtctcac
540
cgtcctgcac caggactggc tgaacggcaa ggagtacaag tgcaaggctc ccaacaaagc
600
cctccacgcc cccatcgaga aaaccatctc caaagccaaa ggacagcccn nnnnnnnnnn
660
nnnnnnnnnn nnnnnnnnnn nnnnngagga gatgaccaag aaccaagta gctgacctg
720
cctggtcaaa ggccttctac ccagcgacat cgcgtggag tgggagagca atgggcagcc
780
ggagaacaac tacaacacca cgcctcccat gctggaactc gacggtcctc tcttctctc
840
cagcaagctc accgtggaca agagcagggt gcagcagggg aacatcttct catgctcgt
```

ANNOTATION

*Example of entry from a nucleotide sequence DB (GenBank)*

# The structure of databases: flat-file format

In general a database structure consists of files or tables each containing **records** and **fields**

(A)

NAME	TELEPHONE	ADDRESS
S. Claus	0203 450	The North Pole, Lapland
M. Mouse	0202 453	Disneyworld, Florida
A. Moonman	0104 459	Craterland, The Moon

fields      *Example of a very simple database table: a single page with a contact list with 3 records and 3 fields per record*

# The structure of databases: flat-file format

In general a database structure consists of files or tables each containing **records** and **fields**

(A)

NAME	TELEPHONE	ADDRESS
S. Claus	0203 450	The North Pole, Lapland
M. Mouse	0202 453	Disneyworld, Florida
A. Moonman	0104 459	Craterland, The Moon

records

*Example of a very simple database table: a single page with a contact list with 3 records and 3 fields per record*

# The structure of databases: flat-file format

In general a database structure consists of files or tables each containing **records** and **fields**

## (B) GenBank Flat-File Format

LOCUS	SCU49845	5028 bp	DNA
DEFINITION	Saccharomyces cerevisiae TCP1-beta gene, partial cds, and Axl2p (AXL2) and Rev7p (REV7) genes, complete cds.		
ACCESSION	U49845		
VERSION	U49845.1 GI:1293613		
KEYWORDS	.		
SOURCE	Saccharomyces cerevisiae (baker's yeast)		
ORGANISM	Saccharomyces cerevisiae Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes; Saccharomycetales; Saccharomycetaceae; Saccharomyces.		

fields

*Example of a  
GenBank record in a  
flat-file format*

# The structure of databases: relational databases

In relational DBs, the most commonly used for biological information, data is stored within a number of tables

Each table consists of **records** and **fields** and is linked to other tables by a shared field called **key**, unique to each record

protab1			
Protein-code	Protein-name	Length	Species-origin
P1001	Hemoglobin	145	Bovine
P1002	Hemoglobin	136	Ovine
P1003	Eye Lens Protein	234	Human
.....			

fields

records

protab2	
Protein-code	Protein-sequence
P1001	MDRTHGFDLKLSPRTVNQWLMLALFFGHS...
P1002	MDKTSHGFEIKLLTPKKLQQWLMIAIYFGHT...
P1003	SRTHEEEGKLMQWPPRPLYIALFTEPPYP...
.....	

fields

records

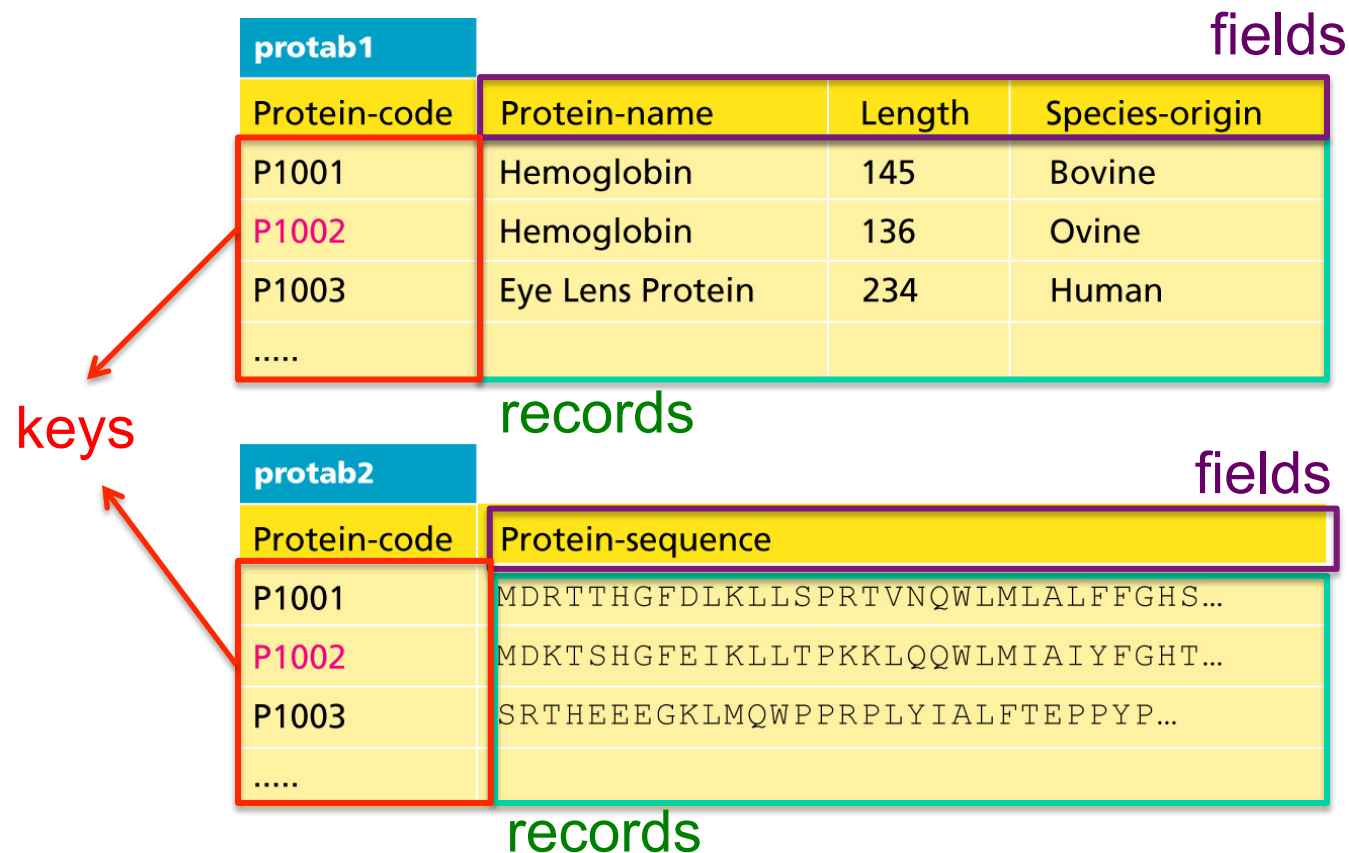
*This example of relational database consists of two tables with different fields, where records are connected by an identifier (**key**)*



# The structure of databases: relational databases

In relational DBs, the most commonly used for biological information, data is stored within a number of tables

Each table consists of **records** and **fields** and is linked to other tables by a shared field called **key**, unique to each record



*This example of relational database consists of two tables with different fields, where records are connected by an identifier (**key**)*

# The structure of databases: relational databases

In relational DBs, a set of operators is provided that allows to manipulate and analyse the data (mathematical, logical etc.)

Usually, results produced by application of operators on the data are displayed in new tables

The diagram illustrates the difference between a flat file and a database using two tables, **protab1** and **protab2**.

**protab1** (Flat File):

Protein-code	Protein-name	Length	Species-origin
P1001	Hemoglobin	145	Bovine
P1002	Hemoglobin	136	Ovine
P1003	Eye Lens Protein	234	Human
....			

**protab2** (Database):

Protein-code	Protein-sequence
P1001	MDRTHGFDLKLSPRTVNQWLMLALFFGHS...
P1002	MDKTSHGFEIKLLTPKKLQQWLMIAIYFGHT...
P1003	SRTHEEEGKLMQWP RPPLYIALFTEPPYP...
....	

**Annotations:**


- keys:** Indicated by red arrows pointing to the **Protein-code** column in both tables.
- records:** Indicated by green arrows pointing to the rows in both tables.
- fields:** Indicated by purple arrows pointing to the column headers in both tables.

*This example of relational database consists of two tables with different fields, where records are connected by an identifier (**key**)*

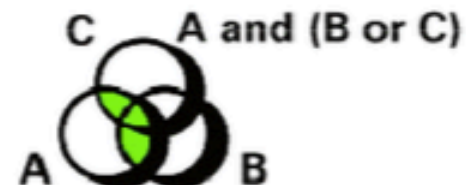
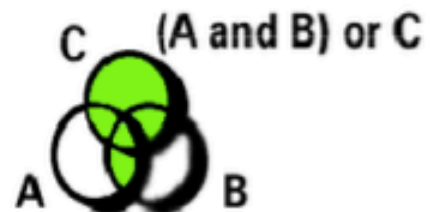
# Searching DBs

Boolean operators help to focalize DB searches:

## Boolean Operators

- AND (&)  A and B
- OR (|)  A or B
- NOT(!)  A and not B

## Complex queries



# Searching DBs

Boolean operators help to focalize DB searches:

*e.g. “human” will list all the entries containing any information on homo sapiens*

*“human liver alcohol dehydrogenase” (human **AND** liver **AND** alcohol **AND** dehydrogenase) will list only entries specific for the enzyme*

Databases usually provide **fields** which can be selected in order to focalize a search, e.g. author ([au]), organism [OS], etc

# The structure of databases: SQL

In relational DBs management systems, operators are written in query-specific languages, such as the **Structured Query Language (SQL)**

Examples of SQL application on a relational database

*Query 1*

*example-1*

```
SELECT protein-code, protein-name  
FROM protab1  
WHERE species-origin = 'Bovine';
```

P1001	Hemoglobin
-------	------------

*Query 2*

*example-2*

```
SELECT protab1.protein-name, protab2.protein-sequence  
FROM protab1, protab2  
WHERE protab1.protein-code = protab2.protein-code  
AND protab1.protein-code = 'P1002';
```

Hemoglobin	MDKTSHGFEIKLLTPKKLQQWLMIAIYFGHT...
------------	------------------------------------

# The structure of databases: XML

**eXtensible Markup Language (XML)** is a powerful system for marking up (annotating) data

It is one of many markup languages, including HTML (hypertext markup language), commonly used to write web pages

The hallmark of these languages is the use of identifiers called tags, which can enclose sections of data, e.g. `<language>XML</language>`

XML has mechanisms allowing arbitrary tags to be used, thus having the flexibility to define bespoke (tailored) data classifications

XML uses plain file format, which makes it portable and accessible

Many bioinformatics DBs are being made available in XML format, although a master copy is often maintained e.g. as relational DBs

# Primary and secondary (derived) biological DBs

PRIMARY DBs: contain primary data, i.e. experimental data & annotations

e.g. *nucleotide sequence* *coding, non coding*

SECONDARY DBs: contain secondary data, i.e. *analyses* of general interest derived from primary DBs & annotations

e.g. *collections of conserved sequence motifs* *functional, non functional*

A number of centers have been funded to provide access to a large number of major databases in an integrated environment

# More on DNA sequences in DBs

There are different types of DNA sequences

## **Genomic (chromosomal) sequences:**

*from genome sequencing projects – in GenBank*

include noncoding regions, introns, control regions etc.

## **cDNA (complementary DNA) sequences:**

*obtained from mRNA sequences by reverse transcription (copying mRNA in DNA)*

represent genes actually expressed in a specific cell/tissue at a given condition

do not include gene regions not transcribed in mRNA (introns, control regions etc.)

## **ESTs (expressed sequence tags):**

*obtained from mRNA sequences by reverse transcription (copying mRNA in DNA)*

are partial cDNA sequences, representing expressed genes – useful for genomes scanning



# Examples of primary biological databases

..aatgcatgccaatg  
ccatccgcatcgat..

EMBL

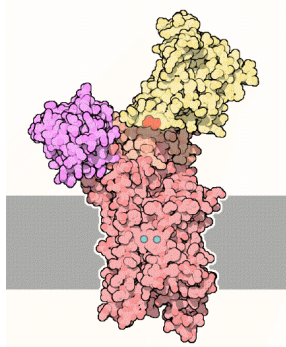
Nucleotide Sequence Database

GenBank

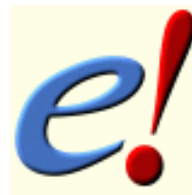
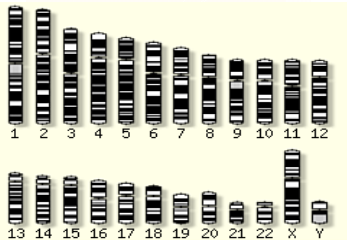
..QEDARTSCG  
AILNQRYWI..

UniProt/Swiss-Prot

Protein Information Resource (**PIR**)



RCSB **PDB**  
PROTEIN DATA BANK



ensembl



PubMed

# Examples of primary biological databases: GenBank

The screenshot displays the NCBI (National Center for Biotechnology Information) homepage as viewed in Microsoft Internet Explorer. The browser's address bar shows the URL <http://www.ncbi.nlm.nih.gov/>. The page features a navigation menu with links to PubMed, All Databases, BLAST, OMIM, Books, TaxBrowser, and Structure. A search bar is prominently displayed, with the word "Gene" entered and a "Go" button. The page is organized into several sections: a left sidebar with a "SITE MAP" and links to "About NCBI", "GenBank", "Literature databases", and "Molecular databases"; a main content area with a "What does NCBI do?" section describing the center's mission and an "Influenza Virus Resource" section; and a right sidebar titled "Hot Spots" listing various resources like the Assembly Archive, Clusters of orthologous groups, and Entrez Home. The bottom of the page shows the "Entrez Gene" link and the status bar indicating "Internet".

NCBI HomePage - Microsoft Internet Explorer

File Edit View Favorites Tools Help

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

Genetics Biology DNA Biotech Health Sci-Fi

Search... Hotbar Meet

msn Search Highlight Options Pop-ups Blocked (97) Hotmail Messenger My MSN

**NCBI** National Center for Biotechnology Information  
National Library of Medicine National Institutes of Health

PubMed All Databases BLAST OMIM Books TaxBrowser Structure

Search Gene for Go

**SITE MAP**  
Alphabetical List  
Resource Guide

**About NCBI**  
An introduction to NCBI

**GenBank**  
Sequence submission support and software

**Literature databases**  
PubMed, OMIM, Books, and PubMed Central

**Molecular databases**  
Sequences, structures, and taxonomy

**What does NCBI do?**

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

**Influenza Virus Resource**

The Influenza Virus Resource enables comparison of influenza virus strains and provides a reference for viral sequences. The resource contains data from the NIAID Influenza Genome Sequencing Project and GenBank, as well as pre-computed alignments of flu sequences.

**Hot Spots**

- Assembly Archive
- Clusters of orthologous groups
- Coffee Break, Genes & Disease, NCBI Handbook
- Electronic PCR
- Entrez Home
- Entrez Tools
- Gene expression omnibus (GEO)
- Human genome resources

**Entrez Gene**

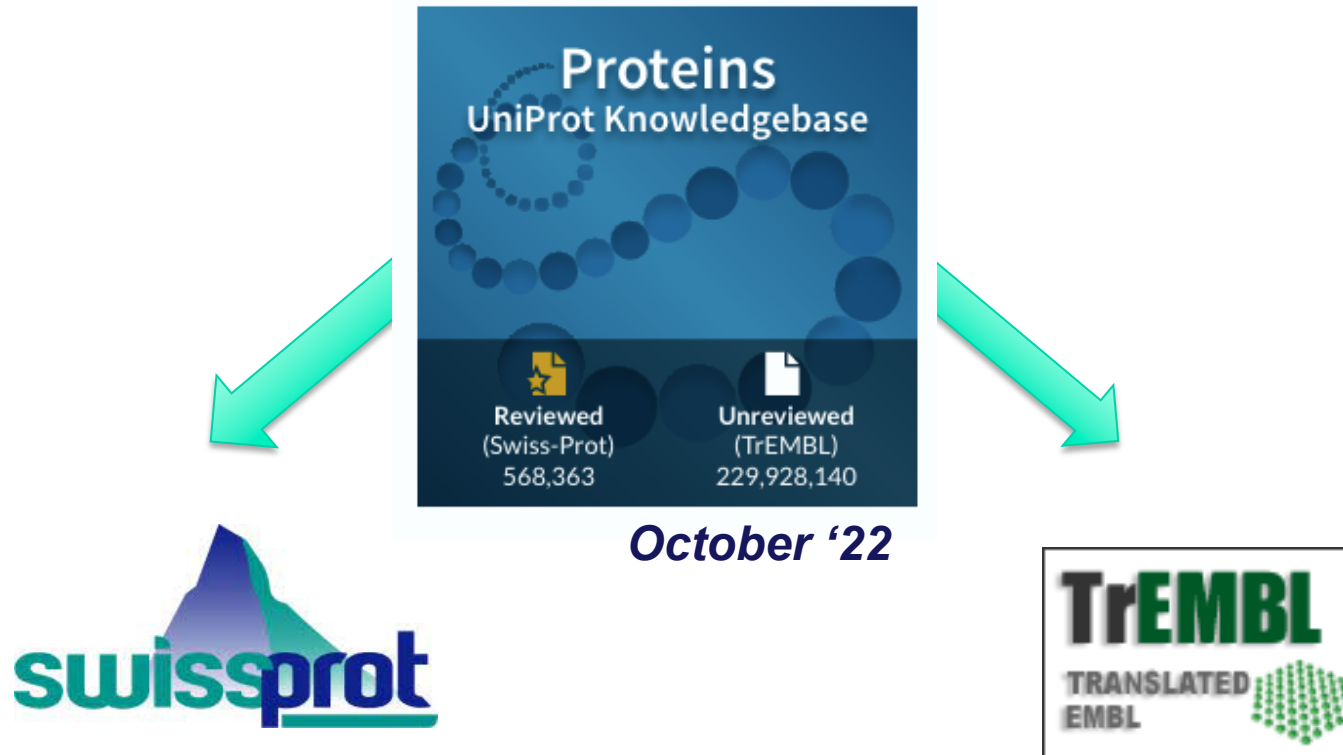
Internet

# Examples of primary biological databases: UniProt

The screenshot shows the UniProt Knowledgebase website in a Microsoft Internet Explorer browser window. The title bar reads "Text Search UniProt Knowledgebase - UniProt [the Universal Protein Resource] - Microsoft Internet Explorer". The address bar shows the URL "http://www.ebi.uniprot.org/uniprot-srv/index.do". The browser's toolbar includes buttons for Back, Forward, Stop, Home, Search, Favorites, and other standard functions. The website's header features the UniProt logo, which is a stylized "UniProt" with the tagline "the universal protein knowledgebase" below it. To the right of the logo, there is a navigation bar with links to Home, Database, and UniProt Knowledgebase Search. Below the navigation bar, there is a search bar with the text "Text Search UniProt Knowledgebase" and a search button. The main content area is titled "Text Search UniProt Knowledgebase" and contains a search form with a text input field labeled "Insert a query". Below the input field, it displays "1929823 Entries In UniProt Release 5.1" and a link to "Help On Your Query". There are also buttons for "Search & View" and "Reset". At the bottom of the main content area, there is a box with the text "A common question: How can I build queries with logical operators?". The left sidebar contains a list of links: Text Search, Power Search, Warehouse, Prediction Search, InterPro Search, CluSTR Search, Entry List Search, Data Set Manager, BLAST, FAQ, Help Desk, and Download. The browser's status bar at the bottom shows "Internet".

# Examples of primary biological databases: UniProt

UniProt consists of two parts:



*1 entry per protein*

*Non redundant, high-quality  
manual annotation -  
reviewed*

*1 entry per nucleotide  
submission*

*Redundant, automatically  
annotated –  
unreviewed*



# UniProt

```
ID   Q9XSK1      PRELIMINARY;      PRT;   142 AA.
AC   Q9XSK1;
DT   01-NOV-1999 (TrEMBLrel. 12, Created)
DT   01-NOV-1999 (TrEMBLrel. 12, Last sequence update)
DT   01-NOV-1999 (TrEMBLrel. 12, Last annotation update)
DE   II ALPHA4 HAEMOGLOBIN CHAIN.
GN   II ALPHA4.
OS   Bubalus bubalis (Water buffalo).
OC   Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Mammalia;
OC   Eutheria; Cetartiodactyla; Ruminantia; Pecora; Bovoidea; Bovidae;
OC   Bovinae; Bubalus.
RN   [1]
RP   SEQUENCE FROM N.A.
RA   Ferranti P., Rullo R., Zappacosta F., Vincenti D., Masala B.,
RA   Di Luccia A.;
RT   "River buffalo (Bubalus bubalis L.) alpha globin gene and chain
RT   sequences: evolutionary structural relationship among some ruminant
RT   species.";
RL   Submitted (MAY-1999) to the EMBL/GenBank/DDBJ databases.
DR   EMBL; AJ242734; CAB43765.1; -.
SQ   SEQUENCE   142 AA;  15171 MW;  99F24F6E285C758F CRC64;
      MVLSAADKSN VKAAWGKVGG HAADYGAEAL ERMFLSFPTT KTYFPHFDLS HGSAQVKGHG
      AKVANALTKA VGHLDDLPGA LSELSDLHAH KLRVDPVNFK LLSHSLLVTL ASHLPNDFTP
      AVHASLDKFL ASVSTVLTSK YR
//
```

## Sequence in FASTA format

**>Q9XSK1**

**MVLSAADKSNVKAAGWKVGGHAADYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG  
AKVANALTKAVGHLDDLPGALSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPNDFTP  
AVHASLDKFLASVSTVLTSKYR**

# UniProt

For a protein  
sequence

Annotation may  
also report:

regions playing a  
function, interacting  
with other  
molecules, featuring  
a certain secondary  
structure, involved in  
disulfide bridges,  
*etc. etc.*

```
ID   Q9XSK1      PRELIMINARY;      PRT;   142 AA.
AC   Q9XSK1;
DT   01-NOV-1999 (TrEMBLrel. 12, Created)
DT   01-NOV-1999 (TrEMBLrel. 12, Last sequence update)
DT   01-NOV-1999 (TrEMBLrel. 12, Last annotation update)
DE   II ALPHA4 HAEMOGLOBIN CHAIN.
GN   II ALPHA4.
OS   Bubalus bubalis (Water buffalo).
OC   Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Mammalia;
OC   Eutheria; Cetartiodactyla; Ruminantia; Pecora; Bovoidea; Bovidae;
OC   Bovinae; Bubalus.
RN   [1]
RP   SEQUENCE FROM N.A.
RA   Ferranti P., Rullo R., Zappacosta F., Vincenti D., Masala B.,
RA   Di Luccia A.;
RT   "River buffalo (Bubalus bubalis L.) alpha globin gene and chain
RT   sequences: evolutionary structural relationship among some ruminant
RT   species.";
RL   Submitted (MAY-1999) to the EMBL/GenBank/DDBJ databases.
DR   EMBL; AJ242734; CAB43765.1; -.
SQ   SEQUENCE   142 AA;  15171 MW;  99F24F6E285C758F CRC64;
      MVLSAADKSN VKAAWGKVG GHAADYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG
      AKVANALTKA VGHLDDLPGLSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPNDFTEP
      AVHASLDKFLASVSTVLTSKYR
//
```

## Sequence in FASTA format

>Q9XSK1

MVLSAADKSNVKAAWGKVG GHAADYGAEALERMFLSFPTTKTYFPHFDLSHGSAQVKGHG  
AKVANALTKAVGHLDDLPGLSELSDLHAHKLRVDPVNFKLLSHSLLVTLASHLPNDFTEP  
AVHASLDKFLASVSTVLTSKYR

- Home
- Tutorial About This Site
- Getting Started
- Download Files
- Deposit and Validate
- Structural Genomics
- Dictionaries & File Formats
- Software Tools
- Educational Resources
- General Information
- Acknowledgements
- Frequently Asked Questions
- Known Problems
- Report Bugs/Comments

## Welcome to the RCSB PDB

The RCSB PDB provides a variety of tools and resources for studying the structures of biological macromolecules and their relationships to sequence, function, and disease.

The RCSB is a member of the wwPDB whose mission is to ensure that the PDB archive remains an international resource with uniform data.

This site offers tools for browsing, searching, and reporting that utilize the data resulting from ongoing efforts to create a more consistent and comprehensive archive.

Information about compatible browsers can be found [here](#).

A [narrated tutorial](#) illustrates how to search, navigate, browse, generate reports and visualize structures using this new site. [This requires the Macromedia [Flash player download](#).]

Comments? [info@rcsb.org](mailto:info@rcsb.org)

### Molecule of the Month: Tissue Factor



Blood performs many essential jobs in your body: it transports oxygen and nutrients, it protects your cells from infection, and it carries hormones and other messages from place to place in your body. But since blood is a liquid that is pumped under pressure, we must protect

### NEWS

- Complete News
- Newsletter
- Discussion Forum

07-Mar-2006

#### RCSB PDB Focus: Frequently Asked Questions

The [Frequently Asked Questions](#) page answers a number of common queries about the new RCSB PDB site

- Full Story ...

28-Feb-2006

#### RCSB PDB Exhibit News

21-Feb-2006

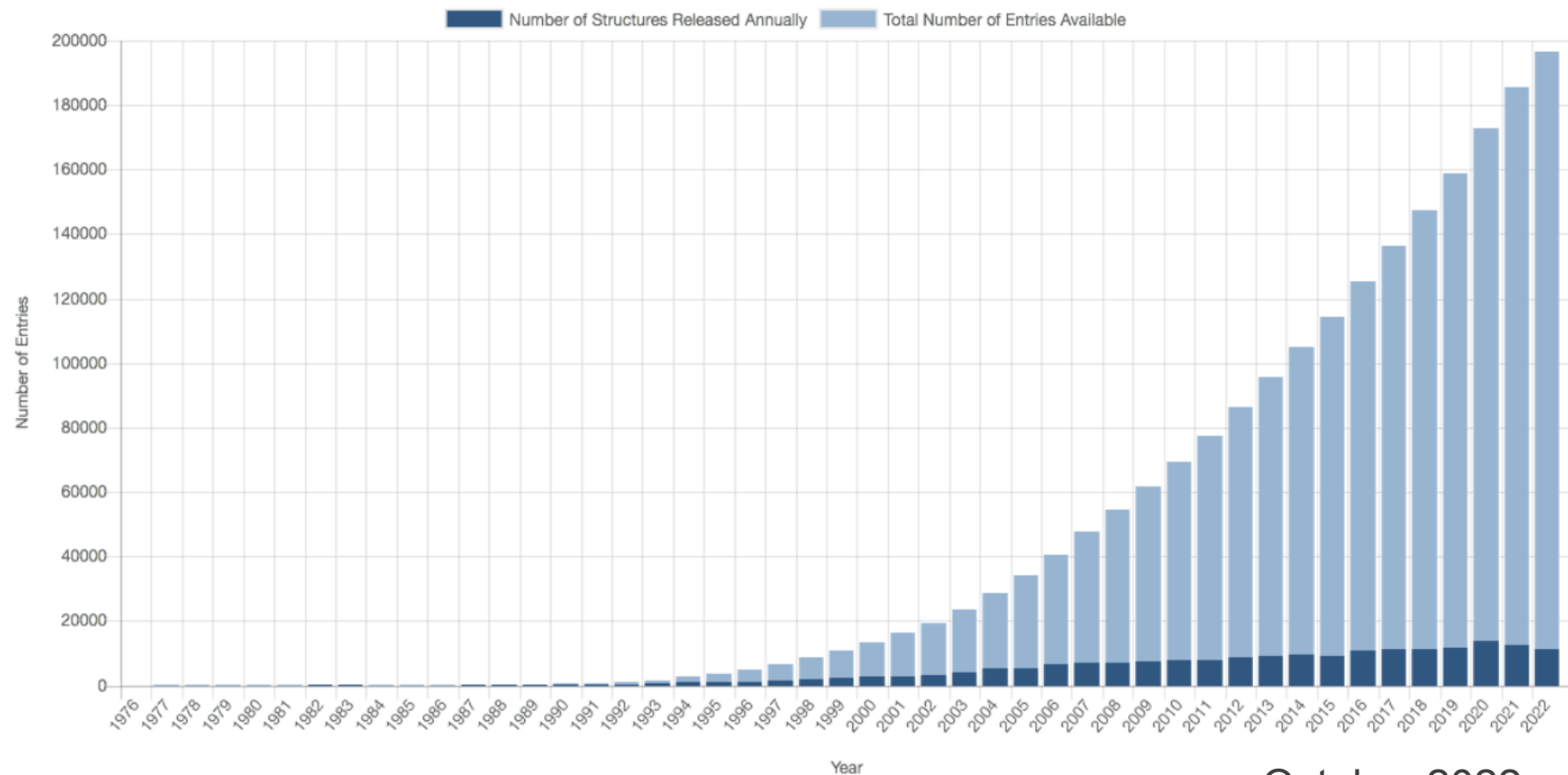
#### Virtual Reality Environment Highlights PDB Structures

# Some statistics on the Protein Data Bank (PDB)

The Protein Data Bank has been established in 1971. Since then, the number of available 3D structures has grown exponentially

PDB Statistics: Overall Growth of Released Structures Per Year

All Statistics



on October 2022

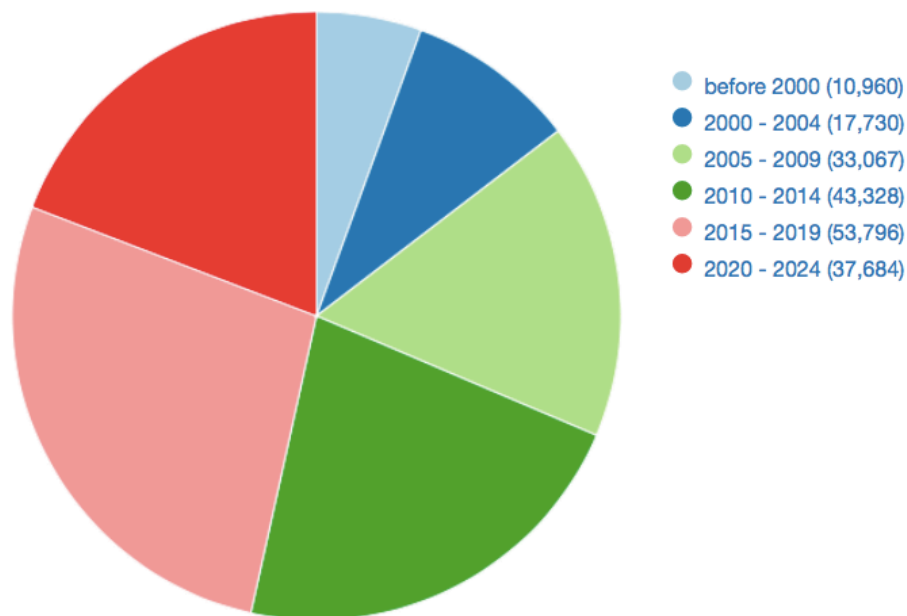


# Some statistics on the Protein Data Bank (PDB)

The Protein Data Bank has been established in 1971

## By Release Date

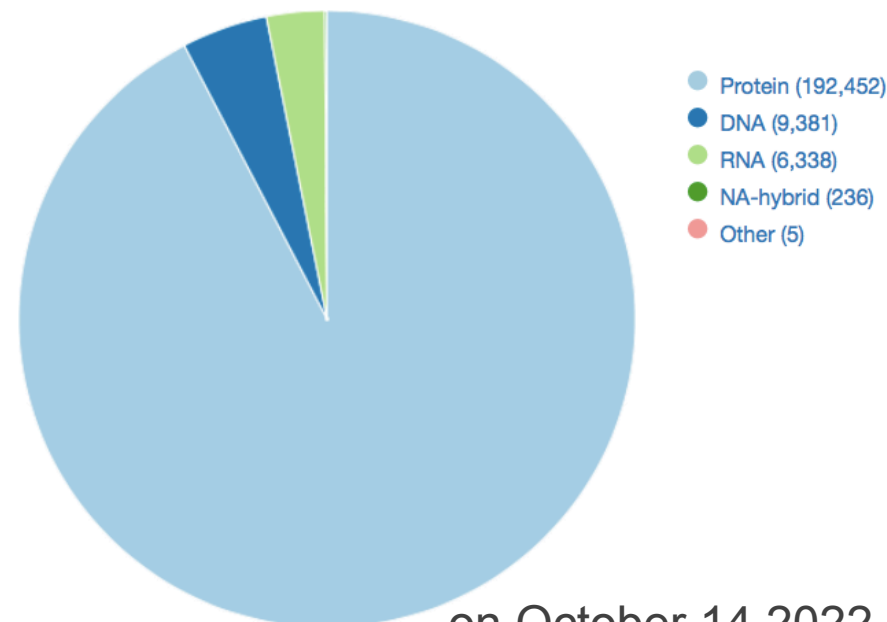
The year a structure entry was released in the PDB archive.



The 3D structures released in the last 3 years (2020-2022) are roughly 4 times those released in the first 30 years (1971-2000)

## By Polymer Entity Type

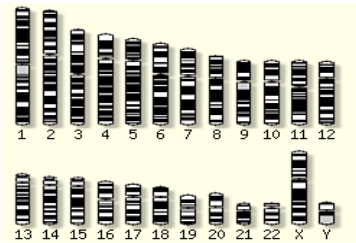
PDB structures may contain multiple entities of different macromolecular types.



on October 14 2022

The 3D structures in the PDB are mostly for proteins but also for DNA and RNA molecules

# Examples of primary biological databases: Ensembl



**Universal source of information on the human genome and other eukaryote genomes (over 200 species to date)**

**DATA:** genomes: genes, SNPs (Single Nucleotide Polymorphisms), repetitive sequences

**ANNOTATIONS:** coding regions, for each gene proofs supporting its identification

**CONNECTIONS:** to other DBs and software

## Examples of secondary (derived) biological databases

<i>Name</i>	<i>Derived from</i>	<i>Contains</i>
EMEST	EST Database	Collection and alignment of EST seqs
DSSP	Protein Data Bank (PDB)	Protein secondary structure assignments
HSSP	Protein Data Bank (PDB) & protein sequences DBs	Alignments of protein seqs of known structure with all similar seqs
FSSP, SCOP, CATH	Protein Data Bank (PDB)	Structural classification
3Dee	Protein Data Bank (PDB)	Protein domains definition
Pfam, Prints, BLOCKS	Protein sequences DBs	Alignments of homologous protein families & domains
Prodom	Protein sequences DBs	Alignments of similar protein domains
PROSITE	Protein sequences DBs	Sequence patterns
OMIM	Genomic databases	Genes & associated genetic diseases
LocusLink	Genomic databases	Genetic loci

## Protein domains

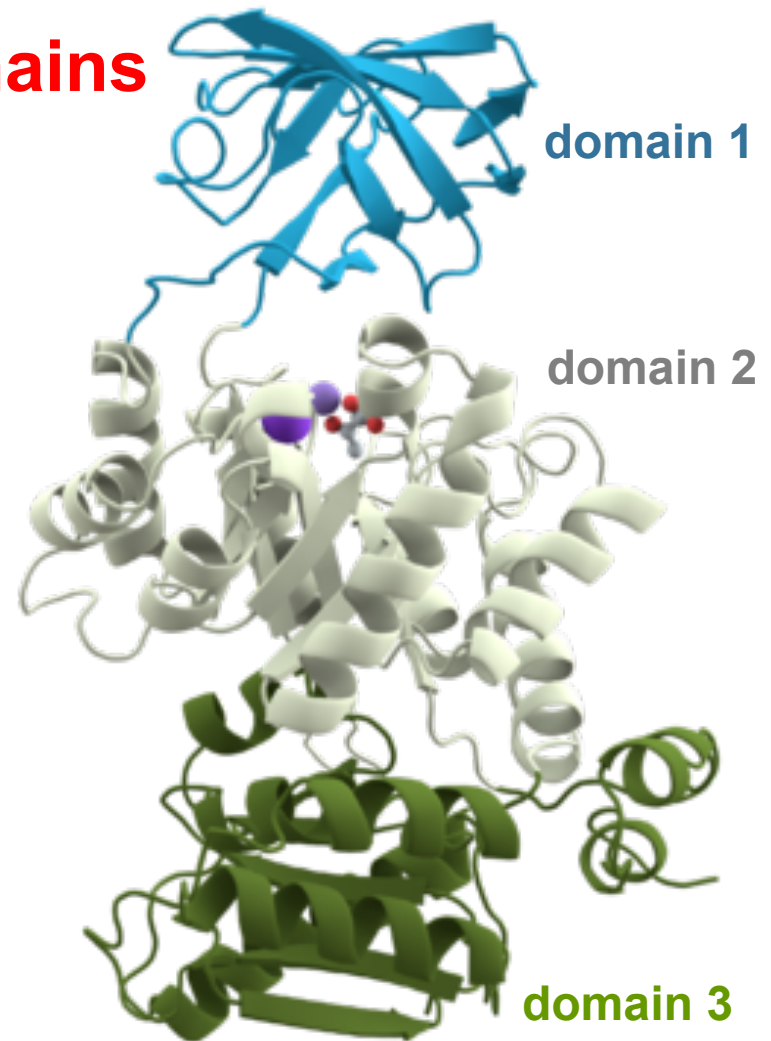
Many proteins consist of several **domains**, which usually form **functional units**

A protein domain is a region of a protein chain that is self-stabilizing and that **folds independently from the rest**

A protein domain is usually  $\approx 50$  to  $\approx 250$  amino acid long

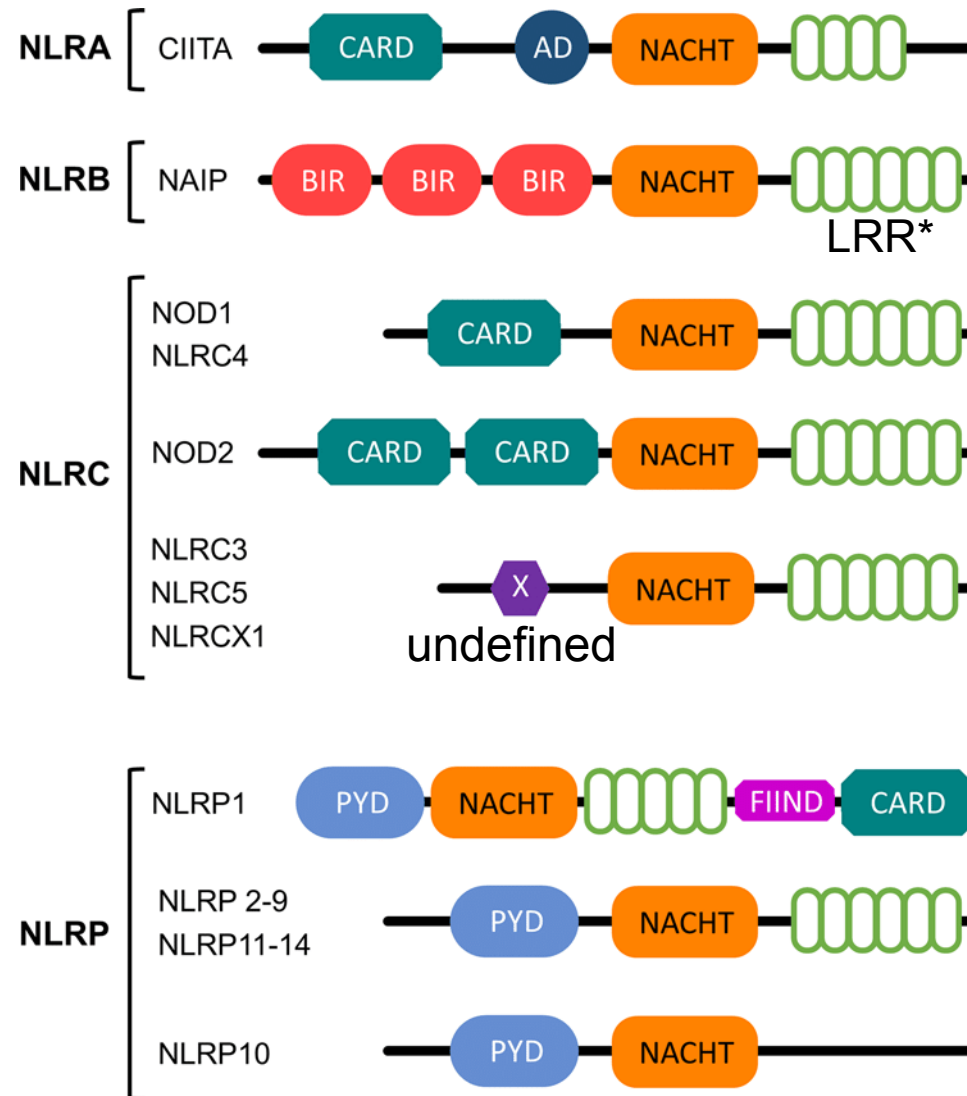
A domain may appear in a variety of different proteins

Molecular evolution uses domains as building blocks in different proteins



***Pyruvate kinase:** contains an all- $\beta$  nucleotide binding domain (blue), an  $\alpha/\beta$ -substrate binding domain (grey) and an  $\alpha/\beta$ -regulatory domain (green), connected by several linkers.*

# Protein domains as building blocks for different proteins



\*Leucine-rich repeats

*Example of domains combination in the **NOD-Like Receptors** proteins*

# Examples of secondary (derived) biological databases

**Pfam: Pfam Home Page - Microsoft Internet Explorer**

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print

Address <http://www.sanger.ac.uk/Software/Pfam/> Go Links

Search X  
Ne >>

Search The Web

Find a Web page containing

Search

Brought to you by MSN Search

Search for other items:  
[Files or Folders](#)  
[Computer](#)  
[People](#)

© 2006 Microsoft  
MSN  
Privacy

**wellcome trust sanger institute**

**Pfam**


RSS Pfam Home Search by Browse by FTP iPfam Help About

Pfam is a large collection of multiple sequence alignments and hidden Markov models covering many common protein domains and families. For each family in Pfam you can:

- Look at multiple alignments
- View protein domain architectures
- Examine species distribution
- Follow links to other databases
- View known protein structures

For more information on Pfam, on using this site, or on the changes between Pfam releases 18.0 and 19.0, click [here](#).

Pfam can be used to view the domain organisation of proteins. A typical example is shown below. Notice that a single protein can belong to several Pfam families.




75% of protein sequences have at least one match to Pfam. This number is called the sequence coverage and is shown in the pie chart on the right.

Pfam is a database of two parts, the first is the curated part of Pfam containing over 7973 protein families. To give Pfam a more comprehensive coverage of known proteins we automatically generate a supplement called Pfam-B. This contains a large number of small families taken from the [PRODOM](#) database that do not overlap with Pfam-A. Although of lower quality Pfam-B families can be useful when no Pfam-A families are found.

**Version 19.0**

December 2005, **8183** families



■ Sequence coverage Pfam-A : 75%  
■ Sequence coverage Pfam-B : 19%  
■ Other

**Web feed**

You can use the RSS feed to keep updated about Pfam releases  
[XML](#) [RSS](#)

**Enter your keyword(s) here**

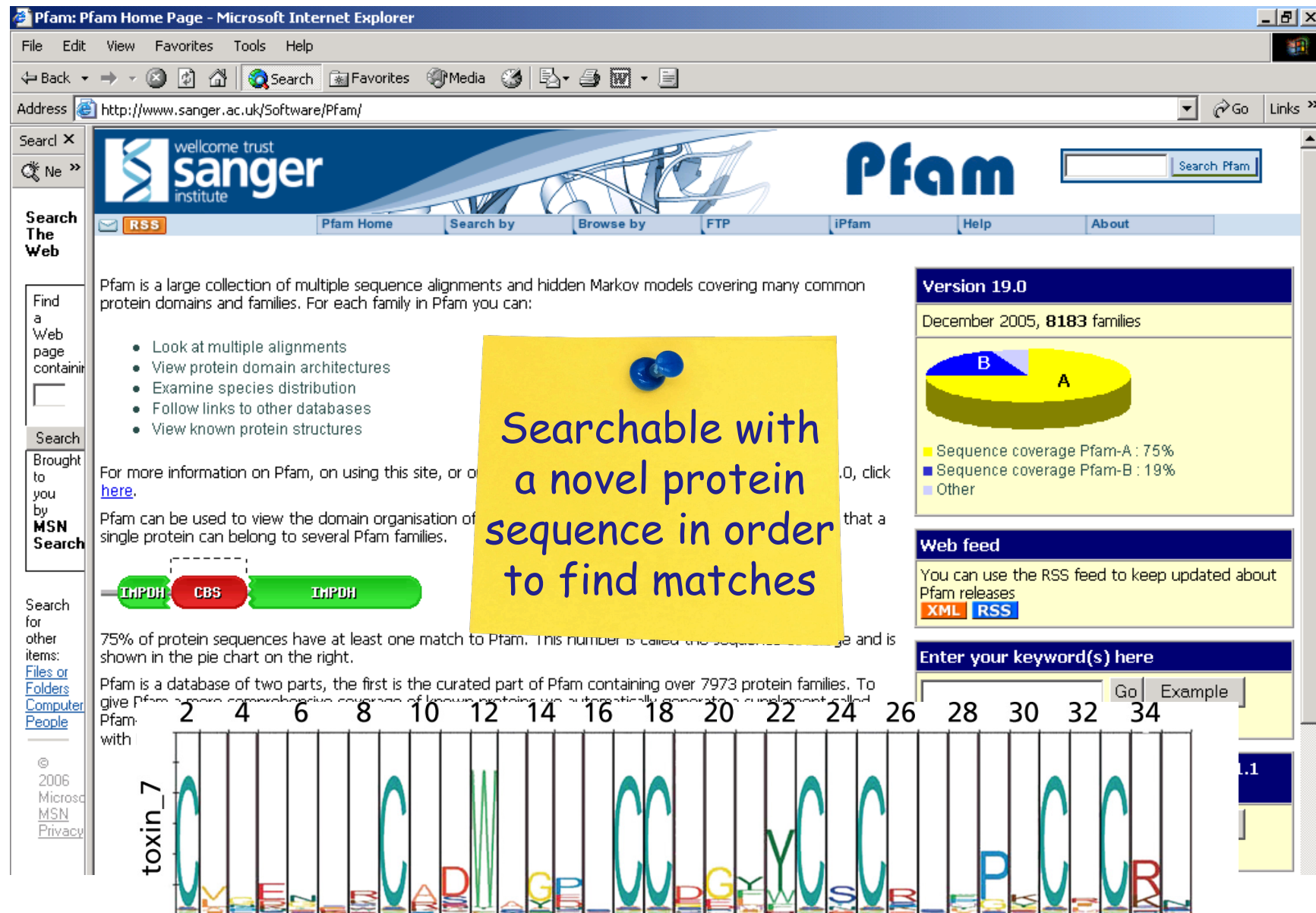
Go Example

**Enter a SWISS-PROT 48.1 or TrEMBL 31.1 name or accession number**

Go Example

Pfam is a collection of protein families and domains  
Contains multiple alignments and HMM profiles

# Examples of secondary (derived) biological databases



*Example of a HMM profile for the Toxin\_7 family, from the Pfam database*



# Examples of secondary (derived) biological databases

ExPASy - PROSITE - Microsoft Internet Explorer

File Modifica Visualizza Preferiti Strumenti ?

Indirizzo <http://www.expasy.ch/prosite/>

Google  Cerca

Segnalibri 23 bloccati Controllo Traduci Invia a

ExPASy Home page Site Map Search ExPASy Contact us Swiss-Prot ENZYME

Search  for  Go Clear

**proSite** Database of protein domains, families and functional sites

Home ScanProsite ProRule Documents Downloads Links

PROSITE consists of [documentation entries](#) describing protein domains, families and functional sites as well as associated [patterns](#) and [profiles](#) to identify them [[More details](#) / [References](#) / [Disclaimer](#) / [Commercial users](#)].

PROSITE is complemented by [new](#) **ProRule**, a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [[More details](#)].

Release 20.9, of 03-Apr-2007 (1478 documentation entries, 1322 patterns, 720 profiles and 744 ProRule)

**PROSITE access**

Search

☐ add wildcard '\*'

Browse:

- by [documentation entry](#)
- by [ProRule description](#) [new](#)
- by [taxonomic scope](#) [new](#)
- by [number of positive hit](#) [new](#)

SRS - Sequence Retrieval System

**PROSITE tools**

Scan a sequence against PROSITE - quick scan  
Enter your sequence or a [UniProtKB](#) ([Swiss-Prot](#) or [TrEMBL](#)) ID or AC [[help](#)]:

Output includes graphical view and feature detection.

Scan Clear

FZ KRINGLE\_2 PROTEIN\_KINASE\_DOM

**proSite** is a database of protein families, domains and functional sites



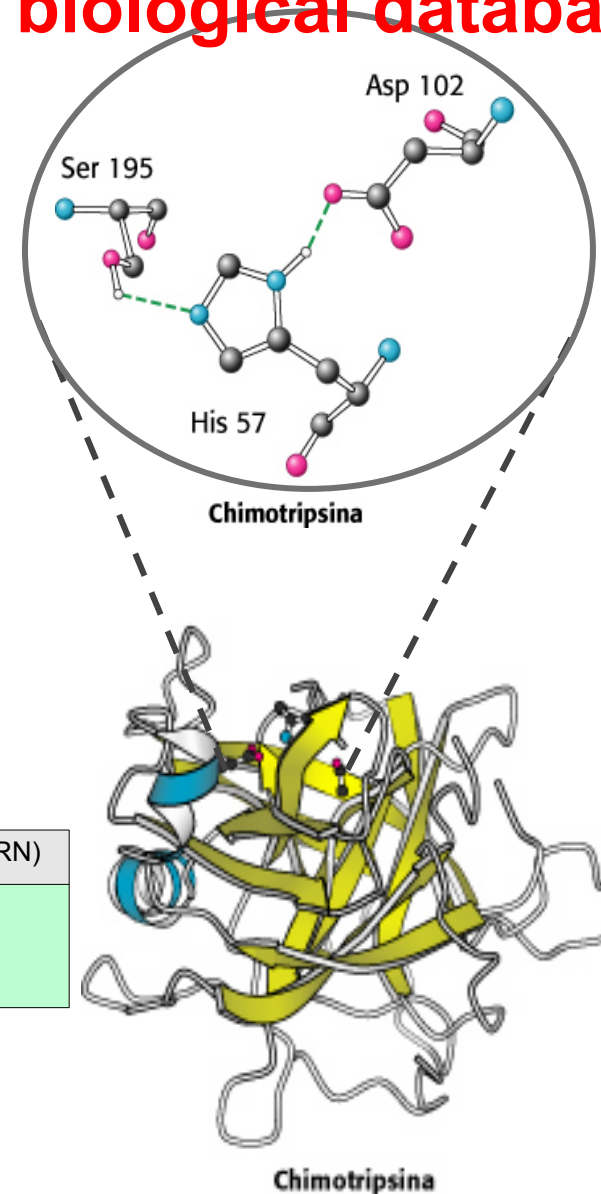
# Examples of secondary (derived) biological databases

Example of a  pattern

Histidine (**H**) and Serine (**S**) belonging to the active site of trypsin-like Ser-protease enzymes appear in sequence patterns of the type reported below

TRYPSIN_HIS, <a href="#">PS00134</a> ; Serine proteases, trypsin family, histidine active site (PATTERN)	
Consensus pattern:	[LIVM] - [ST] - A - [STAG] - <b>H</b> - C <i>H is the active site residue</i>

TRYPSIN_SER, <a href="#">PS00135</a> ; Serine proteases, trypsin family, serine active site (PATTERN)	
Consensus pattern:	[DNSTAGC] - [GSTAPIMVQH] - x(2) - G - [DE] - <b>S</b> - G - [GS] - [SAPHV] - [LIVMFYWH] - [LIVMFYSTANQH] <i>S is the active site residue</i>



# Examples of secondary (derived) biological databases

ExPASy - PROSITE - Microsoft Internet Explorer

File Modifica Visualizza Preferiti Strumenti ?

Indirizzo <http://www.expasy.ch/prosite/>

Google  prosite

ExPASy Home page Site Map Search ExPASy Contact us Swiss-Prot ENZYME

Search  for

**proSite** Database of protein domains, families and functional sites

Home ScanProsite ProRule Documents Downloads Links

PROSITE consists of [documentation entries](#) describing protein domains, families and functional sites as well as associated [patterns](#) and [profiles](#) to identify them [[More details](#) / [References](#) / [Disclaimer](#) / [Commercial users](#)].

PROSITE is complemented by [new](#) **ProRule**, a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [[More details](#)].

Release 20.9, of 03-Apr-2007 (1478 documentation entries, 1322 patterns, 720 profiles and 744 ProRule)

**PROSITE access**

☐ add wildcard <sup>\*\*\*</sup>

Browse:

- [by documentation entry](#)
- [by ProRule description](#) <sup>new</sup>
- [by taxonomic scope](#) <sup>new</sup>
- [by number of positive hits](#) <sup>new</sup>

SRS - Sequence Retrieval System

**PROSITE tools**

Scan a sequence against PROSITE - quick scan  
Enter your sequence or a [UniProtKB](#) ([Swiss-Prot](#) or [TrEMBL](#)) ID or AC [[help](#)]:

Output includes graphical view and feature detection.

FZ KRINGLE\_2 PROTEIN\_KINASE\_DOM

Searchable with a novel protein sequence in order to find matches

**proSite** is a database of protein families, domains and functional sites

# Examples of secondary (derived) biological databases



## OMIM®

An Online Catalog of Human Genes and Genetic Disorders

Updated October 17, 2022

Search OMIM for clinical features, phenotypes, genes, and more...



**Advanced Search :** [OMIM](#), [Clinical Synopses](#), [Gene Map](#)

**Need help?** : [Example Searches](#), [OMIM Search Help](#), [OMIM Video Tutorials](#)

**Mirror site :** <https://mirror.omim.org>

OMIM is supported by a grant from NHGRI, licensing fees, and [generous contributions from people like you](#).

[Make a donation!](#)



**OMIM** (Online Mendelian Inheritance in Man) is a catalog of human genes and genetic disorders

# Main integrated databases


A number of **centers** were funded to provide access to a large number of major databases in an **integrated environment**

Two main molecular biology database systems exist, which provide integrated access to nucleotide and protein sequence data, gene-centered and genomic mapping information, 3D structure data, PubMed MEDLINE, homology search software and more

**Entrez @ NCBI** (National Center for Biotechnology Information) – based in Bethesda, Maryland (USA)

**EMBL-EBI** (European Bioinformatics Institute) – based in Hinxton, Cambridge (UK)

# Main integrated databases

 **National Library of Medicine**  
National Center for Biotechnology Information

Log in

Search NCBI

hemoglobin

×

Search

Results found in 30 databases

<b>Literature</b>	<b>Genes</b>	<b>Proteins</b>
Bookshelf 11,351	Gene 8,859	Conserved Domains 35
MeSH 923	GEO DataSets 4,310	Identical Protein Groups 51,040
NLM Catalog 1,035	GEO Profiles 292,556	Protein 249,802
PubMed 239,625	HomoloGene 15	Protein Family Models 167
PubMed Central 430,147	PopSet 239	Structure 1,005
<b>Genomes</b>	<b>Clinical</b>	<b>PubChem</b>
Assembly 0	ClinicalTrials.gov 12,847	BioAssays 2,303
BioCollections 0	ClinVar 2,438	Compounds 11
BioProject 305	dbGaP 104	Pathways 31
BioSample 2,200	dbSNP 0	Substances 677
Genome 0	dbVar 8	
Nucleotide 195,607	GTR 485	
SRA 4,553	MedGen 1,247	
Taxonomy 0	OMIM 288	

# Main integrated databases

The EMBL-EBI website has been redesigned. Please [send us feedback](#) about this page.

EMBL's European Bioinformatics Institute

## EMBL-EBI

Unleashing the potential of big data in biology

Find a gene, protein or chemical

Example searches: [blast keratin bfl1](#) | [About EBI Search](#)

**Search**

- ✓ All
- Science search
  - Genomes & metagenomes
  - Nucleotide sequences
  - Protein sequences
  - Small molecules
  - Gene expression
  - Gene-Disease Associations
  - Diseases
  - Molecular interactions
  - Reactions & pathways
  - Protein families
  - Literature
  - Samples & ontologies
  - Search web content

**Find data resources** →

**Train with us** →

**Latest news** →

**Explore our research** →

## ***Focalizing a DB search....***

*Databases usually provide **fields** which can be selected in order to focalize a search, e.g. author ([au]), organism [OS], etc*

Imagine you wish to search a scientific article or a protein sequence by a scientist named **E(lisabetta) Coli**

You will have to specify that “**E Coli**” refers to an author (& not a species)!!

e.g. in PubMed: scientific literature database, let's use:

“E coli”

“Coli E[au]”

Example from A. Lesk, “Introduction to Bioinformatics” - Oxford University Press

Edit View Favorites Tools Help

back → × ↻ ↺ Search Favorites Media ↻ ↺ ↺ ↺ ↺ ↺

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?CMD=search&DB=pubmed

ch ×

new ↻ Next >>

Search The Web

Find a Web page containing:

Search

Sought to you MSN Search

Search for other sites:  
[or Folders](#)  
[computers](#)  
[people](#)

2006 Microsoft.  
[SN Privacy](#)

NCBI

PubMed

A service of the National Library of Medicine and the National Institutes of Health

My NCBI 2  
[\[Sign In\]](#) [\[Register\]](#)

All Databases PubMed Nucleotide Protein Genome Structure OMIM PMC Journals Books

Search PubMed for E Coli Go Clear Save Search

Limits Preview/Index History Clipboard Details

Display Summary Show 20 Sort by Send to

All: 200777 Review: 7313

Items 1 - 20 of 200777

Page 1 of 10039 Next

☐ 1: [Kim SM, Kuzuyama T, Chang YJ, Kim SU](#) Links

☐ Cloning and characterization of 2-C-methyl-D: -erythritol 2,4-cyclodiphosphate synthase (MECS) gene from Ginkgo biloba.  
Plant Cell Rep. 2006 Mar 10; [Epub ahead of print]  
PMID: 16528563 [PubMed - as supplied by publisher]

☐ 2: [Zelic B, Bolf N, Vasic-Racki D](#) Links

☐ Modeling of the pyruvate production with Escherichia coli: comparison of mechanistic and neural networks-based models.  
Bioprocess Biosyst Eng. 2006 Mar 10; [Epub ahead of print]  
PMID: 16528533 [PubMed - as supplied by publisher]

☐ 3: [Lahiri C, Mandal S, Ghosh W, Dam B, Roy P](#) Links

☐ A Novel Gene Cluster soxSRT Is Essential for the Chemolithotrophic Oxidation of Thiosulfate and Tetrathionate by Pseudaminobacter salicylatoxidans KCT001.  
Curr Microbiol. 2006 Mar 9; [Epub ahead of print]  
PMID: 16528465 [PubMed - as supplied by publisher]

☐ 4: [Cookson AL, Taylor SC, Attwood GT](#) Links

☐ The prevalence of Shiga toxin-producing Escherichia coli in cattle and sheep in the lower North Island, New Zealand

Internet



Edit View Favorites Tools Help

back → × Search Favorites Media

http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?CMD=search&DB=pubmed

ch ×

new Next >>

Search The Web

Find a Web page containing:

Search

Sought to you MSN Search

Search for other sites:

or Folders Computers People

2006 Microsoft. See Privacy

NCBI

PubMed

A service of the National Library of Medicine and the National Institutes of Health

My NCBI [2]  
[Sign In] [Register]

All Databases PubMed Nucleotide Protein Genome Structure OMIM PMC Journals Books

Search PubMed for Coli E[au] Go Clear Save Search

Limits Preview/Index History Clipboard Details

Display Summary Show 20 Sort by Send to

All: 9 Review: 1

Items 1 - 9 of 9

One page

1: [Ramacciotti CE, Coli E, Paoli R, Gabriellini G, Schulte F, Castrogiovanni S, Dell'Osso L, Garfinkel PE.](#) Related Articles, Links

The relationship between binge eating disorder and non-purging bulimia nervosa.  
Eat Weight Disord. 2005 Mar;10(1):8-12.  
PMID: 15943166 [PubMed - in process]

2: [Ramacciotti CE, Guidi L, Bondi E, Coli E, Dell'Osso L, Pistoia S, Pucci E.](#) Related Articles, Links

Differential dynamic responses of luteinizing hormone to gonadotropin releasing hormone in patients affected by bulimia nervosa-purging versus non-purging type.  
Eat Weight Disord. 1997 Sep;2(3):150-5.  
PMID: 14655839 [PubMed - indexed for MEDLINE]

3: [Lenzi A, Maltinti E, Poggi E, Fabrizio L, Coli E.](#) Related Articles, Links

Effects of rivastigmine on cognitive function and quality of life in patients with schizophrenia.  
Clin Neuropharmacol. 2003 Nov-Dec;26(6):317-21.  
PMID: 14646612 [PubMed - indexed for MEDLINE]

4: [Ramacciotti CE, Coli E, Biadi O, Dell'Osso L.](#) Related Articles, Links

Silent pericardial effusion in a sample of anorexic patients.  
Eat Weight Disord. 2003 Mar;8(1):68-71.  
PMID: 12762627 [PubMed - indexed for MEDLINE]

5: [Ramacciotti CE, Coli E, Paoli R, Marazziti D, Dell'Osso L.](#) Related Articles, Links

Entrez PubMed  
Overview  
Help | FAQ  
Tutorials  
New/Noteworthy  
E-Utilities

PubMed Services  
Journals Database  
MeSH Database  
Single Citation  
Matcher  
Batch Citation Matcher  
Clinical Queries  
Special Queries  
LinkOut  
My NCBI

Related Resources  
Order Documents  
NLM Mobile  
NLM Catalog  
NLM Gateway

Internet

## Quality of data in DBs

Data and relative annotations in DBs may contain **errors**, so the **accuracy** can be an issue

Some of these errors can be checked automatically, by a computer-based approach (*e.g. that DNA sequences only consist of A/G/T/C*)

Some other errors can only be checked by manual curation

## Quality of data in DBs

In addition, databases can be **redundant**, which means that several entries can have identical data

Ideally, all these entries with same data should be included in a single entry, which refers to all the independent experiments and summarizes all the separate DB entries

In this way, a **nonredundant database** is obtained, which allows to discover all the information by reading the single entry

# Quality of data in DBs

Databases need to be regularly updated:

- to include new entries and
- to update and correct existing ones

It is important to recognize if the current entry differs from a copy made earlier, this can be done by using version numbers for entries or reporting the most recent date on which changes were made. A **unique identifier** for the entry is needed!

## Entries in DBs usually have a maturation

Not annotated ➡ Preliminary ➡ Not reviewed ➡ Standard

*For instance, in early stages of a genome annotation, genes are identified by only applying computational methods; they are thus annotated as “hypothetical” genes (translated proteins as hypothetical proteins)*

*Then, when experimental evidence becomes available confirming or not a gene, the relative entry will be updated*

# Updated list of available and relevant biological DBs

OXFORD  
ACADEMIC

Journals

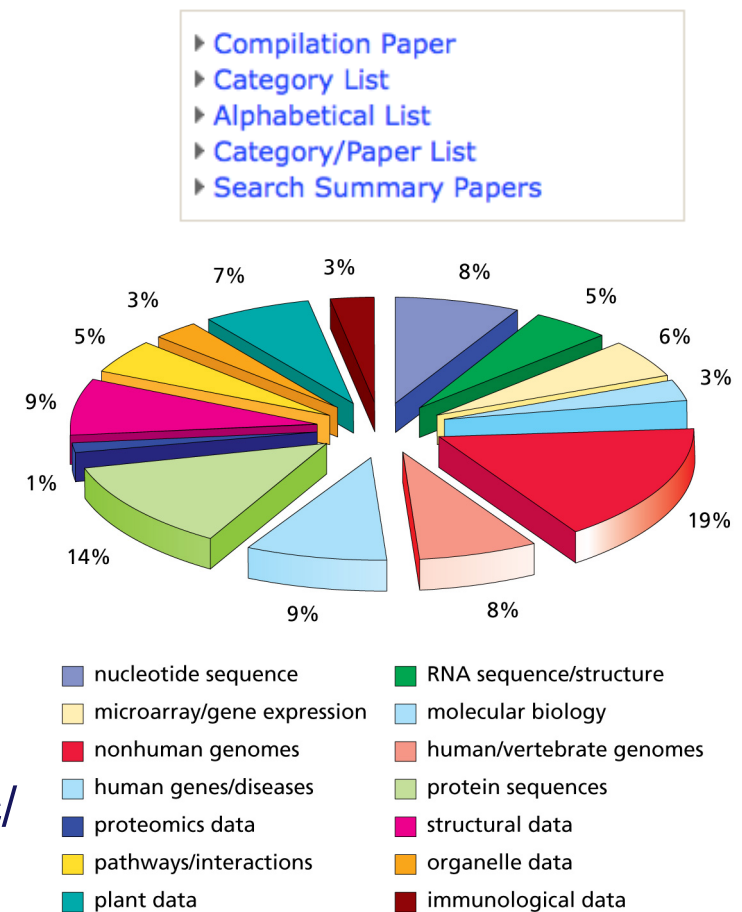
*Updated every January on  
Nucleic Acids Research  
(NAR) – Oxford Journals*

You are here: [NAR Journal Home](#) » Database Summary Paper Categories

## NAR Database Summary Paper Category List

Nucleotide Sequence Databases  
RNA sequence databases  
Protein sequence databases  
Structure Databases  
Genomics Databases (non-vertebrate)  
Metabolic and Signaling Pathways  
Human and other Vertebrate Genomes  
Human Genes and Diseases  
Microarray Data and other Gene Expression Databases  
Proteomics Resources  
Other Molecular Biology Databases  
Organelle databases  
Plant databases  
Immunological databases  
Cell biology

<https://www.oxfordjournals.org/nar/database/c/>



*statistics not updated*

## Lesson 7. Content

1. Biological databases. Constitute the backbone of bioinformatics research. Can be primary or secondary. Largely public and extremely useful although they may contain and propagate errors