

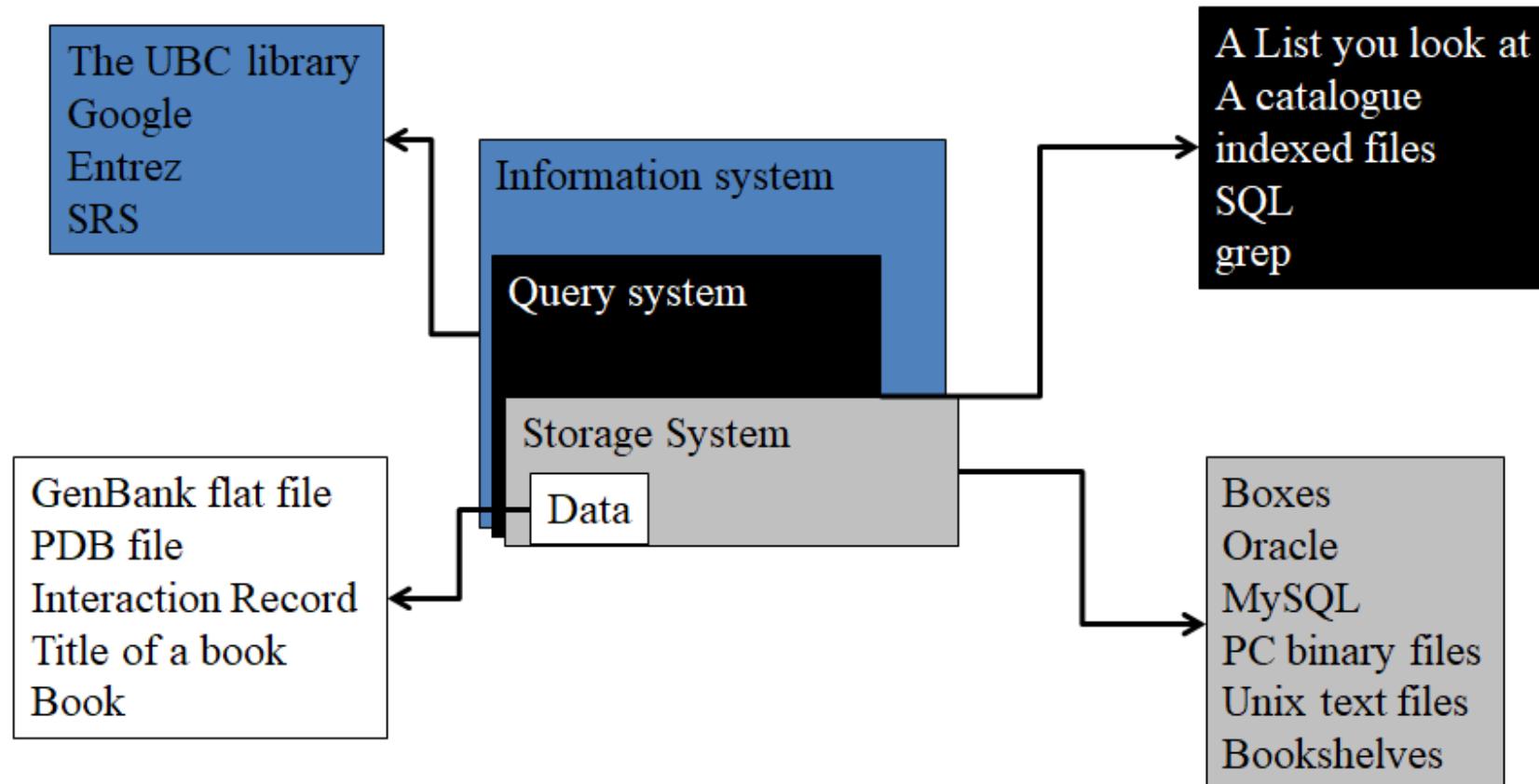
Biological Databases

Outline of discussion

- **Databases and its features**
- **Significance of databases in dealing with large amounts of data**
- **Types of databases**
- **Data integration**
- **Extraction of data from online biological databases**
- **Major online biological databases and their features**

What are databases?

- A database is an organized collection of structured information, or data, typically stored electronically in a computer system.
- It consists of basic units called records or entries.
- For example, a protein database would have protein entries as records and protein properties as fields (e.g., name of protein, length, amino-acid sequence).



Databank vs Database vs Data warehouse?

Types of Database Management System

A database is usually controlled by a database management system (DBMS). Together, the data and the DBMS, along with the applications that are associated with them, are referred to as a database system or just database.

Creating a database means you can remember the rules of the database rather than the locations of individual files and so find your way around more easily.

There are majorly three types of database management systems viz. (1) flat file indexing systems, (2) relational DBMS and (3) object-oriented DBMS.

1. Flat File Databases

A flat file database is not truly a database, it is simply an ordered collection of similar files.

Aarij M Hussain	210106001	B.Tech. (BSBE)	BSBE	Bioinformatics
Abhinav P Singh	210106002	B.Tech. (BSBE)	BSBE	Biophysics
-----	-----	-----	-----	-----
-----	-----	-----	-----	-----
-----	-----	-----	-----	-----
Yashraj Verma	210106085	B.Tech. (BSBE)	BSBE	Immunology

Flat File Databases

The emphasis in **formatting data** for a flat file database is **at the character level** i.e. at the level of how the data would appear if it were printed on a page.



Flat file databases are made useful by **ordering and indexing**. A collection of flat files on a computer file system can be ordered and stored in labeled folders exactly the same way as a collection of printed papers are ordered in a file cabinet Drawer.

The relationship of a flat file to a flat-file database

Flat file databases are often made **searchable by indexing**. An index pulls out a particular attribute from a file and pairs the attribute value in **the index with a filename and location**. It is analogous to a book index, which for example tells you where in a book you will find the word "genome."

Flat file databases in biology

The PDB began by using flat files in the well-known PDB format. The PDB now uses an object-oriented database backend to support database queries and file access.

Beyond the PDB, flat-file databases are still widely used by biologists. Many users of biological sequence data store and access sequences locally using the **Sequence Retrieval System (SRS)**, a flat file indexing system designed with biological data in mind.

Drawbacks of flat file database management systems

However, as flat file collections grow larger and larger, working with them becomes inefficient. An index is one-dimensional, so it is difficult (though not impossible) to make connections between attributes within an indexed flat file database.

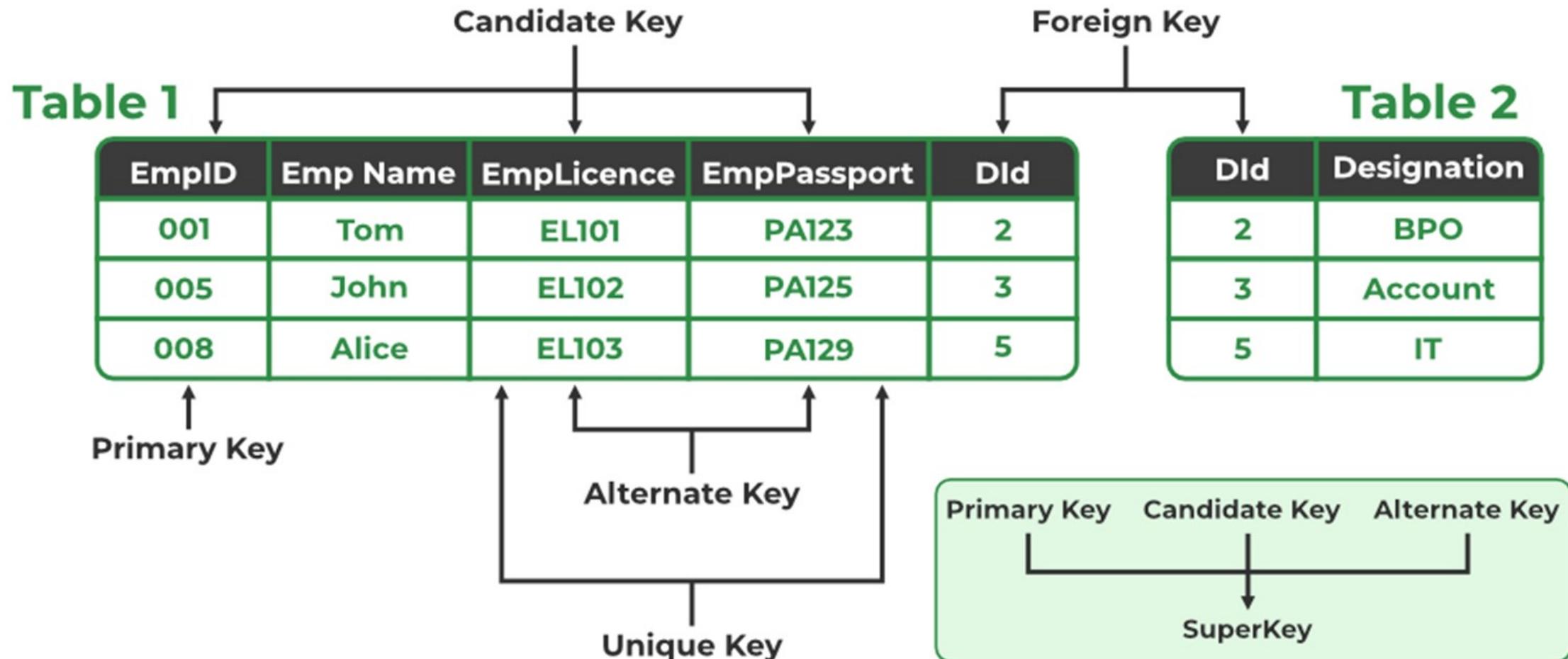
Types of Database Management System

2. Relational Databases: In a relational database, the information is stored in a collection of tables.

Roll No.	Name	Degree	Department	Course taken
210106001	Aarij Modhubhai Hussain	B.Tech. (BSBE)	BSBE	Bioinformatics
210106002	Abhinav Pratap Singh	B.Tech. (BSBE)	BSBE	Biophysics
-----	-----	-----	-----	-----
-----	-----	-----	-----	-----
210106085	Yashraj Verma	B.Tech. (BSBE)	BSBE	Immunology

The form of the tables follows rules that are uniform across the database, so you can access all the tables about student details.

Relational Databases



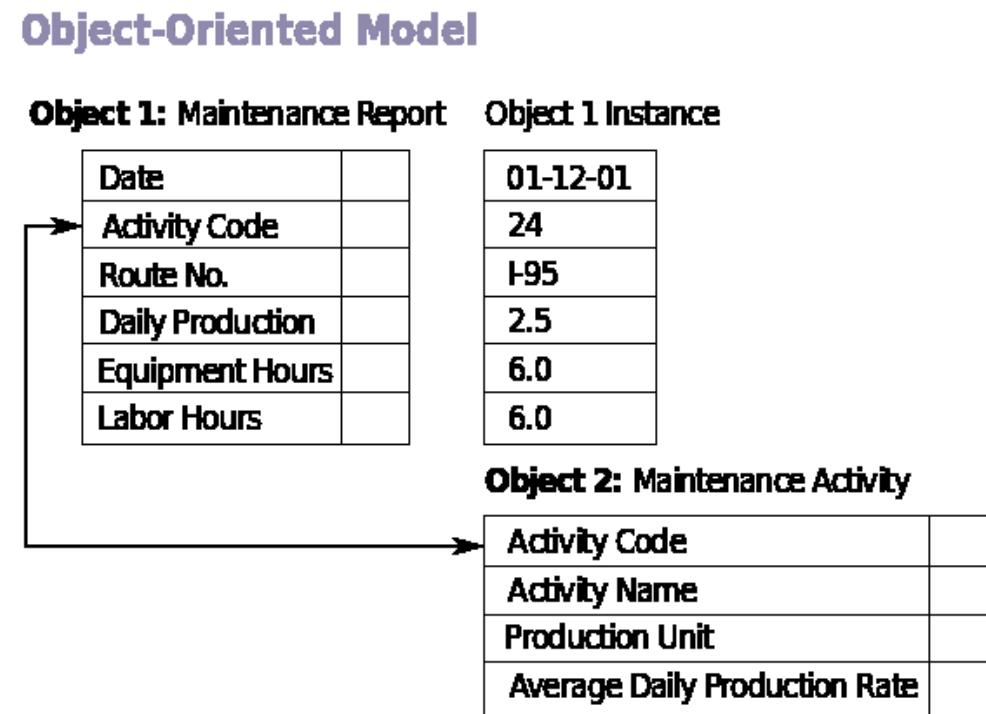
- **Advantages:** Manageability, Flexibility, Avoid errors.

- **Challenges:** Scalability (single server), Performance, Relationships.

Types of Database Management System

3. Object-Oriented DBMS

An OODBMS is a database management system in which information is represented in the form of objects as used in object-oriented programming.

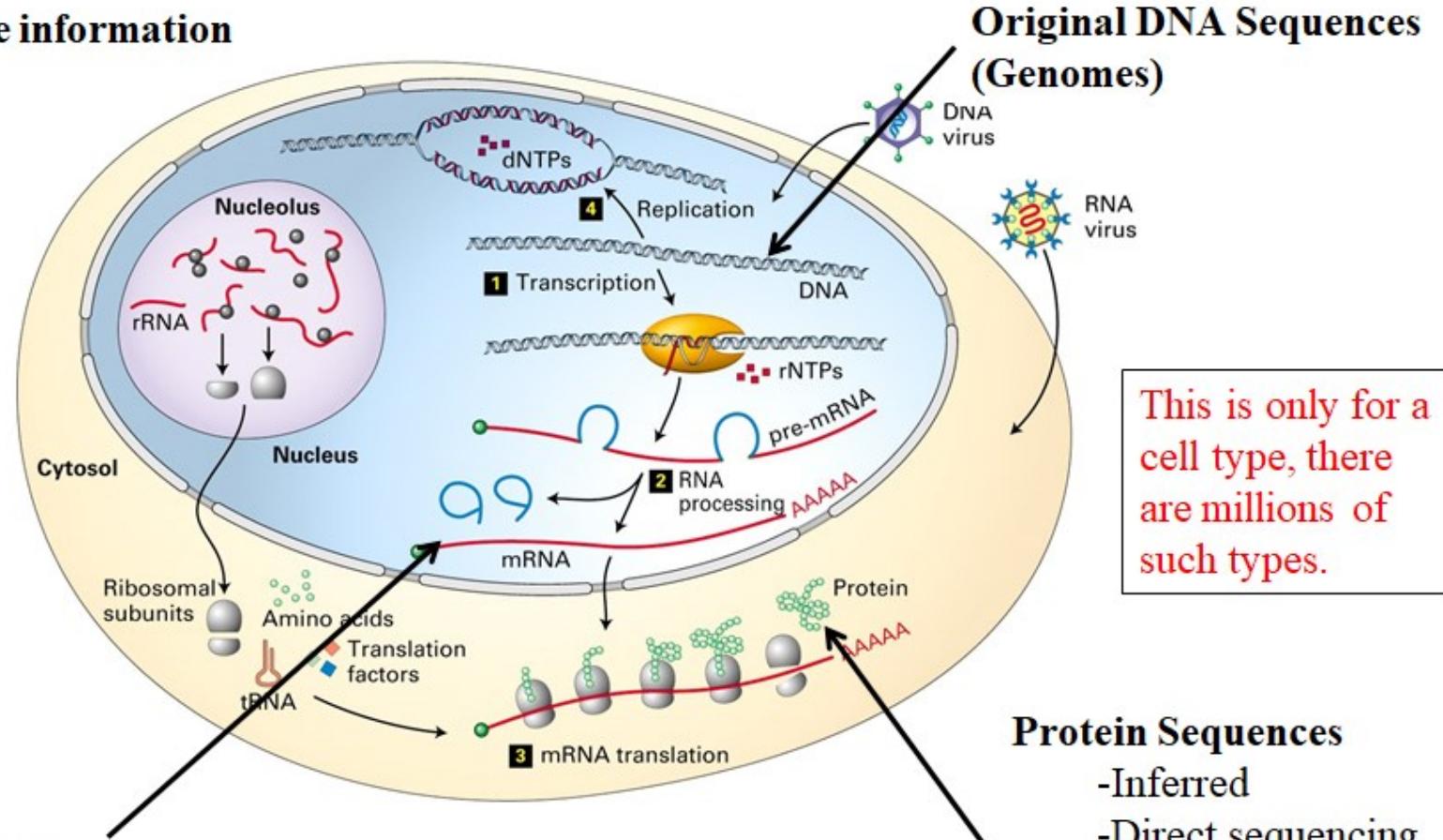


- **Advantages:** Faster access and performance.

- **Challenges:** Versatility, No standard query language, difficult to learn.

Biological databases: data

Literature information



Expressed DNA sequences
(= mRNA, = cDNA sequences)
Expressed Sequence Tags (ESTs)

- Protein Sequences**
- Inferred
 - Direct sequencing
- Protein structures**
- Experiments
 - Models

How many types of cells in humans?

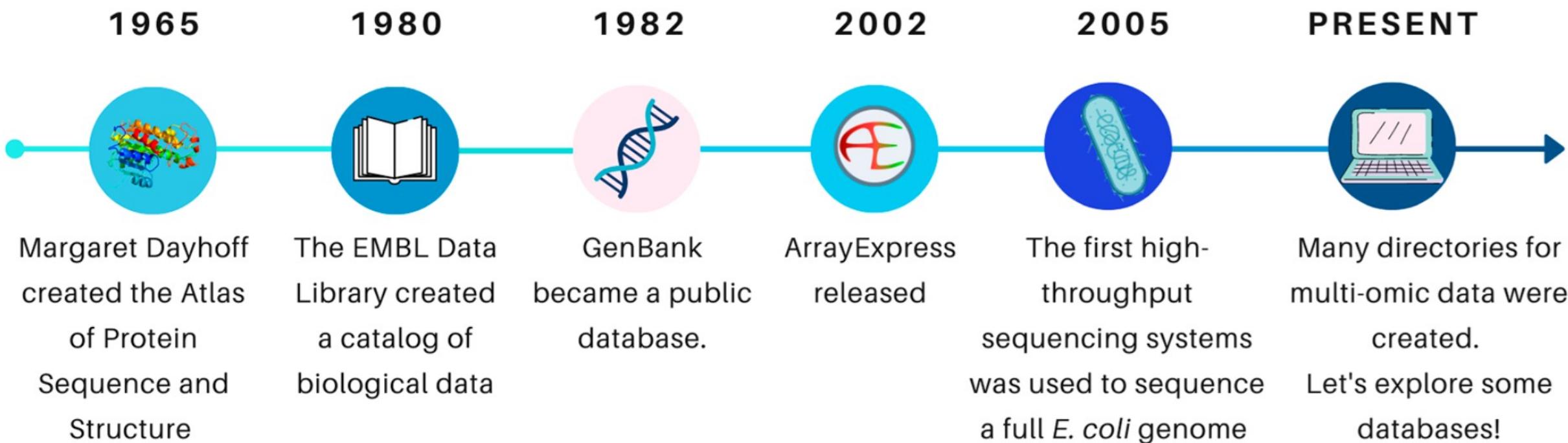
How many cells in humans?

How many genome sequences are available?

How many species are there on the Earth?

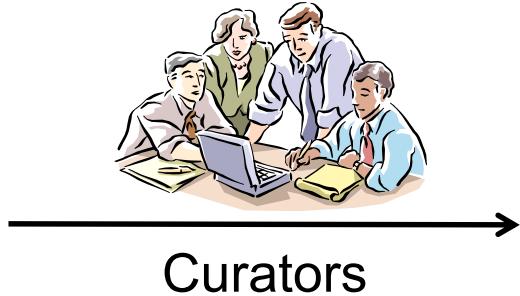
In fact, the scientific community has now generated data beyond the exabyte (10^{18}) level.

A brief history of biological databases



Types of biological databases

- **Primary (archival)**
 - Nucleic acid
 - GenBank
 - EMBL
 - DDBJ
 - Protein
 - Swiss-Prot
 - TrEMBL
 - PIR
 - Structure
 - PDB
 - MMDB
 - ModBase
 - CSD



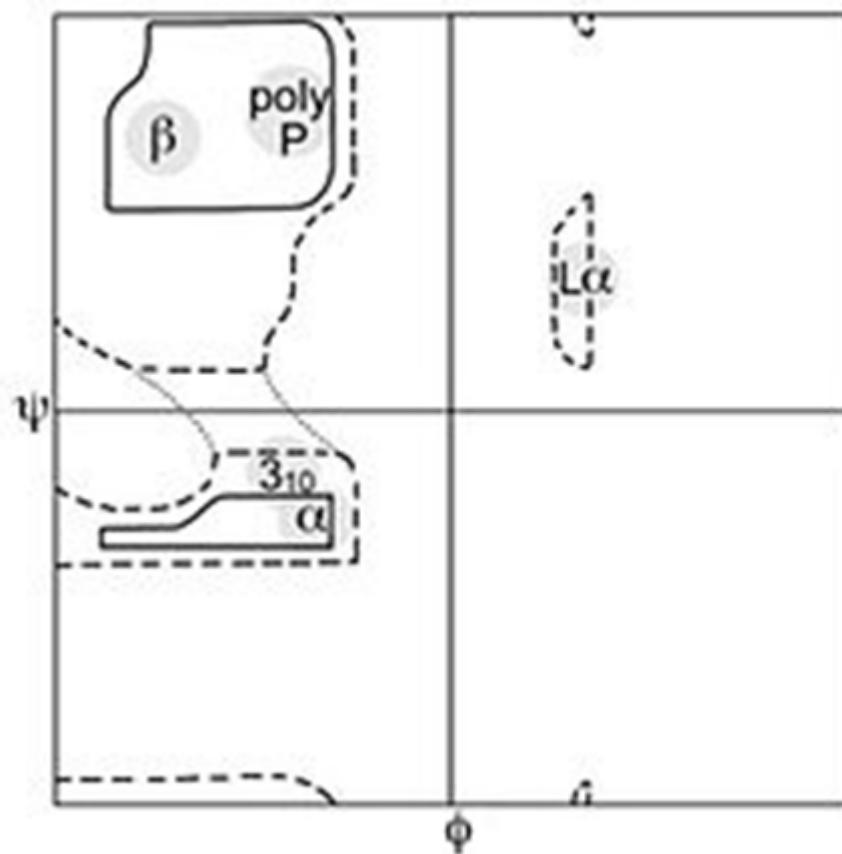
- **Secondary (curated)**
 - Nucleic acid
 - RefSeq
 - FlyBase
 - OMIM
 - Protein
 - Prosite
 - PRINTS
 - Pfam
 - Structure
 - LySDB
 - NDB

Composite databases: the initial data are taken from the primary database, and then they are merged together based on certain conditions. Example: UniProtKB (SwissProt + TrEMBL).

How do you define the ‘Perfect’ database?

- Comprehensive, but easy to search
- Annotated, but not “too annotated”
- A simple, easy to understand its structure
- Cross-referenced
- Minimum redundancy
- Easy retrieval of data
- Accuracy
- Up-to-date
- Good Interface
- Batch search/download options
- Simplify the information space by specialization

Bonus: Allows you to make discoveries



Reasons to search databases

- When obtaining a new DNA sequence, one needs to know whether it has already been deposited in the databanks, or whether they contain any homologous sequences (sequences which are derived from a common ancestry) exist there.
- Given a putative coding ORF, one can search for homologous proteins (similar in their folding or structure of function).
- To find similar non-coding DNA stretches in the database: repeat elements or regulatory sequences, for instance.
- There are other uses for specific purpose, like locating false priming sites for a set of PCR oligonucleotides.

DNA vs Protein searches

Protein sequence should be used for DB similarity search whenever possible because:

- There are very different DNA sequences that code for similar protein sequences.
- When comparing DNA sequences, one gets significantly more random matches than with proteins as (a) DNA is composed of 4 letters hence two unrelated DNA sequences are expected to have 25% similarity, (b) In contrast, proteins are made up of 20 letters, thus sensitivity of the comparison is improved, (c) DNA DBs are much larger and grow faster than protein DBs, implies bigger DB mean more random hits.
- For DNA, usually identity matrices is used while for proteins more sensitive matrices such as PAM or BLOSUM are used, which allows a better search.
- Proteins are rarely mutated during evolution, thus searching them reveals remote evolutionary relationships.

How to perform Database Searching

The amount of biological data is increasing so rapidly that it requires methods to search these information. There are three data retrieval systems in molecular biology (1) Sequence retrieval system (SRS), (2) Entrez and (3) DBGET.

These systems allow text searching of multiple databases and provide links to relevant information for entries that match the search criteria. The three systems differ in the databases they search and the links they have to the other information.

SRS

SRS is a homogeneous interface over 80 DBs developed at EBI, UK. It includes DBs of sequences, metabolic pathways, transcription factors, genomes, mappings, mutations and locus specific mutations.

Entrez

Entrez is a DB and retrieval system at NCBI. It is point to explore distinct but integrated DBs. It is easy to use but offers more limited information to search.

DBGET

DBGET is retrieval system developed at Univ. of Tokyo and provide information to 20 DBs, one at a time.

What can be discovered about a gene by a database search?

- A little or a lot, depending on the gene
 - **Evolutionary information:** homologous genes, taxonomic distributions, allele frequencies, synteny, etc.
 - **Genomic information:** chromosomal location, introns, UTRs, regulatory regions, shared domains, etc.
 - **Structural information:** associated protein structures, fold types, structural domains.
 - **Expression information:** expression specific to particular tissues, developmental stages, phenotypes, diseases, etc.
 - **Functional information:** enzymatic/molecular function, pathway/cellular role, localization, role in diseases.

Where to find biological databases?

The screenshot shows the homepage of the Nucleic Acids Research journal. At the top, there's a navigation bar with links for 'Issues', 'Section browse ▾', 'Advance articles', 'Submit ▾', 'Purchase', and 'About ▾'. To the right of the navigation is a search bar with a magnifying glass icon and a link to 'Advanced Search'. The main title 'Nucleic Acids Research' is displayed prominently with a stylized DNA helix logo to its left. Below the header, there's a large image of a city map with various colored dots representing research centers or databases.



Volume 51, Issue D1
6 January 2023

Article Contents

Abstract

NEW AND UPDATED DATABASES

NAR ONLINE MOLECULAR
BIOLOGY DATABASE
COLLECTION

ACKNOWLEDGEMENTS

FUNDING

JOURNAL ARTICLE

The 2023 Nucleic Acids Research Database Issue and the online molecular biology database collection

Daniel J Rigden , Xosé M Fernández

Nucleic Acids Research, Volume 51, Issue D1, 6 January 2023, Pages D1–D8,

<https://doi.org/10.1093/nar/gkac1186>

Published: 06 January 2023

PDF Split View Cite Permissions Share ▾

Abstract

The 2023 *Nucleic Acids Research* Database Issue contains 178 papers ranging across biology and related fields. There are 90 papers reporting on new databases and 82 updates from resources previously published in the Issue. Six more papers are updates from databases most recently published elsewhere. Major nucleic acid databases reporting updates include Genbank, ENA, ChIPBase, JASPAR, mirDIP and the Issue's first Breakthrough Article, NACDDB for Circular Dichroism data. Updates from BMRB and RCSB cover experimental



Advertisement

CITATIONS



More metrics information

VIEWS



ALTMETRIC

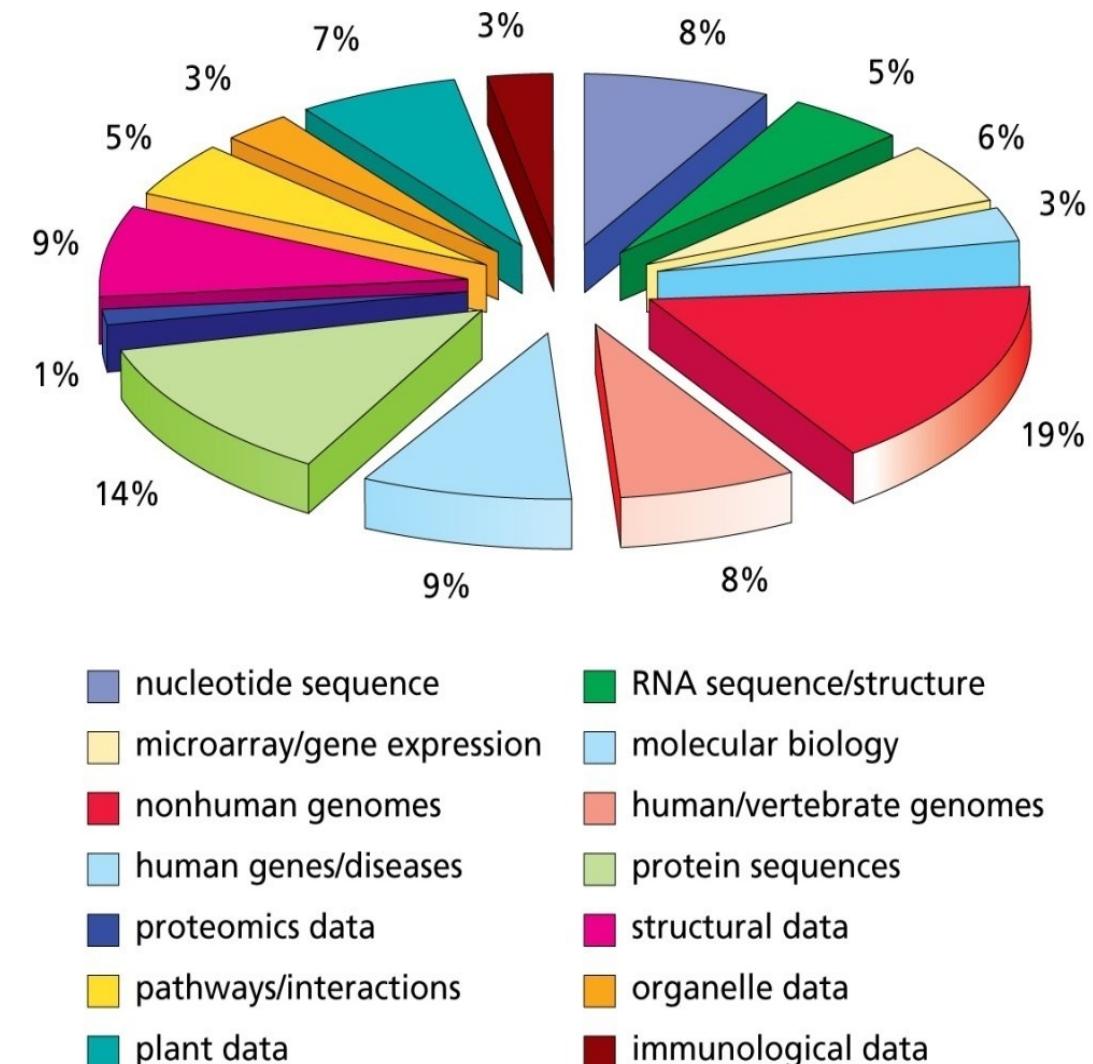


The total number of NAR online databases has been expanded to ~1764 in 2023 (1645, 2022).

Where to find biological databases?

These articles are divided into different categories:

1. Nucleotide Sequence Databases
2. RNA sequence databases
3. Protein sequence databases
4. Structure Databases
5. Genomics Databases (non-vertebrate)
6. Metabolic and Signaling Pathways
7. Human and other Vertebrate Genomes
8. Human Genes and Diseases
9. Microarray Data and other Gene Expression DBs
10. Proteomics Resources
11. Other Molecular Biology Databases
12. Organelle databases
13. Plant databases
14. Immunological databases
15. Cell biology

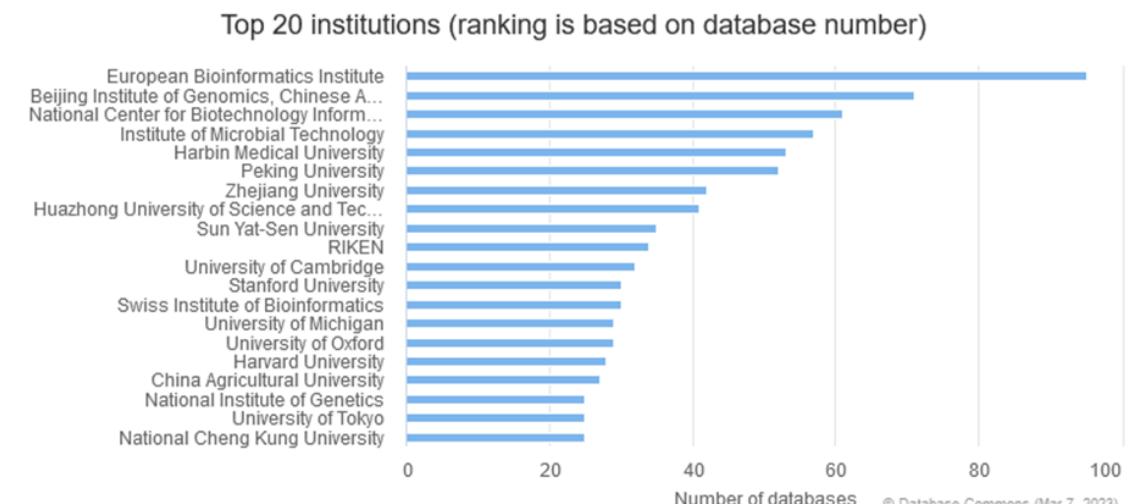
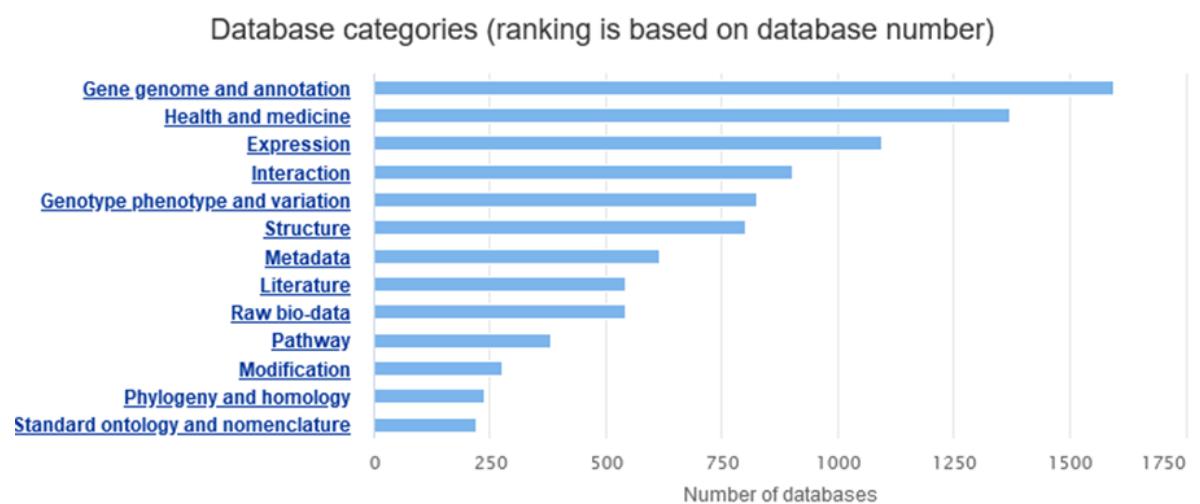
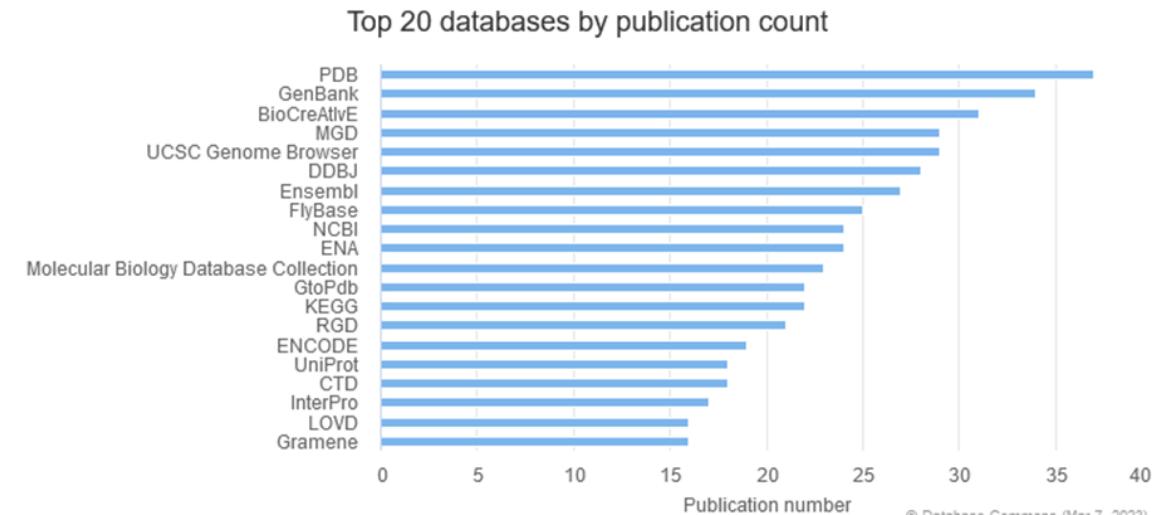
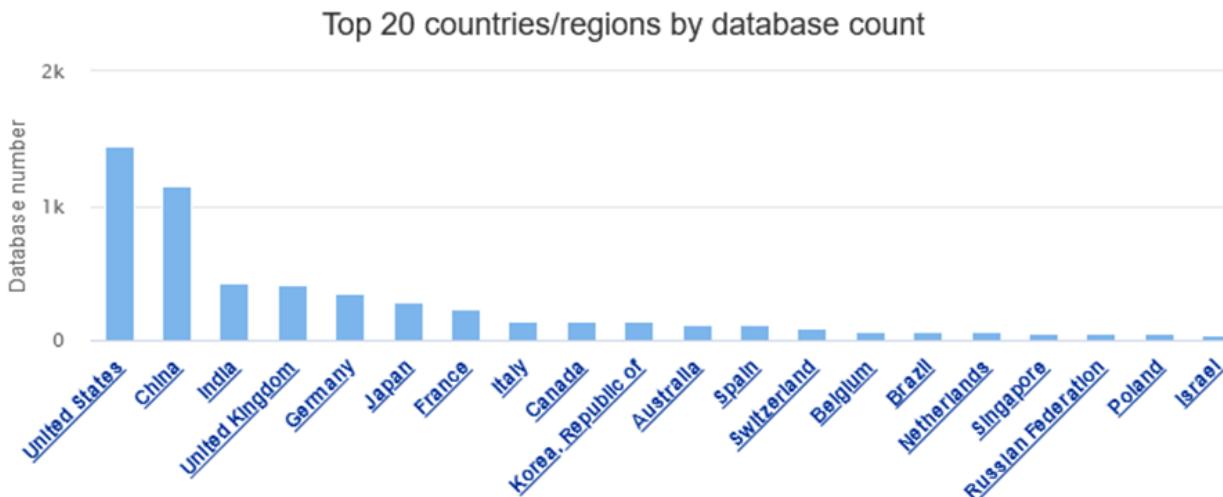


Some interesting online databases

Database name	URL	Brief description
Addgene Vector Database	http://www.addgene.org/vector-database/	Plasmid vectors from publications and commercial sources
AnimalTFDB	http://bioinfo.life.hust.edu.cn/AnimalTFDB/	Contains 125,135 TF genes and 80,060 cofactor genes from 97 animal genomes.
BacFITBase	http://www.tartaglialab.com/bacfitbase/	Bacterial genes relevant to host infection
BBCancer	http://bbcancer.renlab.org	Blood-based biomarkers for cancer
CancerPPD	http://crdd.osdd.net/raghava/cancerppd/	Experimentally validated anticancer peptides
EHFPI	http://biotech.bmi.ac.cn/ehfpi/	Essential Host Factors for Pathogenic Infection
EuRBPDB	http://EuRBPDB.syshospital.org	Eukaryotic RNA binding proteins
MoonProt	http://www.moonlightingproteins.org/	Moonlighting proteins
SuperDRUG2	http://cheminfo.charite.de/superdrug2/	A one stop resource for approved/ marketed drugs

Where to find biological databases?

Database Commons: a catalog of worldwide biological databases: No. of databases: 5910 (2023), No. of publications from: 9023 (2023), Categorized into: 13, No. of species: 1528, No. of countries: 72, No. of institutions: 1997. Highly cited databases: David, KEGG STRING, Pfam, ENCODE.



Where to find biological databases?

Pathguide: It contains information about **702** biological pathway related resources and molecular interaction related resources.

Bio.Tools: It strives to provide a comprehensive registry of software and databases, facilitating researchers from across the spectrum of biological and biomedical science to find, understand, utilize and cite the resources they need in their day-to-day work.

OMICtools: Leveraging life science data to reveal exciting new insights.

Metabases

BioGraph – A knowledge discovery service based on the integration of 21 heterogeneous databases.

Bioinformatic Harvester – Integrating 26 major protein/gene resources.

ConsensusPathDB – A molecular functional interaction database (contains 215541 unique interactions, 4601 pathways from overall 30 dbs).

Enzyme Portal – Integrates enzyme information such as small-molecule chemistry, biochemical pathways and drug compounds.

euGenes – (Genomic information of eukaryotic organisms).

GeneCards – (searchable integrated information of human genes).

mGen – Integrates GenBank, Refseq, EMBL and DDBJ.

PathogenPortal – A repository linking to the Bioinformatics Resource Centers sponsored by the National Institute of Allergy and Infectious Diseases (NIAID).

SOURCE – a unification tool which dynamically collects and compiles data from many scientific databases, and thereby attempts to encapsulate the genetics and molecular biology of genes from the genomes of *Homo sapiens*, *Mus musculus*, *Rattus norvegicus* into easy to navigate GeneReports.

Big Data Centre - Resources with significant updates in the past year include BioProject (a biological project library), BioSample (a biological sample library), Genome Sequence Archive (GSA, a data repository for archiving raw sequence reads), Genome Warehouse (GWH, a centralized resource housing genome-scale data), Genome Variation Map (GVM, a public repository of genome variations), Science Wikis (a catalog of biological knowledge wikis for community annotations) and IC4R (Information Commons for Rice).

Some specific biological databases

GenBank

The first DNA sequence database, established in 1979, was the Gene Sequence Database (GSDB), now known as GenBank.

GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>) is a comprehensive, public database that contains 19.6 trillion base pairs from over 2.9 billion nucleotide sequences for 5,04,000 formally described species.

NCBI, in cooperation with EMBL and other international organizations, provides the most complete collection of DNA sequence data in the world, as well as PubMed, a taxonomy database, and an alternate access point for protein sequence and structure data.

Daily data exchange with the European Nucleotide Archive (ENA) and the DNA Data Bank of Japan (DDBJ) ensures worldwide coverage.

Growth of GenBank

Genetic Sequence Data Bank (June 15, 2023)

NCBI-GenBank Flat File Release 254.0

Sequences: 2,43,560,863

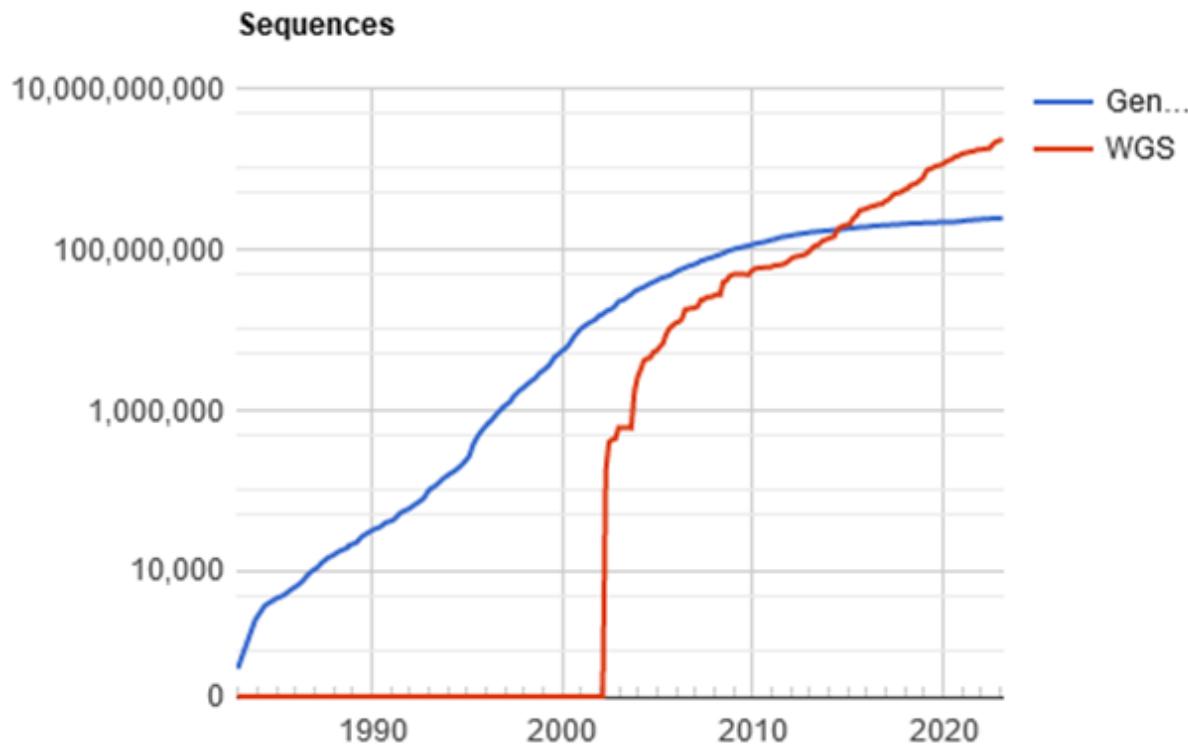


Table: Growth of GenBank Division

Division	Description	2022	2023	Increase	% Increase
BCT	Bacteria	130518385589	166217792419	35699406830.00	27.35
ENV	Environmental samples	7394414660	8516518905	1122104245.00	15.18
EST	Expressed sequence tags	43324455796	43330114068	5658272.00	0.01
GSS	Genome survey sequences	26380049011	26380049011	0.00	0.00
HTC	High-throughput cDNA	737423641	740853492	3429851.00	0.47
HTG	High-throughput genomic	27800219072	27801878633	1659561.00	0.01
INV	Invertebrates	108680334593	269338221858	160657887265.00	147.83
MAM	Other mammals	28568850588	41720029494	13151178906.00	46.03
PAT	Patent sequences	29588418021	30938105095	1349687074.00	4.56
PHG	Phages	935884237	1158493277	222609040.00	23.79
PLN	Plants	350590744188	484803006831	134212262643.00	38.28
PRI	Primates	15165437356	15619743253	454305897.00	3.00
ROD	Rodents	23336550435	66092410483	42755860048.00	183.21
STS	Sequence tagged sites	640923137	640923137	0.00	0.00
SYN	Synthetic	7994601379	8030787249	36185870.00	0.45
TLS	Targeted loci studies	39930167315	43852280645	3922113330.00	9.82
TSA	Transcriptome hot gundata	454757992932	511680950707	56922957775.00	12.52
UNA	Unannotated	4421782	4436341	14559.00	0.33
VRL	Viruses	39351597469	187366647663	148015050194.00	376.13
VRT	Other vertebrates	85320979451	99921122967	14600143516.00	17.11
WGS	Whole genome hot gundata	13888187863722	17511809676629	3623621812907.00	26.09

GenBank sequences are organized into 21 divisions, each of which is represented by a three-letter abbreviation.

Submission of data to GenBank

Direct submissions from scientists. NCBI provides timely and accurate processing and biological review of new entries and updates to existing entries.

BankIt: a WWW-based submission tool with wizards to guide the submission process.

Sequin: NCBI's stand-alone submission tool with wizards to guide the submission process.

Tbl2asn: a command-line program, automates the creation of sequence records for submission to GenBank using many of the same functions as Sequin. It is used primarily for submission of complete genomes and large batches of sequences.

Receiving an Accession Number for your Manuscript

GenBank will provide accession numbers for submitted sequences, usually within two working days.

GenBank submission types

Standard: GenBank accepts mRNA or genomic sequence data directly determined by the submitter.

EST, STS (Sequence-tagged site), GSS (Genome Survey Sequence), HTG (High-Throughput Genomic) Sequences, Complete Microbial Genomes, WGS (Whole Genome Shotgun) Sequences, TSA (Transcriptome Shotgun Assembly) Sequences, TPA (Third Party Annotation)

The following data is not accepted by GenBank

Non-contiguous sequences, Primer sequences, Protein sequences with no underlying nucleotide submission, Sequence containing a mix of genomic and mRNA sequence, Sequences without a physical counterpart (consensus sequences), Sequences with length less than 200 nucleotides.

A sample of GenBank record – abstract syntax notation

```
LOCUS      eIF4E                      2881 bp    DNA     linear  INV 27-OCT-2005
DEFINITION Drosophila melanogaster eukaryotic initiation factor 4E (eIF4E)
            gene, alternative splice products, complete cds.
ACCESSION
VERSION
KEYWORDS .
SOURCE     Drosophila melanogaster (fruit fly)
ORGANISM   Drosophila melanogaster
            Eukaryota; Metazoa; Arthropoda; Hexapoda; Insecta; Pterygota;
            Neoptera; Endopterygota; Diptera; Brachycera; Muscomorpha;
            Ephydrioidea; Drosophilidae; Drosophila.
REFERENCE  1 (bases 1 to 2881)
AUTHORS    Burnett,F.M., van der Waals,J.D. and Szent-Gyorgi,A.
TITLE      Environmental influences on the expansion of germline tandem
            repeats in several species of Galapagos finches
JOURNAL    Unpublished
REFERENCE  2 (bases 1 to 2881)
AUTHORS    Burnett,F.M., van der Waals,J.D. and Szent-Gyorgi,A.
TITLE      Direct Submission
JOURNAL    Submitted (27-OCT-2005) Evolutionary Biology Department, Oxbridge
            University, 1859 Tennis Court Lane, Camford OX1 2BH, United Kingdom
FEATURES   Location/Qualifiers
source      1..2881
            /organism="Drosophila melanogaster"
            /mol_type="genomic DNA"
            /strain="Oregon R"
gene        join(201..224,1550..1920,1986..2085,2317..2404,2466..2629)
            /gene="eIF4E"
CDS         join(201..224,1550..1920,1986..2085,2317..2404,2466..2629)
            /gene="eIF4E"
            /codon_start=1
            /product="eukaryotic initiation factor 4E-II"
            /translation="M V V L E T E K T S A P S T E Q G R P E P P T S A A A P A E A K D V K P K E D P Q E T G
            E P A G N T A T T A P A G D D A V R T E H L Y K H P L M N V U T L W Y L E N D R S K S W E D M Q N E I T S F D T V
            E D F W S L Y N H I K P P S E I K L G S D Y S L F K K N I R P M W E D A A N K Q G G R W V I T L N K S S K T D L D N
            L W L D V L L C L I G E A F D H S D Q I C G A V I N I R G K S N K I S I W T A D G N N E E A A L E I G H K L R D A L
            R L G R N N S L Q Y Q L H K D T M V K Q G S N V K S I Y T L"
```

Access to GenBank

- Search GenBank for sequence identifiers and annotations with **Entrez Nucleotide**, which is divided into three divisions: **CoreNucleotide** (the main collection), **dbEST** (Expressed Sequence Tags) and **dbGSS** (Genome Survey Sequences).
- Search and align GenBank sequences to a query sequence using BLAST.
- Search, link and download sequences programmatically using NCBI e-utilities.

SARS-CoV-2 resources

NCBI Datasets provides downloads for over 1.5 million complete SARS-CoV-2 genomes (<https://www.ncbi.nlm.nih.gov/datasets/coronavirus/genomes>) as well as a new taxonomy page for this virus (<https://www.ncbi.nlm.nih.gov/data-hub/taxonomy/2697049/>).



SARS-CoV-2 Data Hub in the NCBI Virus resource.

An interesting observation

An interesting observation

NCBI Resources How To

Genome **Genome** Search Limits Advanced Help

! COVID-19 is an emerging, rapidly evolving situation. X

[Public health information \(CDC\)](#) | [Research information \(NIH\)](#) | [SARS-CoV-2 data \(NCBI\)](#) | [Prevention and treatment information \(HHS\)](#)



Genome

This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.

Using Genome

[Help](#)
[Browse by Organism](#) UPDATED
[Download / FTP](#)
[Download FAQ](#)
[Submit a genome](#)

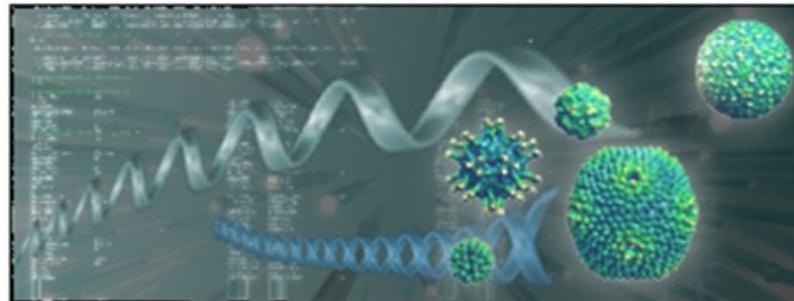
Custom resources

[Human Genome](#)
[Microbes](#)
[Organelles](#)
Viruses
[Prokaryotic reference genomes](#)

Other Resources

[Assembly](#)
[BioProject](#)
[BioSample](#)
[Genome Data Viewer](#)
[NCBI Datasets](#) NEW

An interesting observation



Viral Genomes

This resource provides viral genome sequence data and related information.

Explore Viral Genome Sequences

[Viral genome browser](#)

[Browse viral genomes by family](#)

[View all RefSeq and Neighbor nucleotide records](#)

Resource Tools

[Retrovirus Resource](#)

[Virus Variation Resource](#)

[Pairwise Sequence Comparison Tool \(PASC\)](#)

[Protein Clusters](#)

[Viral Genotyping Tool](#)

Virus Variation Resource

[Influenza virus](#)

[Dengue virus](#)

[West Nile virus](#)

[MERS coronavirus](#)

[Ebolavirus](#)

[Rotavirus](#)

[Zika virus](#)

An interesting observation

HOME | SEARCH | SITE MAP | Viral Genomes Home | Taxonomy groups | All viruses | Help | Contact us

Viruses - 11553 complete genomes

All viruses

Retrieve sequences: -- Select data set from the list --

Genome	Accession	RefSeq type	Source information	Segm	Length	Protein	Neighbors	Host	Created	Updated
Acidianus filamentous virus 1	NC_005830	complete		-	20869 nt	40	-	archaea	03/23/2004	12/20/2020
Acidianus filamentous virus 2	NC_009884	incomplete		-	31787 nt	52	-	archaea	04/01/2005	12/20/2020
Acidianus filamentous virus 3	NC_010155	incomplete		-	40449 nt	68	-		09/16/2007	12/20/2020
Acidianus filamentous virus 6	NC_010152	complete		-	39577 nt	66	-	archaea	09/16/2007	12/20/2020
Acidianus filamentous virus 7	NC_010153	complete		-	36895 nt	57	-	archaea	09/16/2007	12/20/2020
Acidianus filamentous virus 8	NC_010154	complete		-	38179 nt	61	-	archaea	09/16/2007	12/20/2020
Acidianus filamentous virus 9	NC_010537	complete		-	41172 nt	73	-	archaea	03/31/2008	12/20/2020
Acidianus rod-shaped virus 1	NC_009965	complete		-	24655 nt	41	-	archaea	05/10/2005	10/12/2021
Acidianus rod-shaped virus 2	NC_029314	complete	strain:ARV2	-	29763 nt	43	-	archaea	03/01/2016	12/21/2020
Pyrobaculum filamentous virus 1	NC_029548	complete	isolate:1	-	17714 nt	39	-	archaea	11/17/2016	12/21/2020
Stygiolobus rod-shaped virus	NC_025375	complete		-	28096 nt	37	-		10/03/2008	10/10/2021
Sulfolobales Beppu filamentous virus 2	NC_048128	incomplete		-	38091 nt	68	-	archaea	05/15/2020	12/21/2020
Sulfolobales Beppu filamentous virus 3	NC_048127	incomplete		-	31324 nt	54	-	archaea	05/15/2020	12/21/2020
Sulfolobales Mexican rudivirus 1	NC_019413	complete		-	27431 nt	37	-		11/13/2012	12/21/2020
Sulfolobus filamentous virus 1	NC_048037	complete	isolate:S48	-	37311 nt	66	-	archaea	05/15/2020	12/21/2020
Sulfolobus islandicus filamentous virus	NC_003214	complete		-	40900 nt	73	1	archaea	11/06/2001	12/20/2020
Sulfolobus islandicus rod-shaped virus 1	NC_004087	complete		-	32308 nt	45	1	archaea	03/05/2002	10/11/2021
Sulfolobus islandicus rod-shaped virus 10	NC_034625	incomplete		-	32735 nt	49	-	archaea	05/24/2017	10/11/2021
Sulfolobus islandicus rod-shaped virus 11	NC_034624	incomplete		-	33356 nt	50	-	archaea	05/24/2017	10/11/2021

An interesting observation

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy BioCollections

Search for as complete name lock Go Clear

Display 0 levels using filter: none

Human immunodeficiency virus 1

Taxonomy ID: 11676 (for references in articles please use NCBI:txid11676)

current name

Human immunodeficiency virus 1, ICTV accepted

genbank acronym: HIV-1
acronym: HIV
equivalent: human immunodeficiency virus 1 HIV-1

NCBI BLAST name: viruses

Rank: species

Genetic code: [Translation table 1 \(Standard\)](#)

Host: human|vertebrates

Other names:

heterotypic synonym

Human immunodeficiency virus type 1

heterotypic synonym

human immunodeficiency virus type 1 HIV-1

in-part

AIDS virus

[Lineage \(full\)](#)

Viruses; [Riboviria](#); [Paramaviridae](#); [Artverviricota](#); [Reoviraviricetes](#); [Ortervirales](#); [Retroviridae](#); [Orthoretrovirinae](#); [Lentivirus](#)

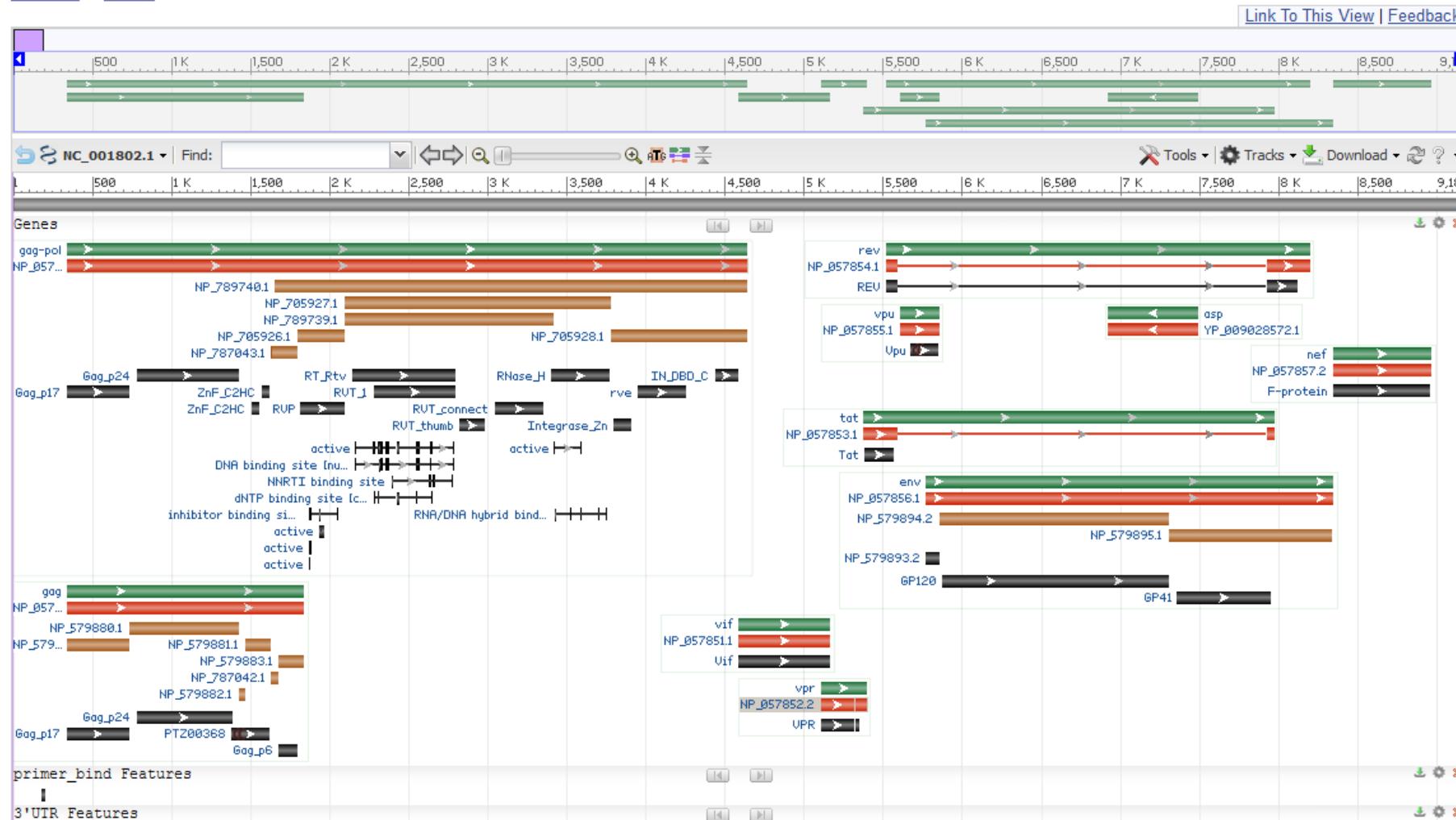
Entrez records		
Database name	Subtree links	Direct links
Nucleotide	1,004,161	998,869
Protein	1,381,716	1,374,779
Structure	2,378	1,853
Genome	1	1
Popset	5,906	5,906
GEO Datasets	768	768
PubMed Central	4,535	4,533
Gene	11	11
SRA Experiments	14,833	11,349
Identical Protein Groups	885,371	880,409
Bio Project	251	237
Bio Sample	13,675	11,484
Assembly	500	492
Probe	90	74
PubChem BioAssay	13,355	13,127
Taxonomy	2,793	1

An interesting observation

Human immunodeficiency virus 1, complete genome

NCBI Reference Sequence: NC_001802.1

[GenBank](#) [FASTA](#)



An interesting observation

Tax BLAST report

[Taxonomy report description](#)

RID HPD0HD1G016 (Expires on 05-18 01:17 am)
Query ID BAX30165.2
Description proton ATPase A [Polyandrocarpa misakiensis]
Molecule type amino acid
Query Length 617

Database Name SMARTBLAST/landmark
Description Landmark database for SmartBLAST
Program BLASTP 2.6.1+ > [Citation](#)

Lineage Report [Organism Report Taxonomy Report](#)

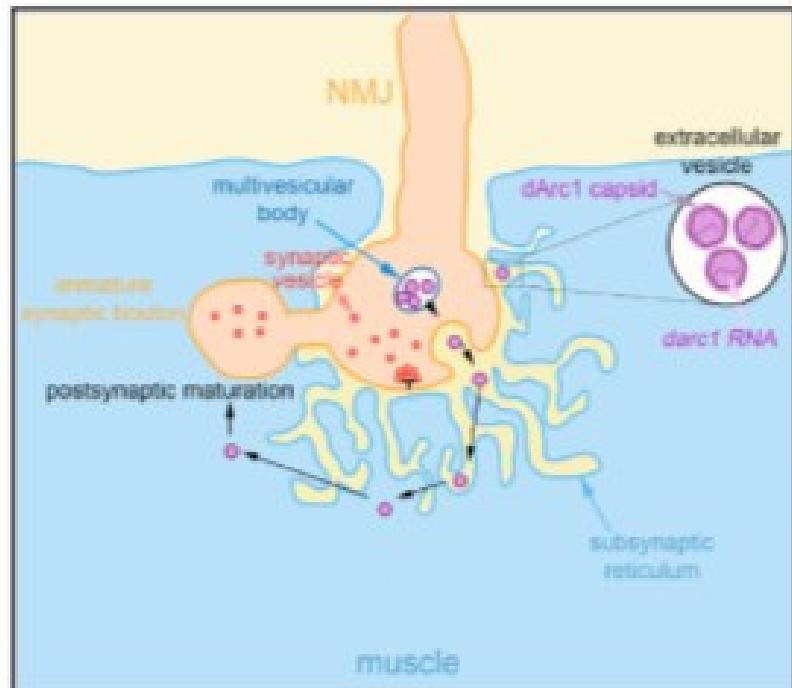
Organism	Blast Name	Score	Number of Hits	Description
cellular organisms			115	
• Eukaryota	eukaryotes		88	
• • Opisthokonta	eukaryotes		52	
• • • Bilateria	animals		46	
• • • • Euteleostomi	vertebrates		26	
• • • • • Danio rerio	bony fishes	1076	7	Danio rerio hits
• • • • • Mus musculus	rodents	1057	8	Mus musculus hits
• • • • • Homo sapiens	primates	1055	11	Homo sapiens hits
• • • • • Drosophila melanogaster	flies	1047	16	Drosophila melanogaster hits
• • • • • Caenorhabditis elegans	nematodes	1004	4	Caenorhabditis elegans hits
• • • Schizosaccharomyces pombe 972h-	ascomycetes	860	3	Schizosaccharomyces pombe 972h- hits
• • • Saccharomyces cerevisiae S288C	ascomycetes	448	3	Saccharomyces cerevisiae S288C hits
• • Glycine max	eudicots	885	13	Glycine max hits
• • Arabidopsis thaliana	eudicots	876	15	Arabidopsis thaliana hits
• • Dictyostelium discoideum AX4	cellular slime molds	876	4	Dictyostelium discoideum AX4 hits
• • Leishmania donovani	kinetoplastids	794	2	Leishmania donovani hits
• • Plasmodium falciparum 3D7	apicomplexans	767	2	Plasmodium falciparum 3D7 hits
• Methanothermobacter thermautotrophicus	euryarchaeotes	619	1	Methanothermobacter thermautotrophicus hits
• Deinococcus radiodurans R1	bacteria	588	2	Deinococcus radiodurans R1 hits
• Clostridioides difficile 630	firmicutes	577	4	Clostridioides difficile 630 hits
• Sulfolobus acidocaldarius	crenarchaeotes	572	2	Sulfolobus acidocaldarius hits
• Streptococcus pneumoniae R6	firmicutes	129	1	Streptococcus pneumoniae R6 hits
• Thermotoga maritima MSB8	thermotogales	125	2	Thermotoga maritima MSB8 hits
• Pseudomonas aeruginosa PAO1	g-proteobacteria	123	5	Pseudomonas aeruginosa PAO1 hits
• Streptomyces coelicolor A3(2)	high GC Gram+	124	1	Streptomyces coelicolor A3(2) hits
• Mycobacterium tuberculosis H37Rv	high GC Gram+	115	2	Mycobacterium tuberculosis H37Rv hits
• Shewanella oneidensis MR-1	g-proteobacteria	111	1	Shewanella oneidensis MR-1 hits
• Synechocystis sp. PCC 6803	cyanobacteria	108	1	Synechocystis sp. PCC 6803 hits
• Microcystis aeruginosa	cyanobacteria	104	1	Microcystis aeruginosa hits
• Escherichia coli str. K-12 substr. MG1655	enterobacteria	95.1	1	Escherichia coli str. K-12 substr. MG1655 hits
• Methanothermobacter	euryarchaeotes	89.7	1	Methanothermobacter hits
• Neisseria meningitidis MC58	b-proteobacteria	75.1	2	Neisseria meningitidis MC58 hits

Cells hack virus-like protein to communicate

Cell

Retrovirus-like Gag Protein Arc1 Binds RNA and Traffics across Synaptic Boutons

Graphical Abstract



Authors

James Ashley, Benjamin Cordy,
Diandra Lucia, Lee G. Fradkin,
Vivian Budnik, Travis Thomson

Correspondence

vivian.budnik@umassmed.edu (V.B.)
travis.thomson@umassmed.edu (T.T.)

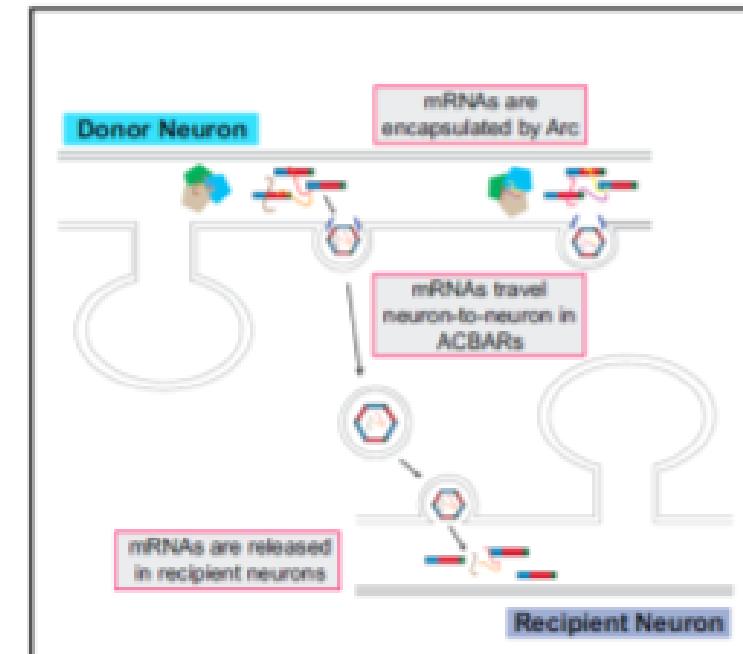
In Brief

The neuronal protein Arc is evolutionarily related to retrotransposon Gag proteins and forms virus-like capsid structures to transmit mRNA between cells in the nervous system.

Cell

The Neuronal Gene Arc Encodes a Repurposed Retrotransposon Gag Protein that Mediates Intercellular RNA Transfer

Graphical Abstract



Authors

Elissa D. Pastuzyn, Cameron E. C.
Rachel B. Kearns, ..., John A.G. E.
Cédric Feschotte, Jason D. Shep

Correspondence

jason.shepherd@neuro.utah.edu

In Brief

The neuronal protein Arc is evolutionarily related to retrotransposon Gag proteins and forms virus-like capsid structures that can transfer mRNA between neurons in the nervous system.

Endogenous retroviruses (ERVs) are remnants of ancient retroviral infections, and comprise nearly 8% of the human genome.

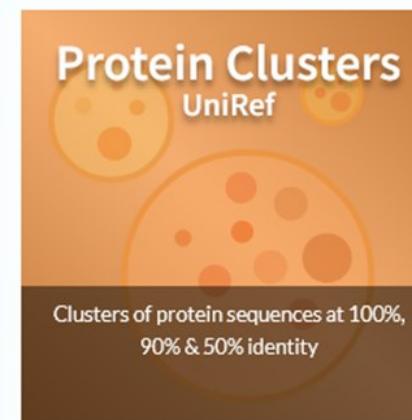
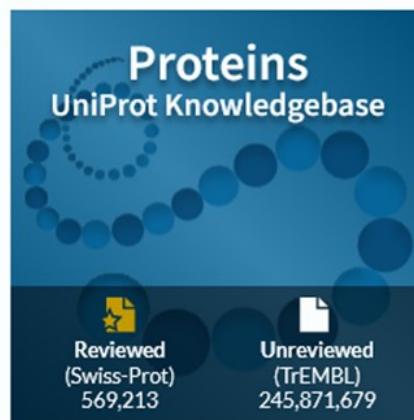
UniProtKB

Find your protein

[UniProtKB ▾](#)[Advanced](#) | [List](#) [Search](#)

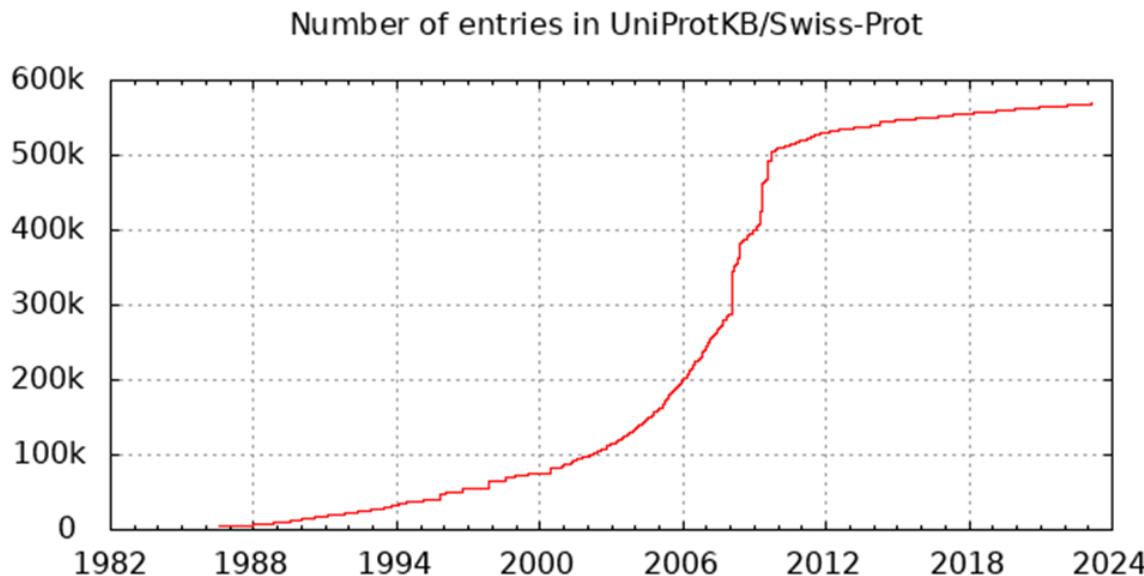
Examples: Insulin, APP, Human, P05067, organism_id:9606

UniProt is the world's leading high-quality, comprehensive and freely accessible resource of protein sequence and functional information. [Cite UniProt](#) 

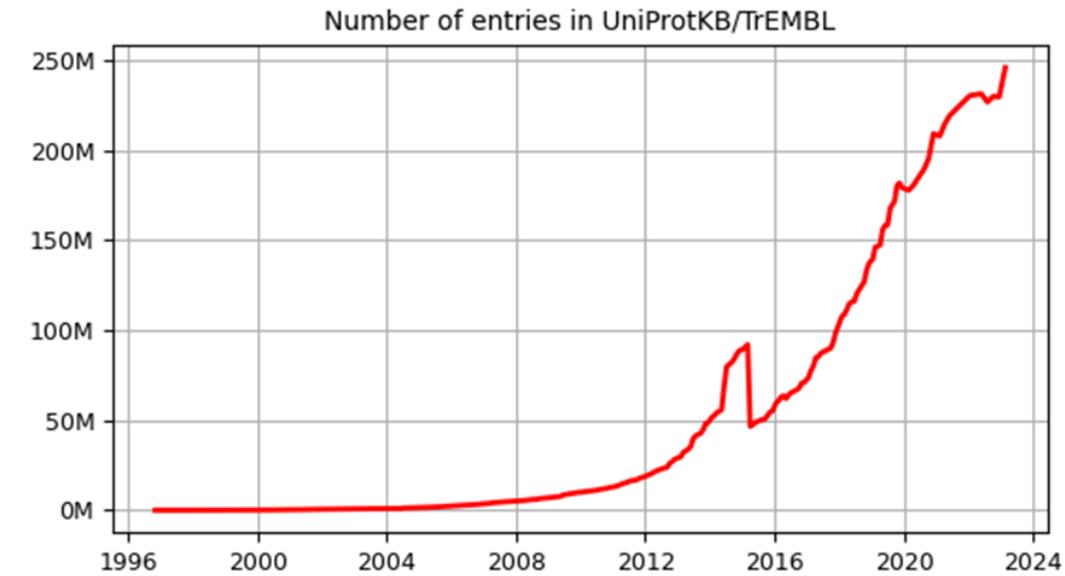


Growth of UniProtKB

Release **2023_03** of **28-Jun-2023** of **UniProtKB/Swiss-Prot** contains **5,69,793** sequence entries, curated from **2,93,323** unique references and comprising **206,004,162** amino acids.



Release **2023_03** of **28-Jun-2023** of **UniProtKB/TrEMBL** contains **248,272,897** sequence entries, comprising **86,599,718,967** amino acids.



UniProtKB: Search

UniProt BLAST Align Peptide search ID mapping SPARQL UniProtKB ▾ homo sapiens Advanced | List Search Help

Status

- Reviewed (Swiss-Prot) (26,806)
- Unreviewed (TrEMBL) (1,646,365)

Popular organisms

- Human (207,049)
- Mouse (8)
- A. thaliana (2)
- Fruit fly (2)
- Rat (2)

Taxonomy

Filter by taxonomy

Proteins with

- 3D structure (9,416)
- Active site (17,073)
- Activity regulation (24,334)

UniProtKB 1,673,171 results

BLAST Align Map IDs Download Add View: Cards Table Customize columns Share

Entry	Entry Name	Protein Names	Gene Names	Organism	Length	3D structures
Q8WY91	THAP4_HUMAN	Peroxynitrite isomerase THAP4[...]	THAP4, CGI-36, PP238	Homo sapiens (Human)	577 AA	X-ray: 1
Q53XC5	Q53XC5_HUMAN	Bone morphogenetic protein 4[...]	BMP4, hCG_20967	Homo sapiens (Human)	408 AA	
A8K571	A8K571_HUMAN	Bone morphogenetic protein 7 (Osteogenic protein 1), isoform CRA_b[...]	BMP7, hCG_40100	Homo sapiens (Human)	431 AA	
A8K660	A8K660_HUMAN	Adiponectin[...]	ADIPOQ, hCG_1784052	Homo sapiens (Human)	244 AA	
Q5TCX1	Q5TCX1_HUMAN	Triggering receptor expressed on myeloid cells 2[...]	TREM2	Homo sapiens (Human)	230 AA	
Q6FH53	Q6FH53_HUMAN	EDN1 protein[...]	EDN1, hCG_37405	Homo sapiens (Human)	212 AA	
Q75ME0	Q75ME0_HUMAN	STX1A protein[...]	STX1A, hCG_96107	Homo sapiens (Human)	288 AA	
Q9Y6N6	LAMC3_HUMAN	Laminin subunit gamma-3[...]	LAMC3	Homo sapiens	1,575 AA	

UniProtKB: customize results

UniProt BLAST Align Peptide search ID mapping SPARQL UniProtKB ▾ homo sapiens Advanced | List Search Help

Status
Reviewed (Swiss-Prot) (26,806) X

Popular organisms
Human (20,422)
Mouse (6)
Fruit fly (2)
A. thaliana (1)
Rat (1)

Taxonomy
Filter by taxonomy

Proteins with
3D structure (8,632)
Active site (3,003)
Activity regulation (1,996)
Allergen (6)
Alternative products

UniProtKB 26,806 results

BLAST Align Map IDs Download Add View: Cards Table Customize columns Share

Entry	Entry Name	Protein Names	Gene Names	Organism	Length	3D structures
Q8WY91	THAP4_HUMAN	Peroxynitrite isomerase THAP4[...]	THAP4, CGI-36, PP238	Homo sapiens (Human)	577 AA	X-ray: 1
Q9Y6N6	LAMC3_HUMAN	Laminin subunit gamma-3[...]	LAMC3	Homo sapiens (Human)	1,575 AA	
Q9H987	SYP2L_HUMAN	Synaptopodin 2-like protein	SYNPO2L	Homo sapiens (Human)	977 AA	
Q9UKP3	ITBP2_HUMAN	Integrin beta-1-binding protein 2[...]	ITGB1BP2, MSTP015	Homo sapiens (Human)	347 AA	
Q5T230	UTF1_HUMAN	Undifferentiated embryonic cell transcription factor 1	UTF1	Homo sapiens (Human)	341 AA	
A0A087X1C5	CP2D7_HUMAN	Putative cytochrome P450 2D7[...]	CYP2D7	Homo sapiens (Human)	515 AA	
A0A0B4J2F0	PIOS1_HUMAN	Protein PIGBOS1[...]	PIGBOS1	Homo sapiens (Human)	54 AA	
A0A0K2S4Q6	CD3CH_HUMAN	Protein CD300H[...]	CD300H	Homo sapiens	201 AA	

UniProtKB: customize results

Customize columns

Reviewed x Entry Name x Protein names x Gene Names x Organism x Length x 3D x

Search for available columns Search

UniProt Data

- Names & Taxonomy 7
- Sequences 4
- Function 1
- Miscellaneous 1
- Interaction
- Expression
- Gene Ontology (GO)
- Pathology & Biotech
- Subcellular location
- PTM / Processing
- Structure

Reset to default Cancel Save

Advanced | List Share

Gene Names	Organism	Length	3D structures
THAP4, CGI-36, PP238	Homo sapiens (Human)	577 AA	X-ray: 1
LAMC3	Homo sapiens (Human)	1,575 AA	
SYNPO2L	Homo sapiens (Human)	977 AA	
ITGB1BP2, MSTP015	Homo sapiens (Human)	347 AA	
UTF1	Homo sapiens (Human)	341 AA	
CYP2D7	Homo sapiens (Human)	515 AA	
PIGBOS1	Homo sapiens (Human)	54 AA	
CD300H	Homo sapiens	201 AA	

UniProtKB: advanced search

Advanced Search x

Searching in UniProtKB

Add Field All All homo sapiens Remove

Search for field All UniProtKB AC Entry Name [ID] Protein Name [DE] Gene Name [GN] Organism [OS] Taxonomy [OC] Virus host Protein Existence [PE] Function ▶ Subcellular location ▶ Pathology & Biotech ▶

Cancel Search

Advanced | List Search

Customize columns Share ▾

Gene Names ▾	Organism ▾	Length ▾	3D structures
THAP4, CGI-36, PP238	Homo sapiens (Human)	577 AA	X-ray: 1
LAMC3	Homo sapiens (Human)	1,575 AA	
SYNPO2L	Homo sapiens (Human)	977 AA	
ITGB1BP2, MSTP015	Homo sapiens (Human)	347 AA	
Description	UTF1	Homo sapiens (Human)	341 AA
CYP2D7	Homo sapiens (Human)	515 AA	
PIGBOS1	Homo sapiens (Human)	54 AA	
CD300H	Homo sapiens	201 AA	

UniProtKB: refine search

UniProt BLAST Align Peptide search ID mapping SPARQL UniProtKB ▾ homo sapiens Advanced | List Search Help

Status
Reviewed (Swiss-Prot) (7,816) X

Popular organisms
Human (7,816) X

Taxonomy
Filter by taxonomy

Proteins with
3D structure (7,816) X

Active site (1,348)

Activity regulation (1,073)

Allergen (5)

Alternative products
(isoforms) (4,735)

More items

Protein existence

Protein level (7,815)

UniProtKB 7,816 results

BLAST Align Map IDs Download Add View: Cards Table Customize columns Share

Entry	Entry Name	Protein Names	Gene Names	Organism	Length	3D structures
Q8WY91	THAP4_HUMAN	Peroxynitrite isomerase THAP4[...]	THAP4, CGI-36, PP238	Homo sapiens (Human)	577 AA	X-ray: 1
A0A5B9	TRBC2_HUMAN	T cell receptor beta constant 2	TRBC2, TCRBC2	Homo sapiens (Human)	178 AA	X-ray: 61 Other: 1
A0AVK6	E2F8_HUMAN	Transcription factor E2F8[...]	E2F8	Homo sapiens (Human)	867 AA	X-ray: 1
A0JLT2	MED19_HUMAN	Mediator of RNA polymerase II transcription subunit 19[...]	MED19, LCMR1	Homo sapiens (Human)	244 AA	EM: 8
A4D126	ISPD_HUMAN	D-ribitol-5-phosphate cytidylyltransferase[...]	CRPPA, ISPD	Homo sapiens (Human)	451 AA	X-ray: 1
A6ND01	JUNO_HUMAN	Sperm-egg fusion protein Juno[...]	IZUMO1R, FOLR4, JUNO	Homo sapiens (Human)	250 AA	X-ray: 7
A6NIH7	U119B_HUMAN	Protein unc-119 homolog B	UNC119B	Homo sapiens (Human)	251 AA	X-ray: 2
A6NJ78	MET15_HUMAN	12S rRNA N4-methylcytidine (m4C)	METTL15, METT5D1	Homo sapiens	407 AA	EM: 3

UniProtKB: exploring a protein

UniProt BLAST Align Peptide search ID mapping SPARQL UniProtKB Advanced | List Search Help

Function Q8WY91 · THAP4_HUMAN

Names & Taxonomy Proteinⁱ Peroxynitrite isomerase THAP4 Amino acids 577
Subcellular Location Geneⁱ THAP4 Protein existenceⁱ Evidence at protein level
Disease & Variants Statusⁱ UniProtKB reviewed (Swiss-Prot) Annotation scoreⁱ 5/5
PTM/Processing Organismⁱ Homo sapiens (Human)

Expression Entry Feature viewer Publications External links History

Interaction BLAST Align Download Add Add a publication Entry feedback

Structure

Family & Domains

Sequence & Isoform

Similar Proteins

Functionⁱ

Heme-binding protein able to scavenge peroxynitrite and to protect free L-tyrosine against peroxynitrite-mediated nitration, by acting as a peroxynitrite isomerase that converts peroxynitrite to nitrate. Therefore, this protein likely plays a role in peroxynitrite sensing and in the detoxification of reactive nitrogen and oxygen species (RNS and ROS, respectively). Is able to bind nitric oxide (NO) in vitro, but may act as a sensor of peroxynitrite levels in vivo, possibly modulating the transcriptional activity residing in the N-terminal region. 2 Publications

Catalytic activity

peroxynitrite = nitrate 1 Publication

This reaction proceeds in the forward direction. 1 Publication

Source: Rhea 63116 ▾

Hide Rhea reaction

UniProtKB: exploring a protein

UniProt BETA BLAST Align Peptide search ID mapping SPARQL UniProtKB ▾ Advanced | List Search 📁 🗂️ 🗑️ Help

Function
Names & Taxonomy
Subcellular Location
Disease & Drugs
PTM/Processing
Expression
Interaction
Structure
Family & Domains
Sequence & Isoform
Similar Proteins

A detailed ribbon diagram of a protein structure, showing its three-dimensional fold. The structure is composed of several beta-sheets represented by grey ribbons and alpha-helices represented by grey coils. The overall shape is roughly spherical or oval.

SOURCE	IDENTIFIER	METHOD	RESOLUTION	CHAIN	POSITIONS	LINKS
-- Select --		-- Select --				
PDB	3IA8	X-ray	1.79 Å	A/B	415-577	PDBe · RCSB-PDB · PDBj · PDBsum
AlphaFold	AF-Q8WY91-F1	Predicted			1-577	AlphaFold

Feedback Help

Tools on the UniProtKB web site

The screenshot shows the UniProtKB BLAST search interface. At the top, there is a navigation bar with links for BLAST, Align, Peptide search, ID mapping, SPARQL, UniProtKB (with a dropdown menu), Advanced, List, Search, and Help. On the right side of the page, there are two buttons: 'Feedback' (blue) and 'Help' (green).

BLAST

Find a protein sequence to run BLAST sequence similarity search by UniProt ID (e.g. P05067 or A4_HUMAN or UPI0000000001).

UniProt IDs

OR

Enter one or more sequences (20 max). You may also [load from a text file](#).

Protein or nucleotide sequence(s) in FASTA format.

Target database: UniProtKB reference proteomes + Swiss-Prot

Restrict by taxonomy: Enter taxon names or IDs to include

Name your BLAST job: "my job title"

Protein Data Bank

Protein Data Bank (PDB)

RCSB PDB Deposit Search Visualize Analyze Download Learn About Documentation Careers COVID-19 MyPDB Contact us

RCSB PDB PROTEIN DATA BANK 207,791 Structures from the PDB 1,068,577 Computed Structure Models (CSM) ▾ 3D Structures Enter search term(s), Entry ID(s), or sequence Include CSM Advanced Search | Browse Annotations Help

PDB-101 wwPDB EMDataResource NAKB wwPDB Foundation PDB-Dev [f](#) [t](#) [y](#) [o](#)

New: More Computed Structure Models (CSM) available [Learn more](#)

Welcome

Deposit

Search

Visualize

Analyze

Download

Learn

RCSB Protein Data Bank (RCSB PDB) enables breakthroughs in science and education by providing access and tools for exploration, visualization, and analysis of:

- Experimentally-determined 3D structures from the **Protein Data Bank (PDB)** archive
- Computed Structure Models (CSM) from AlphaFold DB and ModelArchive

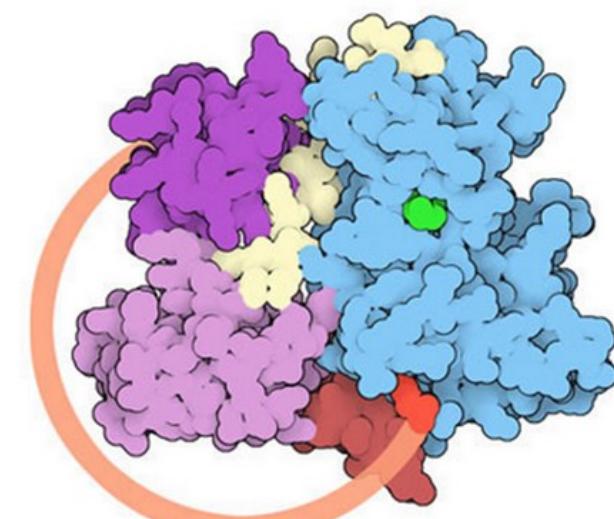
These data can be explored in context of external annotations providing a structural view of biology.

Explore NEW Features

PDB-101 Training Resources

July Molecule of the Month

c-Abl Protein Kinase and Imatinib



PDB Current Holdings Breakdown

Molecular type	X-ray	NMR	EM	Multiple methods	Neutron	Other	Total
Protein (only)	1,56,853	12,236	10,865	192	72	32	1,80,250
Protein + Oligosaccharide	9,180	34	1,906	7	1	0	11,128
Protein + Nucleic acid	8,200	283	3,420	7	0	0	11,910
Nucleic acid (only)	2,698	1,454	108	13	2	1	4,276
Other	164	32	9	0	0	0	205
Oligosaccharide (only)	11	6	0	1	0	4	22
Total	1,77,106	14,045	16,308	220	75	37	2,07,791

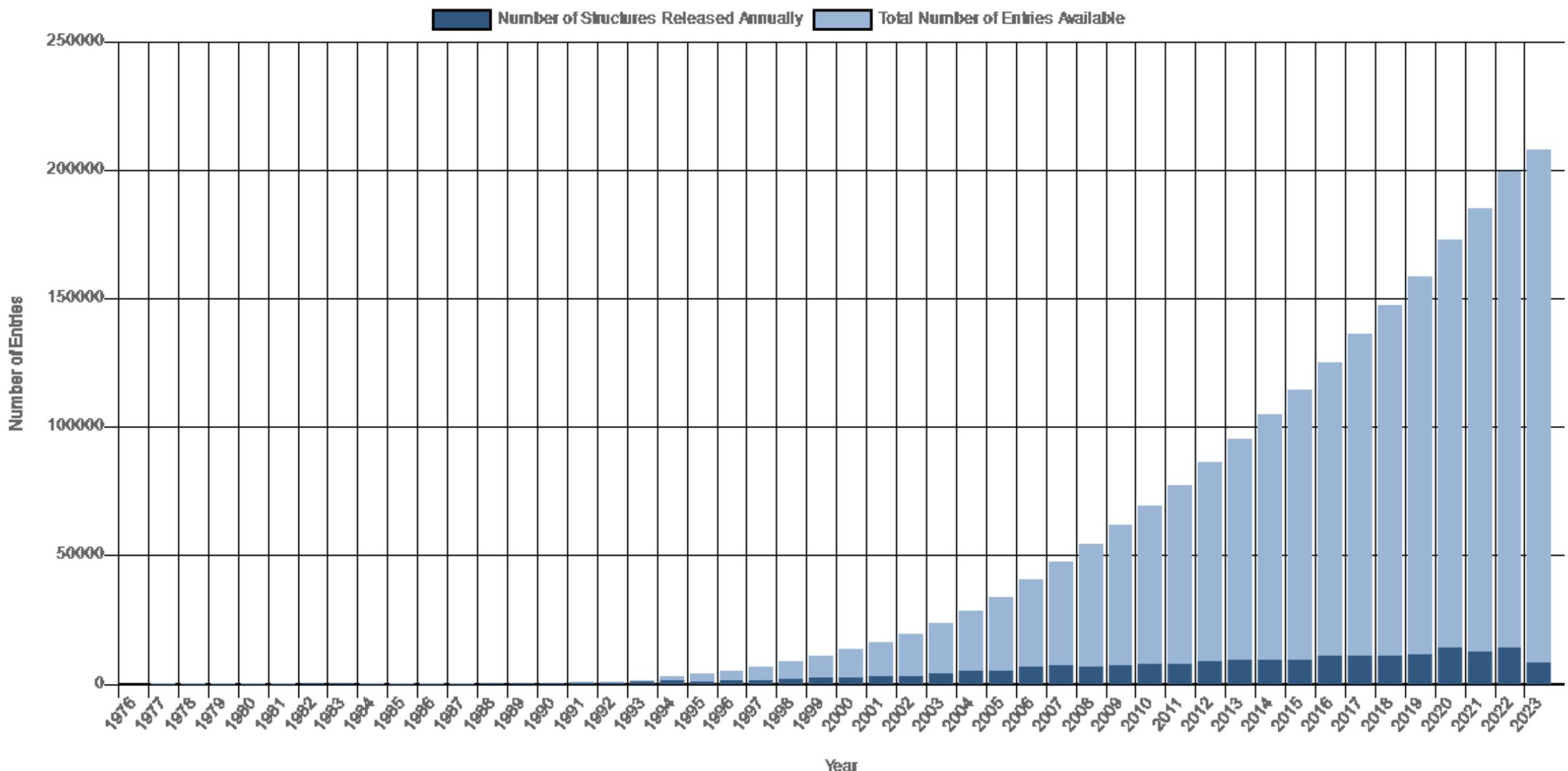
63,708 Structures of Human Sequences

16,385 Nucleic Acid Containing Structures

999,251 AlphaFoldDB

As of July 25, 2023

Growth of PDB



Search suggestions

RCSB PDB Deposit Search Visualize Analyze Download Learn About Documentation Careers MyPDB Contact us

PDB PROTEIN DATA BANK 201,979 Structures from the PDB 1,068,577 Computed Structure Models (CSM) ▾ 3D Structures sar-cov-2 Include CSM Help Advanced Search | Browse Annotations

PDB-101 wwPDB EMDDataResource NUCLEIC ACID DATABASE wwwPDB Foundation PDB-Dev Help

f t y o

Search Query History Browse Annotations MyPDB

QUERY: Full Text = "sars-cov-2" JSON MyPDB Login

Advanced Search Query Builder Help

Search Summary This query matches 3,322 Structures.

Refinements Structure Determination Methodology All Selected Download File View File

experimental (3,322)

Scientific Name of Source Organism 1 to 25 of 3,322 Structures Page 1 of 133 Sort by Score

Severe acute respiratory syndrome coronavirus 2 (2,964)
Homo sapiens (1,028)
Mus musculus (104)
synthetic construct (86)
Lama glama (66)
Severe acute respiratory syndrome-related coronavirus (46)
Vicugna pacos (37)

6XRZ Download File View File

The 28-kDa Frameshift Stimulation Element from the SARS-CoV-2 RNA Genome
Zhang, K., Zheludev, I., Hagey, R., Wu, M., Haslecker, R., Hou, Y., Kretsch, R., Pintilie, G., Rangan, R., Kladwang, W., Li, S., Pham, E., Souibgui, C., Baric, R., Sheahan, T., Souza, V., Glenn, J., Chiu, W., Das, R. (2020) Biorxiv

Released 2020-08-19
Method ELECTRON MICROSCOPY 6.9 Å
Organisms Severe acute respiratory syndrome coronavirus 2
Macromolecule Frameshift Stimulation Element from the SARS-CoV-2 RNA Genome (nucleic acid)

3D View

Query refinements

Refinements

SCIENTIFIC NAME OF SOURCE ORGANISM

- Homo sapiens (56066)
- Mus musculus (7860)
- Escherichia coli (6605)
- synthetic construct (5991)
- Escherichia coli K-12 (3929)
- Rattus norvegicus (3524)
- Bos taurus (3345)
- Saccharomyces cerevisiae (2987)
- Gallus gallus (2179)
- Saccharomyces cerevisiae S288C (2151)
- [More...](#)

TAXONOMY

- Eukaryota (102323)
- Bacteria (64823)
- Riboviria (10652)
- other sequences (6052)
- Archaea (5561)
- Duplodnaviria (2708)
- Varidnaviria (617)
- Monodnaviria (509)
- unclassified sequences (375)
- Naldaviricetes (48)
- [More...](#)

EXPERIMENTAL METHOD

- X-RAY DIFFRACTION (162628)
- SOLUTION NMR (12875)
- ELECTRON MICROSCOPY (10072)
- NEUTRON DIFFRACTION (194)
- ELECTRON CRYSTALLOGRAPHY (174)
- SOLID-STATE NMR (145)
- SOLUTION SCATTERING (73)
- FIBER DIFFRACTION (34)
- POWDER DIFFRACTION (20)
- EPR (8)
- [More...](#)

POLYMER ENTITY TYPE

- Protein (183882)
- DNA (7076)
- RNA (5894)
- NA-hybrid (208)
- Other (5)

REFINEMENT RESOLUTION (Å)

- < 0.5 (2)
- 0.5 - 1.0 (846)
- 1.0 - 1.5 (15654)
- 1.5 - 2.0 (55453)
- 2.0 - 2.5 (50596)
- 2.5 - 3.0 (29180)
- 3.0 - 3.5 (12218)
- 3.5 - 4.0 (4587)
- 4.0 - 4.5 (1755)
- > 4.5 (2816)

RELEASE DATE

- 1975 - 1979 (53)
- 1980 - 1984 (111)
- 1985 - 1989 (176)
- 1990 - 1994 (2376)
- 1995 - 1999 (7779)
- 2000 - 2004 (17422)
- 2005 - 2009 (32825)
- 2010 - 2014 (43063)
- 2015 - 2019 (53563)
- 2020 - 2024 (28665)

ENZYME CLASSIFICATION NAME

- Hydrolases (40914)
- Transferases (35252)
- Oxidoreductases (16895)
- Lyases (8147)
- Isomerase (4053)
- Ligases (3403)
- Translocases (1359)

MEMBRANE PROTEIN ANNOTATION

- PDBTM (6525)
- MemProtMD (5825)
- OPM (5430)
- mpstruc (5113)

SYMMETRY TYPE

- Asymmetric (111333)
- Cyclic (61288)
- Dihedral (14611)
- Icosahedral (1205)
- Helical (616)
- Octahedral (598)
- Tetrahedral (530)

SCOP CLASSIFICATION

- Alpha and beta proteins (a/b) (33561)
- All beta proteins (29560)
- Alpha and beta proteins (a+b) (30499)
- All alpha proteins (19133)
- Artifacts (19319)
- Multi-domain proteins (alpha and beta) (3462)
- Membrane and cell surface proteins and peptides (2051)
- Small proteins (3765)
- Coiled coil proteins (981)
- Low resolution protein structures (202)
- [More...](#)

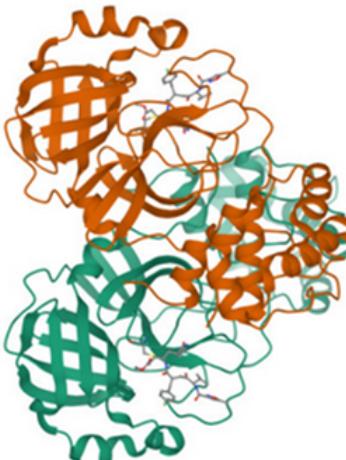


Molecular descriptions

RCSB PDB Deposit Search Visualize Analyze Download Learn About Documentation Careers MyPDB Contact us

Structure Summary 3D View Annotations Experiment Sequence Genome Ligands Versions

Biological Assembly 1 ?



3D View: Structure | 1D-3D View | Electron Density | Validation Report | Ligand Interaction

Global Symmetry: Cyclic - C2 (3D View)

Global Stoichiometry: Homo 2-mer - A2

Find Similar Assemblies

Biological assembly 1 assigned by authors and RCSB PDB

7P35

Structure of the SARS-CoV-2 3CL protease in complex with rupintrivir

PDB DOI: [10.22110/pdb7P35/pdb](https://doi.org/10.22110/pdb7P35/pdb)

Classification: VIRAL PROTEIN

Organism(s): Severe acute respiratory syndrome coronavirus 2

Expression System: Escherichia coli

Mutation(s): No

Deposited: 2021-07-07 Released: 2021-07-21

Deposition Author(s): Fabrega-Ferrer, M., Perez-Saavedra, J., Herrera-Morande, A., Coll, M.

Funding Organization(s): Spanish Ministry of Science, Innovation, and Universities, Spanish National Research Council

Experimental Data Snapshot

Method: X-RAY DIFFRACTION

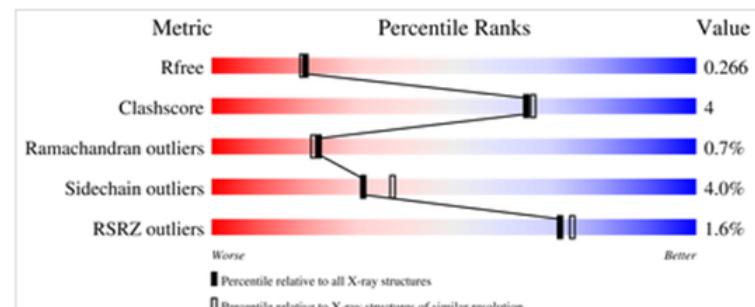
Resolution: 2.26 Å

R-Value Free: 0.258

R-Value Work: 0.200

wwPDB Validation

3D Report Full Report



Ligand Structure Quality Assessment

Molecular descriptions

RCSB PDB Deposit ▾ Search ▾ Visualize ▾ Analyze ▾ Download ▾ Learn ▾ About ▾ Documentation ▾ Careers MyPDB ▾ Contact us

Macromolecules

Find similar proteins by: [Sequence](#) ▾ (by identity cutoff) | [3D Structure](#)

Entity ID: 1

Molecule	Chains <small>i</small>	Sequence Length	Organism	Details	Image
3C-like proteinase	A [auth AAA], B [auth BBB]	306	Severe acute respiratory syndrome coronavirus 2	Mutation(s): 0 <small>i</small> Gene Names: ORF1ab EC: 2.7.7.48 (PDB Primary Data), 3.4.19.12 (PDB Primary Data), 3.4.22.69 (PDB Primary Data), 3.6.4.12 (PDB Primary Data), 3.6.4.1 (PDB Primary Data)	

Entity Groups i

Sequence Clusters

[30% Identity](#) [50% Identity](#) [70% Identity](#) [90% Identity](#) [95% Identity](#) [100% Identity](#)

Protein Feature View

[Expand](#)

Reference Sequence

7P35_1

7P35_1
UNMODELED A[auth AAA]
UNMODELED B[auth BBB]
HYDROPATHY



Sequence information

7P35

Usage

Display Files ▾

Download Files ▾

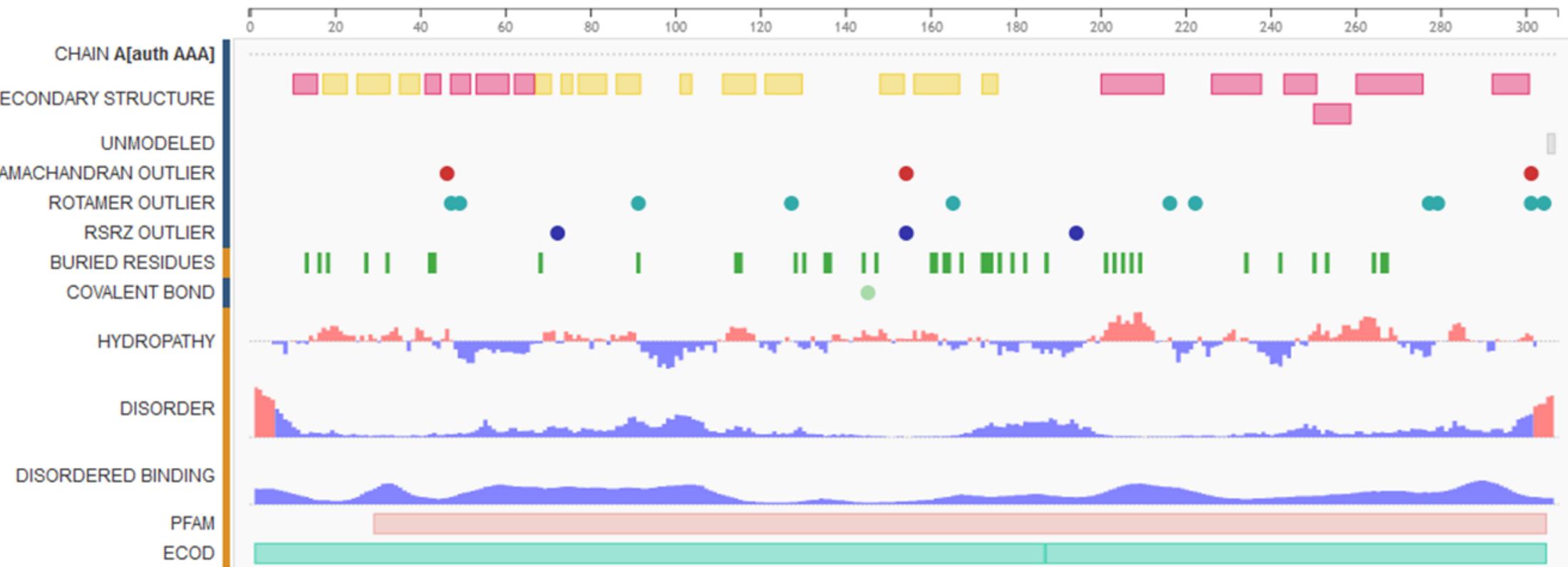
Structure of the SARS-CoV-2 3CL protease in complex with rupintrivir

Help

CHAIN

A [auth...

3c-like proteinase - Severe acute respiratory syndrome coronavirus 2 [View Features in 3D](#)



A sample of PDB file format

HEADER VIRAL PROTEIN/PROTEIN BINDING 01-JUN-20 7C8J
 TITLE STRUCTURAL BASIS FOR CROSS-SPECIES RECOGNITION OF COVID-19 VIRUS SPIKE
 TITLE 2 RECEPTOR BINDING DOMAIN TO BAT ACE2
 COMPD MOL_ID: 1;
 COMPD 2 MOLECULE: ANGIOTENSIN-CONVERTING ENZYME;
 COMPD 3 CHAIN: A;
 COMPD 4 EC: 3.4.-.-;
 COMPD 5 ENGINEERED: YES;
 COMPD 6 MOL_ID: 2;
 COMPD 7 MOLECULE: SARS-COV-2 RECEPTOR BINDING DOMAIN;
 COMPD 8 CHAIN: B;
 COMPD 9 FRAGMENT: UNP RESIDUES 333-527;
 COMPD 10 SYNONYM: S GLYCOPROTEIN,E2,PEPLOMER PROTEIN;
 COMPD 11 ENGINEERED: YES
 SOURCE MOL_ID: 1;
 SOURCE 2 ORGANISM_SCIENTIFIC: RHINOLOPHUS MACROTIS;
 SOURCE 3 ORGANISM_COMMON: BIG-EARED HORSESHOE BAT;
 SOURCE 4 ORGANISM_TAXID: 196889;
 SOURCE 5 GENE: ACE2;
 SOURCE 6 EXPRESSION_SYSTEM: ESCHERICHIA COLI;
 SOURCE 7 EXPRESSION_SYSTEM_TAXID: 562;
 SOURCE 8 MOL_ID: 2;
 SOURCE 9 ORGANISM_SCIENTIFIC: SEVERE ACUTE RESPIRATORY SYNDROME CORONAVIRUS
 SOURCE 10 2;
 SOURCE 11 ORGANISM_COMMON: 2019-NCOV;
 SOURCE 12 ORGANISM_TAXID: 2697049;
 SOURCE 13 GENE: S, 2;
 SOURCE 14 EXPRESSION_SYSTEM: INSECT CELL EXPRESSION VECTOR PTIE1;
 SOURCE 15 EXPRESSION_SYSTEM_TAXID: 266783
 KEYWDS COVID-19, RECEPTOR BINDING DOMAIN (RBD), RHINOLOPHUS MACROTIS, BATS,
 KEYWDS 2 ACE2, VIRAL PROTEIN-PROTEIN BINDING COMPLEX
 EXPDTA X-RAY DIFFRACTION
 AUTHOR K.F.LIU,J.WANG,S.G.TAN,S.NIU,L.L.WU,Y.F.ZHANG,X.Q.PAN,Y.M.MENG,
 AUTHOR 2 Q.CHEN,Q.H.WANG,H.W.WANG,J.X.QI,G.F.GAO
 REVDAT 2 03-FEB-21 7C8J 1 REMARK
 REVDAT 1 27-JAN-21 7C8J 0
 JRNL AUTH K.LIU,S.TAN,S.NIU,J.WANG,L.WU,H.SUN,Y.ZHANG,X.PAN,X.QU,P.DU,
 JRNL AUTH 2 Y.MENG,Y.JIA,Q.CHEN,C.DENG,J.YAN,H.W.WANG,Q.WANG,J.QI,
 JRNL AUTH 3 G.F.GAO
 JRNL TITL CROSS-SPECIES RECOGNITION OF SARS-COV-2 TO BAT ACE2.
 JRNL REF PROC.NATL.ACAD.SCI.USA V. 118 2021
 JRNL REFN ESSN 1091-6490
 PMID 33335073
 DOI 10.1073/PNAS.2020216118

SEQRES	1	A	707	THR	THR	GLU	ASP	GLU	ALA	LYS	LYS	PHE	LEU	ASP	LYS	PHE	HET	ZN	A	801	1			
SEQRES	2	A	707	ASN	SER	LYS	ALA	GLU	ASP	LEU	SER	TYR	GLU	SER	SER	LEU	HETNAM	ZN	ZINC	ION				
SEQRES	3	A	707	ALA	SER	TRP	ASP	TYR	ASN	THR	ASN	ILE	SER	ASP	GLU	ASN	FORMUL	3	ZN	ZN	2+			
SEQRES	4	A	707	VAL	GLN	LYS	MET	ASP	GLU	ALA	GLY	ALA	LYS	TRP	SER	ALA	HELIX	1	AA1	THR	A	20	ASN A 53 1	
SEQRES	5	A	707	PHE	TYR	GLU	GLU	GLN	SER	LYS	ALA	LYS	ASN	TYR	PRO		HELIX	2	AA2	ASP	A	56	ASN A 82 1	
SEQRES	6	A	707	LEU	GLU	GLU	ILE	GLN	ASN	ASP	THR	VAL	LYS	ARG	GLN	LEU	HELIX	3	AA3	ASN	A	90	SER A 103 1	
SEQRES	7	A	707	GLN	ILE	LEU	GLN	GLN	SER	GLY	SER	PRO	VAL	LEU	SER	GLU	HELIX	4	AA4	SER	A	109	GLY A 130 1	
SEQRES	8	A	707	ASP	LYS	SER	LYS	ARG	LEU	ASN	SER	ILE	LEU	ASN	ALA	MET	HELIX	5	AA5	PRO	A	146	SER A 155 1	
SEQRES	9	A	707	SER	THR	ILE	TYR	SER	THR	GLY	LYS	VAL	CYS	LYS	PRO	ASN	HELIX	6	AA6	ASP	A	157	GLU A 171 1	
SEQRES	10	A	707	ASN	PRO	GLN	GLU	CYS	VAL	LEU	GLU	PRO	GLY	LEU	ASP		HELIX	7	AA7	VAL	A	172	TYR A 194 1	
SEQRES	11	A	707	ASN	ILE	MET	GLY	THS	SER	LYS	ASP	TYR	ASN	GLU	ARG	LEU	HELIX	8	AA8	ASP	A	198	ARG A 204 1	
SEQRES	12	A	707	TRP	ALA	TRP	GLU	GLY	TRP	ARG	ALA	GLU	VAL	GLY	LYS	GLN	HELIX	9	AA9	ARG	A	205	GLU A 208 5	
SEQRES	13	A	707	LEU	ARG	PRO	LEU	TYR	GLU	GLU	TYR	VAL	VAL	LEU	LYS	ASN	HELIX	10	AB1	SER	A	218	TYR A 252 1	
SEQRES	14	A	707	GLU	MET	ALA	ARG	GLY	TYR	HIS	TYR	GLU	ASP	TYR	GLY	ASP	HELIX	11	AB2	TRP	A	275	ASN A 277 5	
SEQRES	15	A	707	TYR	TRP	ARG	ARG	ASP	TYR	GLU	THR	GLU	GLU	SER	SER	GLY	HELIX	12	AB3	LEU	A	278	VAL A 283 1	
SEQRES	16	A	707	PRO	GLY	TYR	SER	ARG	ASP	GLN	LEU	MET	LYS	ASP	VAL	ASP	HELIX	13	AB4	VAL	A	293	GLN A 300 1	
SEQRES	17	A	707	ARG	ILE	PHE	THR	GLU	ILE	LYS	PRO	LEU	TYR	GLU	HIS	LEU	HELIX	14	AB5	ASP	A	303	VAL A 318 1	
SEQRES	18	A	707	HIS	ALA	TYR	VAL	ARG	ALA	LYS	LEU	MET	ASP	THR	TYR	PRO	HELIX	15	AB6	THR	A	324	SER A 331 1	
SEQRES	19	A	707	LEU	HIS	ILE	SER	PRO	THR	GLY	CYS	LEU	PRO	ALA	HIS	LEU	HELIX	16	AB7	MET	A	366	TYR A 385 1	
SEQRES	20	A	707	LEU	GLY	ASP	MET	TRP	GLY	ARG	PHE	TRP	THR	ASN	LEU	TYR	HELIX	17	AB8	ALA	A	386	GLN A 388 5	
SEQRES	21	A	707	PRO	LEU	THR	VAL	PRO	PHE	GLY	GLN	LYS	TRP	ASP	ALA	ASN	HELIX	18	AB9	PRO	A	389	ARG A 393 5	
SEQRES	22	A	707	VAL	THR	ASP	ALA	MET	LEU	ASN	GLN	GLY	TRP	ASP	SER	VAL	HELIX	19	AC1	PHE	A	400	ALA A 413 1	
SEQRES	23	A	707	ARG	ILE	PHE	LYS	GLU	ALA	GLU	LYS	PHE	PHE	VAL	SER	VAL	HELIX	20	AC2	THR	A	414	THR A 420 1	
SEQRES	24	A	707	SER	LEU	PRO	LYS	MET	THR	GLU	GLY	TRP	TRP	ASN	LYS	SER	HELIX	21	AC3	ASP	A	431	VAL A 447 1	
SEQRES	25	A	707	MET	LEU	THR	GLU	PRO	GLY	ASP	GLY	ARG	LYS	VAL	VAL	CYS	HELIX	22	AC4	GLY	A	448	GLY A 466 1	
SEQRES	26	A	707	HIS	PRO	THR	ALA	TRP	ASP	LEU	GLY	LYS	GLY	PHE	ASP	ARG	HELIX	23	AC5	PRO	A	469	GLU A 471 5	
SEQRES	27	A	707	ILE	LYS	MET	CYS	THE	LYS	VAL	THR	MET	GLU	ASP	PHE	LEU	HELIX	24	AC6	GLU	A	472	ILE A 484 1	
SEQRES	28	A	707	THR	ALA	HIS	HIS	GLU	MET	GLY	ILE	GLN	TYR	ASP	MET		HELIX	25	AC7	CYS	A	498	SER A 502 5	
SEQRES	29	A	707	ALA	TYR	ALA	SER	GLN	PRO	TYR	LEU	LEU	ARG	ASN	GLY	ALA	HELIX	26	AC8	LEU	A	503	ASN A 508 1	
SEQRES	30	A	707	ASN	GLU	GLY	PHE	HIS	GLU	ALA	VAL	GLY	GLU	VAL	MET	SER	HELIX	27	AC9	PHE	A	512	ALA A 533 1	
SEQRES	31	A	707	LEU	SER	VAL	ALA	THR	PRO	LYS	His	LEU	LYS	THR	MET	GLY	HELIX	28	AD1	PRO	A	538	CYS A 542 5	
SEQRES	32	A	707	LEU	LEU	SER	PRO	PHE	ARG	GLU	ASP	ASP	GLU	THR	GLU		HELIX	29	AD2	SER	A	547	GLY A 561 1	
SEQRES	33	A	707	ILE	ASN	PHE	LEU	LYS	GLN	ALA	LEU	ASN	ILE	VAL	GLY		HELIX	30	AD3	ALA	A	565	ASP A 575 1	
SEQRES	34	A	707	THR	LEU	PRO	PHE	THR	TYR	MET	LEU	GLU	LYS	TRP	ASP	TRP	HELIX	31	AD4	VAL	A	581	ASN A 599 1	
SEQRES	35	A	707	MET	VAL	PHE	LYS	GLY	GLU	ILE	PRO	LYS	GLU	ILE	VAL	GLY	HELIX	32	AD5	ASN	A	636	LYS A 659 1	
SEQRES	36	A	707	LYS	LYS	TRP	TRP	GLU	MET	LYS	ARG	GLU	ILE	VAL	GLY	VAL	HELIX	33	AD6	GLY	A	666	GLU A 668 5	
SEQRES	37	A	707	VAL	GLU	PRO	VAL	PRO	HIS	ASP	GLU	THR	TYR	CYS	ASP	PRO	HELIX	34	AD7	PRO	A	696	ARG A 716 1	
SEQRES	38	A	707	ALA	SER	ILE	PHE	His	VAL	ALA	ASN	ASP	TYR	SER	PHE	ILE	HELIX	35	AD8	PRO	B	337	ASN B 343 1	
SEQRES	39	A	707	ARG	TYR	TYR	THR	ARG	THR	ILE	PHE	GLU	PHE	GLN	PHE	HIS	HELIX	36	AD9	ASP	B	364	ASN B 370 1	
SEQRES	40	A	707	GLU	ALA	LEU	CYS	ARG	ILE	ALA	GLN	HIS	ASN	GLY	PHE	LEU	HELIX	37	AE1	ASP	B	405	ILE B 410 5	
SEQRES	41	A	707	HIS	LYS	CYS	ASP	ILE	SER	ASN	SER	THR	ASP	ALA	GLY	LYS	HELIX	38	AE2	GLY	B	416	ASN B 422 1	
SEQRES	42	A	707	LYS	LEU	HIS	GLN	MET	LEU	SER	VAL	GLY	LYS	SER	GLN	ALA	HELIX	39	AE3	SER	B	438	SER B 443 1	
SEQRES	43	A	707	TRP	THR	LYS	THR	LEU	GLU	ASP	ILE	VAL	ASP	SER	ARG	ASN	HELIX	40	AE4	GLY	B	502	TYR B 505 5	
SEQRES	44	A	707	MET	ASP	VAL	GLY	PRO	LEU	LEU	ARG	TYR	PHE	LYS	PRO	LEU	SHEET	1	AA1	2	LYS	A	131	VAL A 132 0
SEQRES	45	A	707	TYR	THR	TRP	LEU	GLN	GLU	GLN	ASN	ARG	LYS	SER	TYR	VAL	SHEET	2	AA1	2	LEU	A	142	LEU A 143 -1
SEQRES	46	A	707	GLY	TRP	ASN	THR	ASP	TRP	SER	PRO	TYR	ALA	ASP	GLN	SER		O	LEU	A	142	N	VAL A 132	
SEQRES	47	A	707	ILE	LYS	VAL	TRP	ILE	SER	LEU	LYS	SER	ALA	LEU	GLY	GLU								

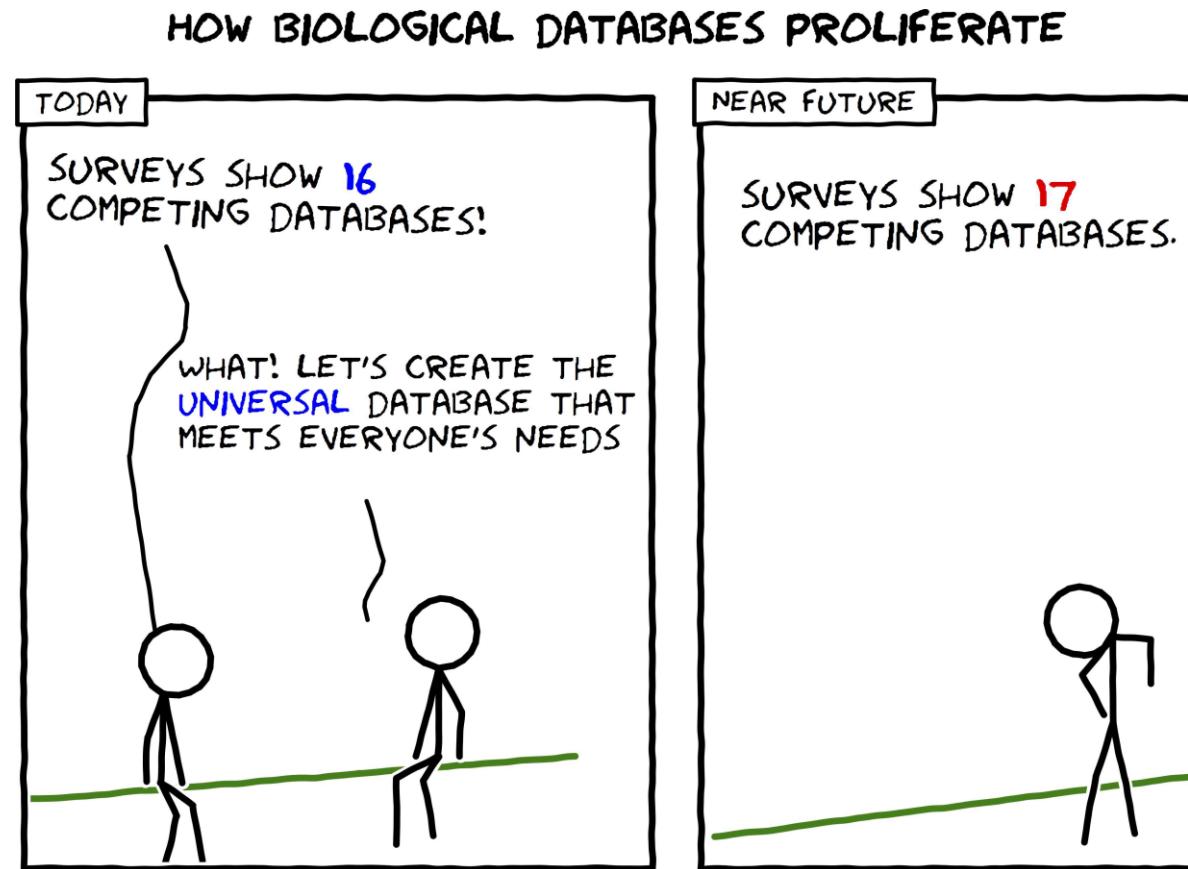
A sample of PDB file format

ATOM	1	N	PRO A	4	6.719	-12.134	26.603	1.00	18.91	N
ATOM	2	CA	PRO A	4	6.735	-10.746	27.122	1.00	18.45	C
ATOM	3	C	PRO A	4	6.209	-9.735	26.108	1.00	16.72	C
ATOM	4	O	PRO A	4	6.701	-9.658	24.983	1.00	16.64	O
ATOM	5	CB	PRO A	4	8.174	-10.427	27.495	1.00	20.82	C
ATOM	6	CG	PRO A	4	8.942	-11.387	26.584	1.00	20.17	C
ATOM	7	CD	PRO A	4	8.093	-12.664	26.557	1.00	22.00	C
ATOM	8	N	LEU A	5	5.207	-8.963	26.521	1.00	16.15	N
ATOM	9	CA	LEU A	5	4.605	-7.937	25.674	1.00	14.51	C
ATOM	10	C	LEU A	5	5.700	-6.960	25.244	1.00	14.38	C
ATOM	11	O	LEU A	5	6.564	-6.600	26.042	1.00	15.34	O
ATOM	12	CB	LEU A	5	3.513	-7.204	26.458	1.00	13.81	C
ATOM	13	CG	LEU A	5	2.639	-6.180	25.737	1.00	14.69	C
ATOM	14	CD1	LEU A	5	1.815	-6.864	24.656	1.00	15.29	C
ATOM	15	CD2	LEU A	5	1.725	-5.506	26.754	1.00	15.24	C

Ten Simple Rules for Developing Public Biological Databases

So, if you are considering developing a new database, and especially if you are a student or postdoc, please, for the love of science, follow these ten simple rules for creating and maintaining biological databases (and also a similar set of great rules for scientific web resources).

Rule 1: Don't reinvent the wheel



Ten Simple Rules for Developing Public Biological Databases

Rule 2: The three most important things in database development are data quality, data quality, and data quality

Rule 3: Know your users

Rule 4: Use modern technology

Rule 5: Put yourself in your user's shoes

Rule 6: Keep search simple and organized

Rule 7: Give users data where they need it

Rule 8: Support open science

Rule 9: Tell the world

Rule 10: Maintain, update, or retire

Thank You