

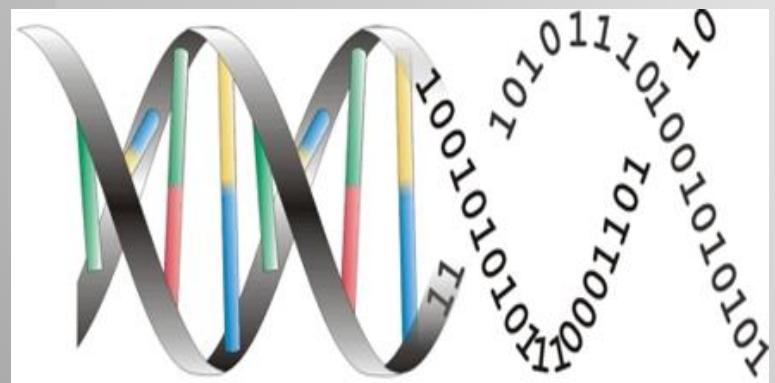


The University of Georgia

Programming for Computational & Systems Biology

Instructor: Paul Xie

Tue. & Thr. 9:35~10:50



Course Format

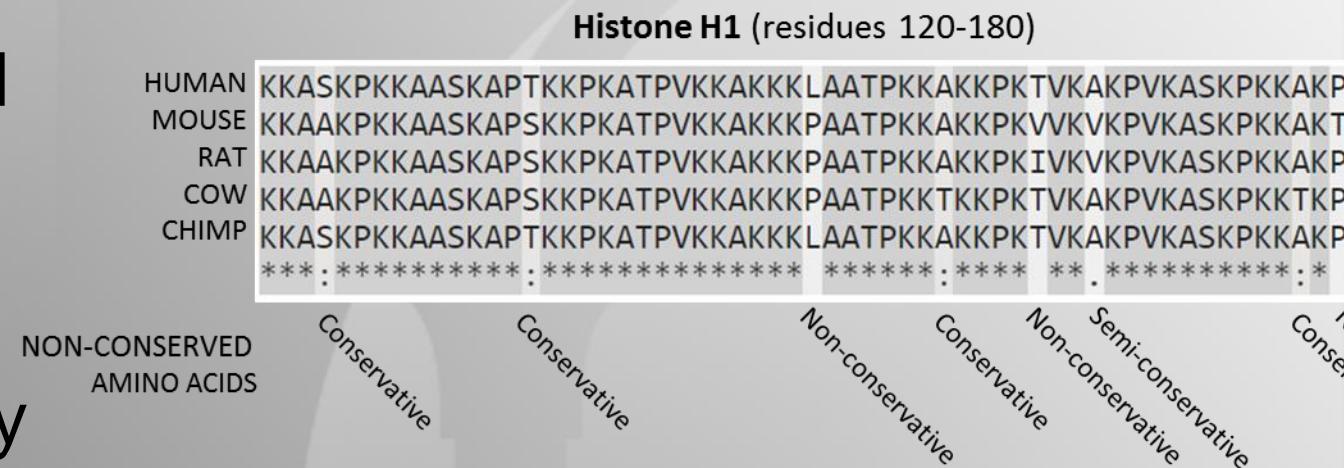
- 1 hour lecture + 2 hour computer lab per week
- 2 topics/week; 1 computational biology + 1 coding concept
- 1 assignment/week, please check eLC
- 1 term paper (50%)

Databases & Data Retrieval

Instructor: Paul Xie (6)

Last Weeks

- Protein alignment
 - Translation
 - PAM
 - BLOSUM
- Coding
 - List
 - Dictionary
 - Function



This Week

- Biological Databases
- Data retrieval
- How to obtain the data you need in your research project

Introduction

- Fast increase in biological information
- Biological science has now turned into a data rich science
- Gene sequences
- Amino acid sequences in proteins
- Motifs and domains in proteins
- Structural data from XRD & NMR
- Metabolic pathways
- Protein-protein interactions
- Gene expression data DNA microarrays

WHY

- What is a database and what are the features of an ideal database?
- What are the relationships/differences between primary and secondary (derived) sequence databases?
- What are the benefits of RefSeq?
- Why is data integration useful?

ADVANTAGES OF DATA INTEGRATION



- More relevant *inter-related* information in one place
- Makes it easier to find additional relevant information related to your initial query
- Potentially find information *indirectly* linked, but *relevant* to your subject of interest
 - uncover *non-obvious* genetic features that explain phenotype or disease
- Easier to build a ‘story’ based on *multiple* pieces of biological evidence

WHAT

- **Structured** collection of information.
- Consists of basic units called records or entries.
- Each record consists of fields, which hold **pre-defined** data related to the record.
- 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)

Biological databases

- Biological database is a collection of data which is structured, searchable, updated periodically and also cross-referenced.
- Some databases are multi functional
- Major purposes of databases is as follows:

Availability of
biological data

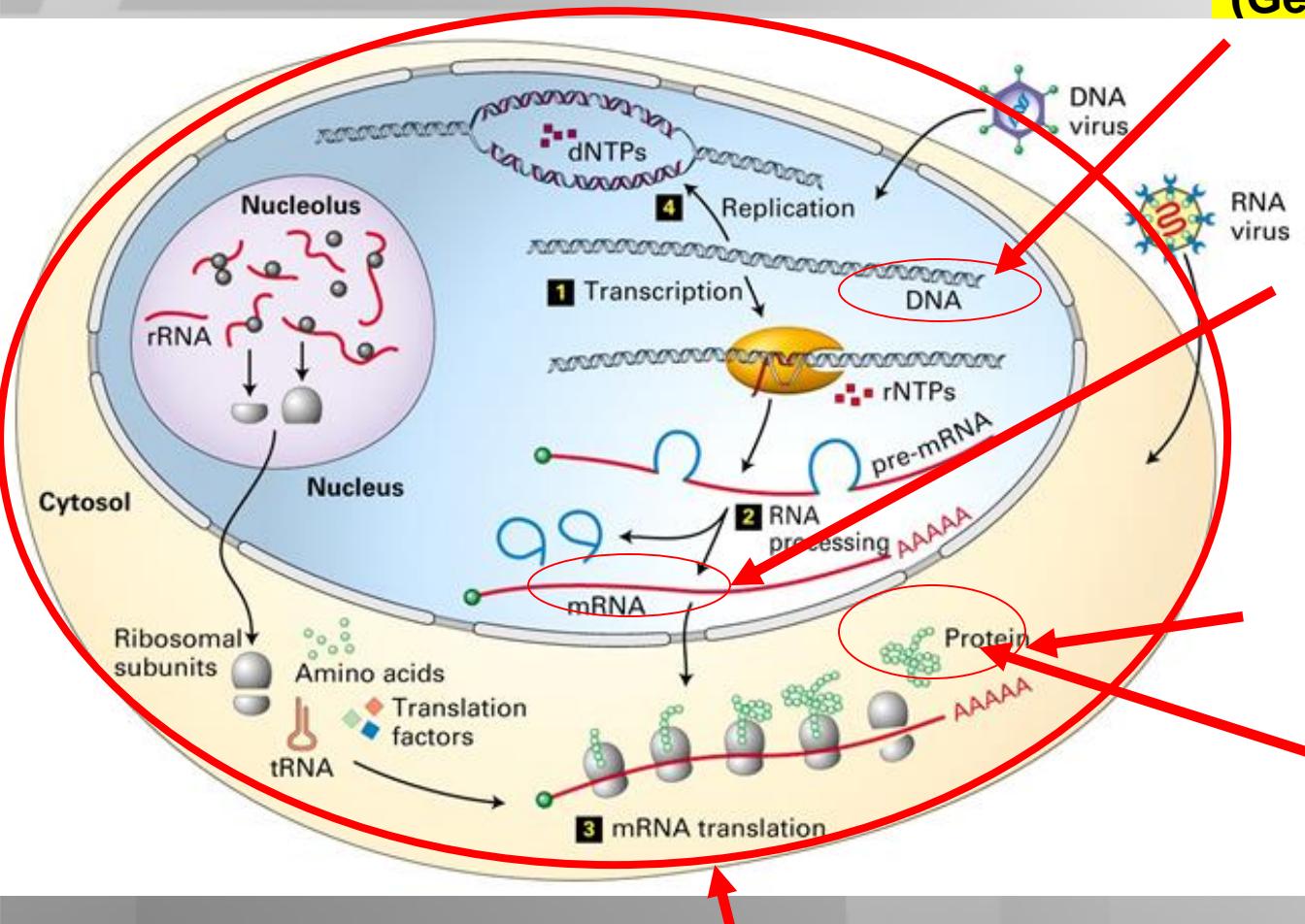
Systemization
of data

Analysis of
computed
biological data

Data type

- Genome database
- Sequence database
- Structure database
- Microarray database
- Chemical database
- Pathway database
- Enzyme database
- Disease database
- Literature database

The Central Dogma & Biological Data



History

- 1956; first sequence database when insulin was sequenced
- 51 amino acids
- Atlas of protein sequences and structures in 1965 by Margaret Day Hoff et al was a printed book.
- Became base for PIR protein information resource
- First nucleotide sequence: yeast tRNA
- 77 bases
- During this time 3D structure of proteins was being studied and renowned PDB was made.

History (2)

- First genome published was of free living virus *haemophilus influenzae* in 1995
- Genome?
- All genes ? Or all DNA?
- Why are complete genome interesting?

Biological Databases



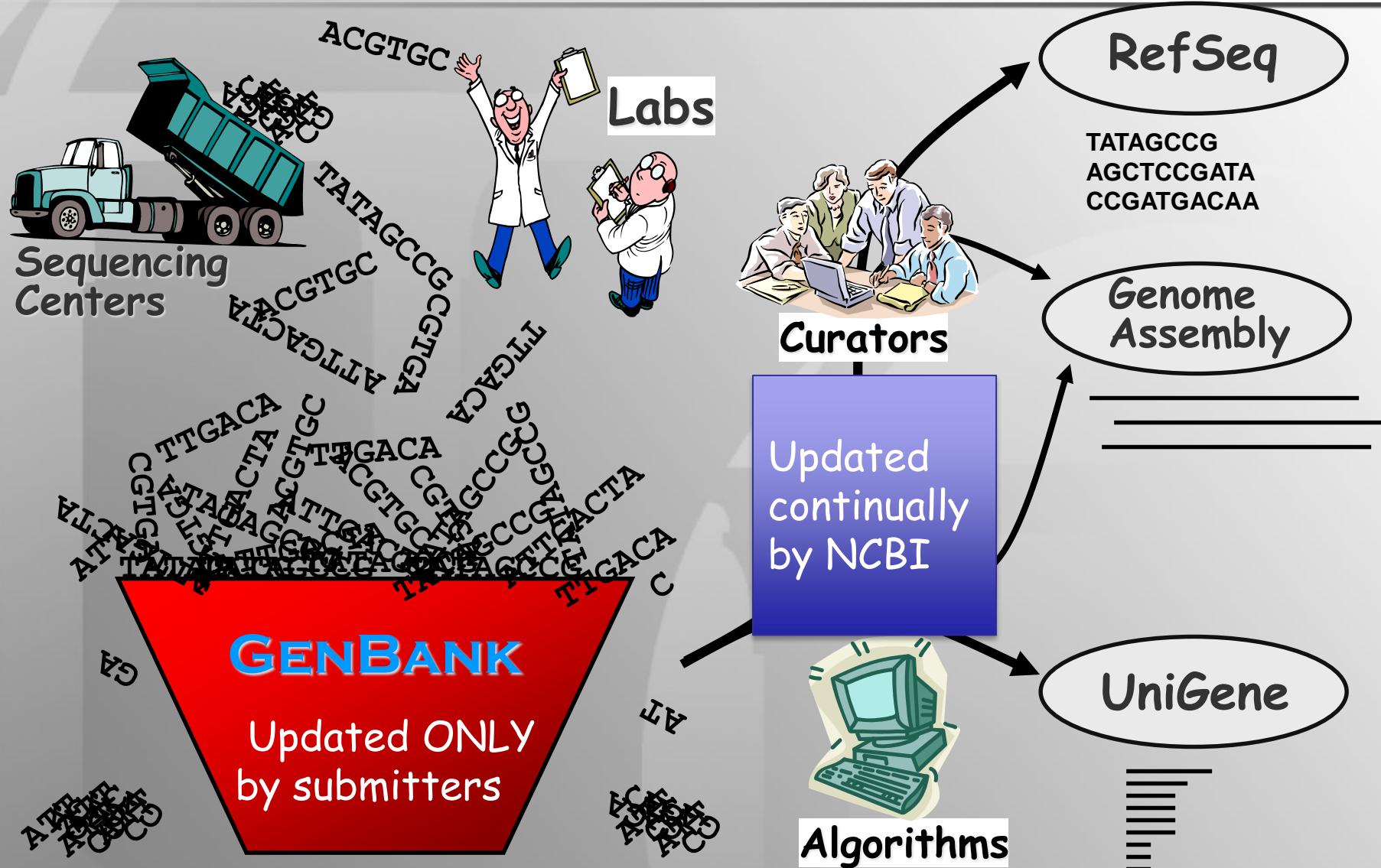
- Primary Databases
- Secondary Databases

Primary and Secondary Databases



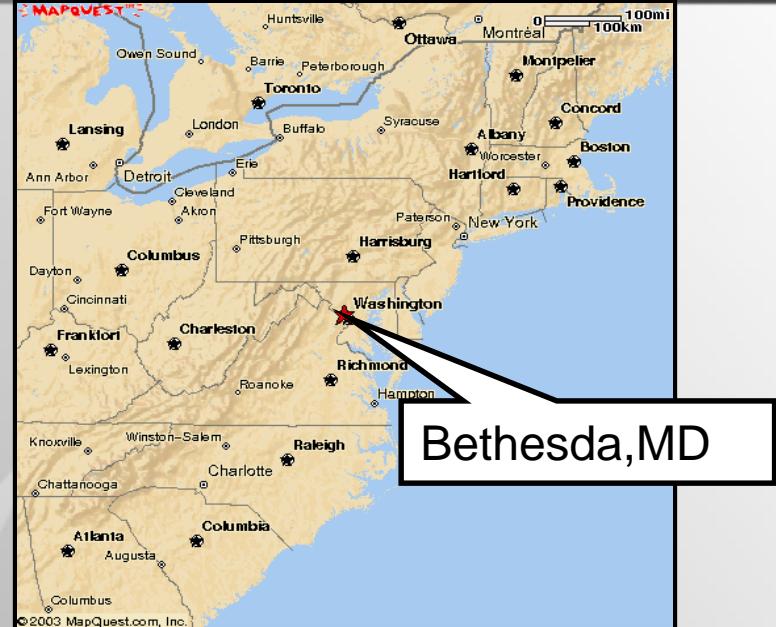
- Primary (archival)
 - GenBank/EMBL/DDBJ
 - UniProt (EBI:European Bioinformatics Institute SIB:Swiss Institute of Bioinformatics, PIR:Protein Information Resource)
 - PDB
 - Medline (PubMed)
 - BIND
- Secondary (curated)
 - RefSeq
 - Taxon
 - UniProtKB
 - OMIM (online Mendelian Inheritance in Man)

PRIMARY VS. DERIVATIVE SEQUENCE DATABASES



The National Center for Biotechnology Information

The University of Georgia



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

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

Web Access: The University of Georgia

www.ncbi.nlm.nih.gov

 NCBI Resources How To My NCBI | Sign In

National Center for Biotechnology Information Search All Databases for New pages!

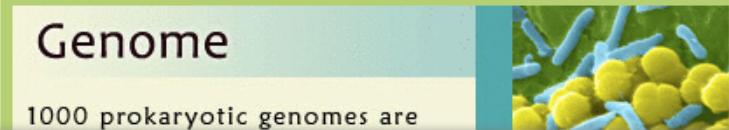
Resources

- NCBI Home
- All Resources (A-Z)
- Literature
- DNA & RNA
- Proteins
- Sequence Analysis
- Genes & Expression
- Genomes

Welcome to NCBI

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

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


1000 prokaryotic genomes are 

Popular Resources

- PubMed
- PubMed Central
- Bookshelf
- BLAST
- Gene
- Nucleotide
- Protein
- GEO
- Conserved Domains
- Structure

Maps & Markers You are here: NCBI Help Desk

GETTING STARTED	RESOURCES	POPULAR	FEATURED	NCBI INFORMATION
Site Map	Literature	PubMed	GenBank	About NCBI
NCBI Help Manual	DNA & RNA	PubMed Central	Reference Sequences	Research at NCBI
NCBI Handbook	Proteins	Bookshelf	Map Viewer	NCBI Newsletter
Training & Tutorials	Sequence Analysis	BLAST	Genome Projects	NCBI FTP Site
	Genes & Expression	Gene	Human Genome	Contact Us
	Genomes	Nucleotide	Mouse Genome	
	Maps & Markers	Protein	Influenza Virus	
	Domains & Structures	GEO	Primer-BLAST	
	Genetics & Medicine	Conserved Domains	Short Read Archive	
	Taxonomy	Structure		
	Data & Software	PubChem		
	Training & Tutorials			
	Homology			
	Small Molecules			
	Variation			

Common footer

NCBI Databases and Services

- GenBank primary sequence database
- Free public access to biomedical literature
 - PubMed free Medline (3 million searches per day)
 - PubMed Central full text online access
- Entrez integrated molecular and literature databases

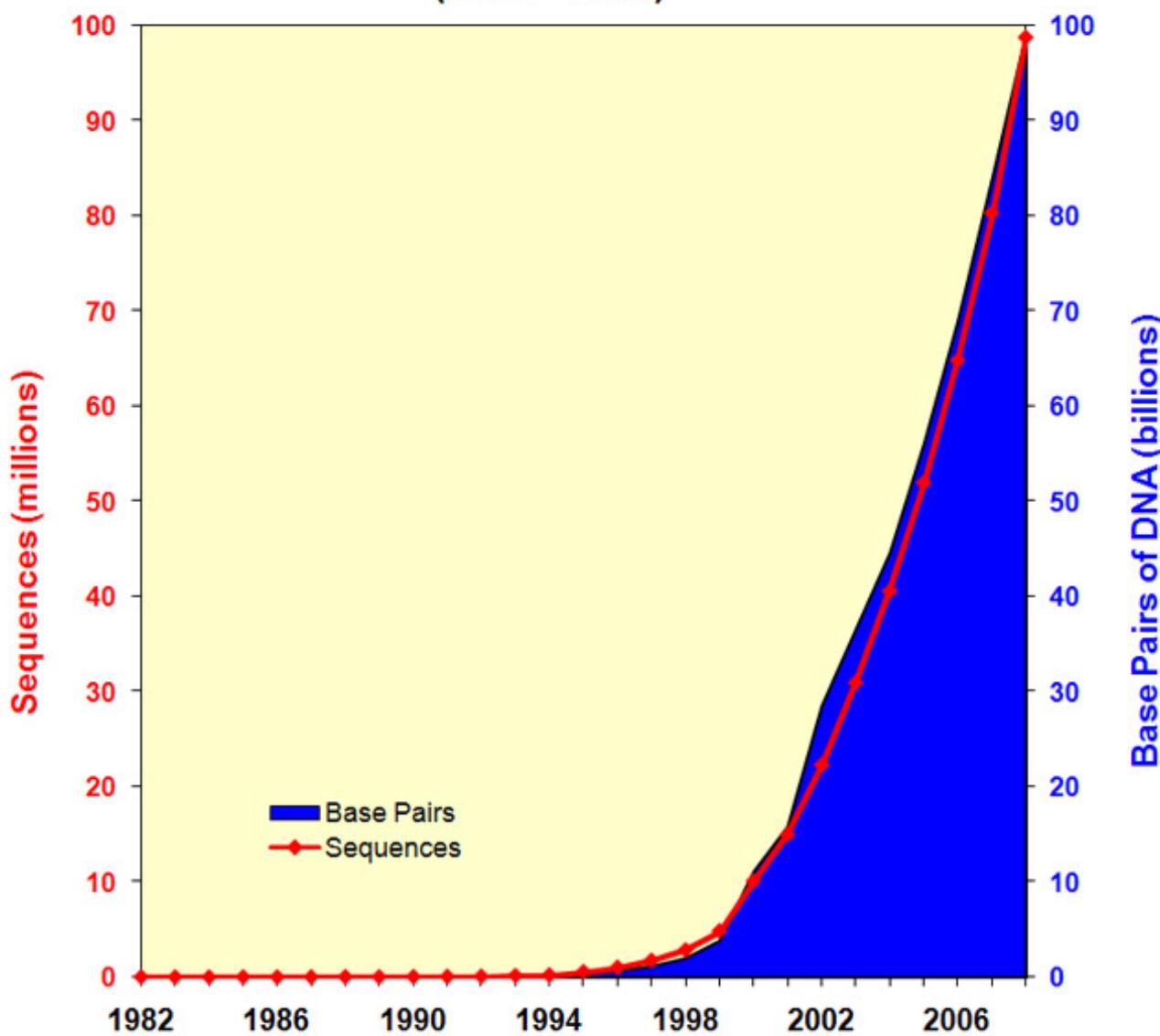
GENBANK - PRIMARY SEQUENCE DB

- **Nucleotide only** sequence database
- **Archival** in nature
 - Historical
 - Reflective of submitter point of view (subjective)
 - **Redundant**
- Data
 - Direct submissions (traditional records)
 - Batch submissions
 - FTP accounts (genome data)

GENBANK - PRIMARY SEQUENCE DB (2)

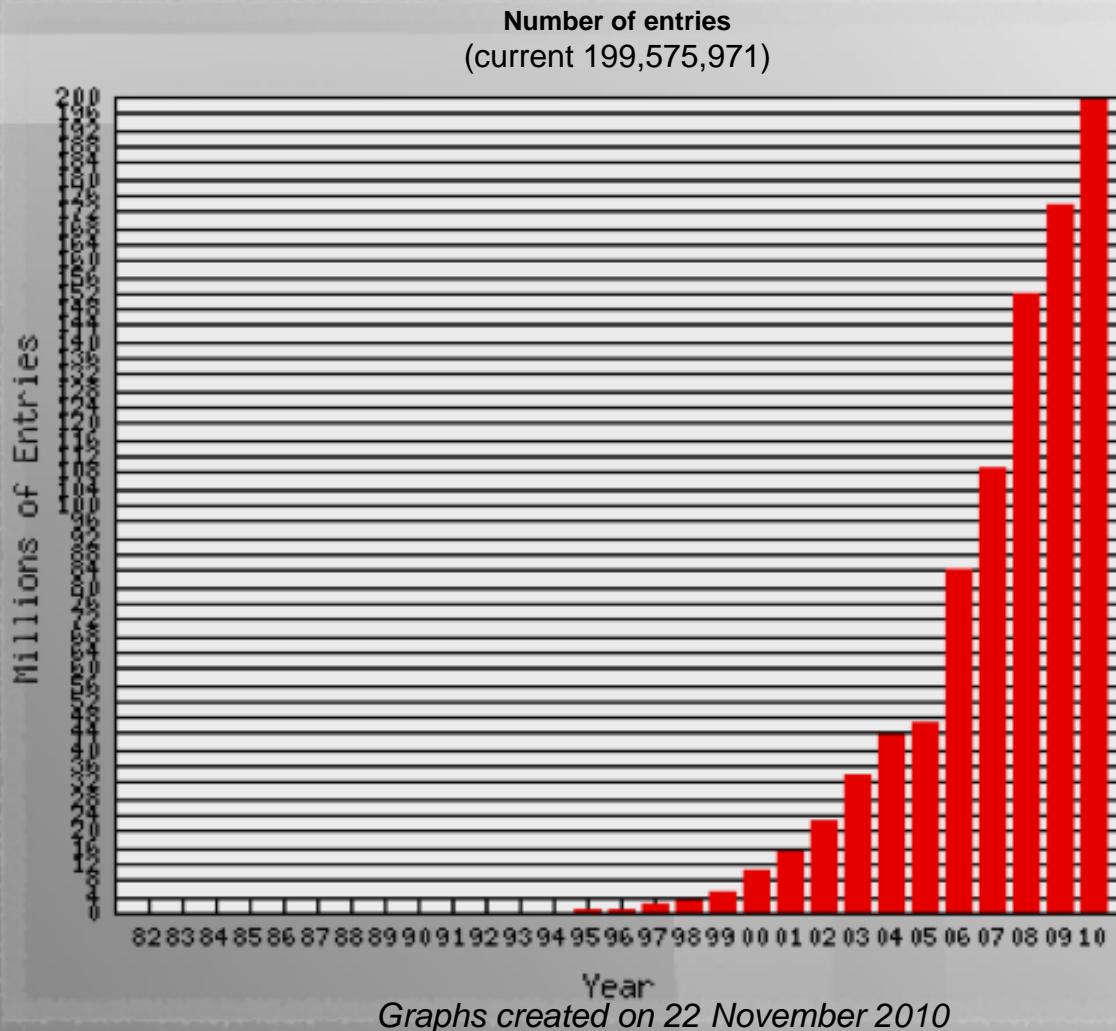
- Three collaborating databases
 1. GenBank
 2. DNA Database of Japan (DDBJ)
 3. European Molecular Biology Laboratory (EMBL) Database

Growth of GenBank (1982 - 2008)





EMBL Database



Graphs created on 22 November 2010

National Center for
Biotechnology Information

NCBI Home

Site Map (A-Z)

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials

Variation

Search All Databases

Search Clear

Welcome to NCBI

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

[About the NCBI](#) | [Mission](#) | [Organization](#) | [Research](#) | [RSS Feeds](#)

Get Started

- » [Tools](#): Analyze data using NCBI software
- » [Downloads](#): Get NCBI data or software
- » [How-To's](#): Learn how to accomplish specific tasks at NCBI
- » [Submissions](#): Submit data to GenBank or other NCBI databases

Human Microbiome Project



NIH Roadmap Initiative designed to characterize the community of microorganisms living on and in the human body.

II 1 2 3 4 5

Popular Resources

- » [BLAST](#)
- » [Bookshelf](#)
- » [Gene](#)
- » [Genome](#)
- » [Nucleotide](#)
- » [OMIM](#)
- » [Protein](#)
- » [PubChem](#)
- » [PubMed](#)
- » [PubMed Central](#)
- » [SNP](#)

NCBI News

[Retirement of Peptidome, SRA & Trace Archive](#)

16 Feb 2011

Due to budget constraints, NCBI will be discontinuing the

[The Bookshelf has a new design & Browsing Tool](#)

09 Feb 2011

Featuring a new homepage, search results display, limits and

[More...](#)



Entrez, The Life Sciences Search Engine.

HOME SEARCH SITE MAP

PubMed

All Databases

Human Genome

GenBank

Map Viewer

BLAS

Search across databases

GO

Clear

Help

Welcome to the Entrez cross-database search page

PubMed: biomedical literature citations and abstracts	Books: online books
PubMed Central: free, full text journal articles	Images: images from full text resources at NCBI
Site Search: NCBI web and FTP sites	OMIM: online Mendelian Inheritance in Man

Nucleotide: Core subset of nucleotide sequence records	dbGaP: genotype and phenotype
EST: Expressed Sequence Tag records	UniGene: gene-oriented clusters of transcript sequences
GSS: Genome Survey Sequence records	CDD: conserved protein domain database
Protein: sequence database	UniSTS: markers and mapping data
Genome: whole genome sequences	PopSet: population study data sets
Structure: three-dimensional macromolecular structures	GEO Profiles: expression and molecular abundance profiles
Taxonomy: organisms in GenBank	GEO DataSets: experimental sets of GEO data
SNP: single nucleotide polymorphism	Epigenomics: Epigenetic maps and data sets
dbVar: Genomic structural variation	Cancer Chromosomes: cytogenetic databases
Gene: gene-centered information	PubChem BioAssay: bioactivity screens of chemical substances

Traditional GenBank Record

The University
of Georgia



LOCUS	HSHMLHI	2503 bp	mRNA	linear	PRI	31-MAR-1994	
DEFINITION	Human DNA misma		FEATURES	Location/Qualifiers			
ACCESSION	U07418		source	1..2503			
VERSION	U07418.1	GI:46		/organism="Homo sapiens"			
KEYWORDS	.			/db_xref="taxon:9606"			
SOURCE	Homo sapiens (h			/chromosome			
ORGANISM	<u>Homo sapiens</u>			/map=			
	Eukaryota; Meta			BASE COUNT 723 a 539 c 599 g 642 t			
	Mammalia; Euthe			ORIGIN			
REFERENCE	1 (bases 1 to 2503)			1 gttgaacatc tagacgttc ctgggcttt ctggcgccaa aatgtcggttgcgtgggg			
A			gene	61 ttattccggc gctggacgag acagtggta accgcattcgccggggggaa gttatccagc			
			CDS	121 ggccagcta tgctatccaa gagatgattt agaactgttt agatgcacaaa tccacaagta			
				181 ttcaatggat tgttaaaggag ggaggctgtt gatccaaagac aatggccaccc			
				241 ggatcggaa aagaatctg gatattgtat gtggaaagggtt cactactgtt aaatgcacgt			
				301 cctttggaga tttagccagt atttctacct atgggtttcg aggtggggct ttggccacca			
				361 taaggcatgt ggctcatgtt actattacaa cggaaacacgc tgatggaaatgtgtgcataca			
				421 gagcaagta ctcagatggaa aactgtaaag cccctcccaa accatgtgttgcgttgcacaa			
				481 ggaccaggat cacgggtggg gaccccccataaacatagc caccggggaa aaagcttta			
				541 aaaaatccaaatg tggaaatataatggggaaatttt tgccggatgtt ccgtatcacaa			
				601 atgcaggcat tagttctca gttaaaaaaaac aaggagagac agtagctgtt gttggacac			
				661 tacccatgc ctcacccgtg gacaaatattc gctccgtttt tgaaaatgtgtt gtttgtcgac			
				721 aactgtatggaaatttggatgtt gaggataaaa cccttagcctt caaaatgttgcgttccat			
				781 ccaatgcacaaatctactgtt aagaatgtca ttccctttact ctccatcaac catcgctgg			
				841 tagaaatcaac ttcccttgaga aaagccatag aaacagtgtt gtcagccat ttggccaaaa			
				901 acacacaccc attccgttac ctcaatgtttag aaatcaatgttcc ccagaatgtt gatgttaatg			
				961 tgcacccccc aaacatgttac gttcaatcttc tgccggggaa gggccatccgttgcgttgc			
				1021 agcggccatcg cgagccatcg ctccctgggtt ccaattccctc caggatgttac ttccatccacaa			
				1081 ctttgcgtacc aggactgttgc gggccctctg gggagatgtt cttatccacaa acaatgttgc			
				1141 cctctgttcc tacttcttgc agtagtgcata aggtctatgc ccaccatgtt gttctgttgc			
				1201 attccccggaa acacaaatgtt gatccatcc ttccctgttcc gggccaaaccc ctgtccatgttgc			
				1261 agccccccgtt cattgttccaa gggatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1321 aagatggaaatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1381 tggggggggatacaacaaatggggacttgc aatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1441 gcaacccccc aaagatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1501 gaaaggaaatgtt gatctgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1561 tgatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1621 accacttccctgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1681 tataccatccca aacaccacc aagatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1741 atttttggccatccatccca aacaccacc aagatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1801 tgcttgcctt aatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1861 ttgctgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1921 ctttggaaatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				1981 tgccccccttggggacttgcataatgttcc ttccctgttcc gttccatgttgc			
				2041 acgaagaaaaaa gggatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2101 ggaaggatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2161 ccattccaaatccctgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2221 acatccatccca aacaccacc aagatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2281 ctgtatcttgcataatgttcc ttccctgttcc gttccatgttgc			
				2341 gttcttttctctgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2401 accaacaatataatgttgcataatgttcc ttccctgttcc gttccatgttgc			
				2461 tacacatgttgcataatgttcc ttccctgttcc gttccatgttgc			

well annotated

the sequence is the data

GenPept: GenBank CDS translations



FEATURES	Location/Qualifiers
source	1..2484 /organism="Homo sapiens" /mol_type="mRNA" /db_xref="taxon:9606" /chromosome="3" /map="3p22-p23"
gene	1..2484 /gene="M"
CDS	>gi 463989 gb AAC50285.1 DNA mismatch repair prote... 22..2292 MSFVAGVIRRLDETVVNRIAAGEVIQRPANAIKEMIENCLDAKSTSIV... /gene="M" EDLDIVCERFTTSKLQSFEDLASISTYGFRGEALASISHVAHVTITTKTAD... /note="homolog of S. cerevisiae PMS1 (Swiss-Prot Accession Number P14242), cerevisiae MLH1 (GenBank Accession Number U07187), coli MUTL (Swiss-Prot Accession Number P23367), Salmonella typhimurium MUTL (Swiss-Prot Accession Number P14161) and Streptococcus pneumoniae (Swiss-Prot Accession Number P14160)" /codon_start=1 /product="DNA mismatch repair protein homolog" /protein_id="AAC50285.1" /db_xref="GI:463989" /translation="MSFVAGVIRRLDETVVNRIAAGEVIQRPANAIKEMIENCLDAKS... TSIQIVKEGLKLIQIQDNGTGIRKEDLDIVCERFTTSKLQSFEDLASISTYGFRGE... ALASISHVAHVTITTKTADGKCAYRASYSDGKLKAPPKPCAGNQGTQITVEDLFYNIA... TRRKALKNPSEEYGKILEVVGRYSVHNAGISFSVKKQGETVADVRTLPNASTVDNIRS..."

REFSEQ: *DERIVATIVE SEQUENCE DATABASE*



- Curated transcripts and proteins
- Model transcripts and proteins
- Assembled Genomic Regions
- [C ftp://ftp.ncbi.nih.gov/refseq/release/](ftp://ftp.ncbi.nih.gov/refseq/release/)
 - microbial
 - organelle

Selected RefSeq Accession Numbers

mRNAs and Proteins

NM_123456

Curated mRNA

NP_123456

Curated Protein

NR_123456

Curated non-coding RNA

XM_123456

Predicted mRNA

XP_123456

Predicted Protein

XR_123456

Predicted non-coding RNA

Gene Records

NG_123456

Reference Genomic Sequence

Chromosome

NC_123455

Microbial replicons, organelle

AC_123455

Alternate assemblies

Assemblies

NT_123456

Contig

NW_123456

WGS Supercontig

GenBank to RefSeq

[Human apolipoprotein E \(epsilon-4 allele\) gene, complete cds](#)
1. 5,515 bp linear DNA
M10065.1 GI:178852

[Human mRNA fragment for apolipoprotein E \(apo E\)](#)
2. 528 bp linear mRNA
X00199.1 GI:28808

[H.sapiens mRNA](#)
3. 275 bp linear mRNA
Z70760.1 GI:12631

[Homo sapiens cDNA](#)
4. 1,023 bp linear mRNA
AK314898.1 GI:16

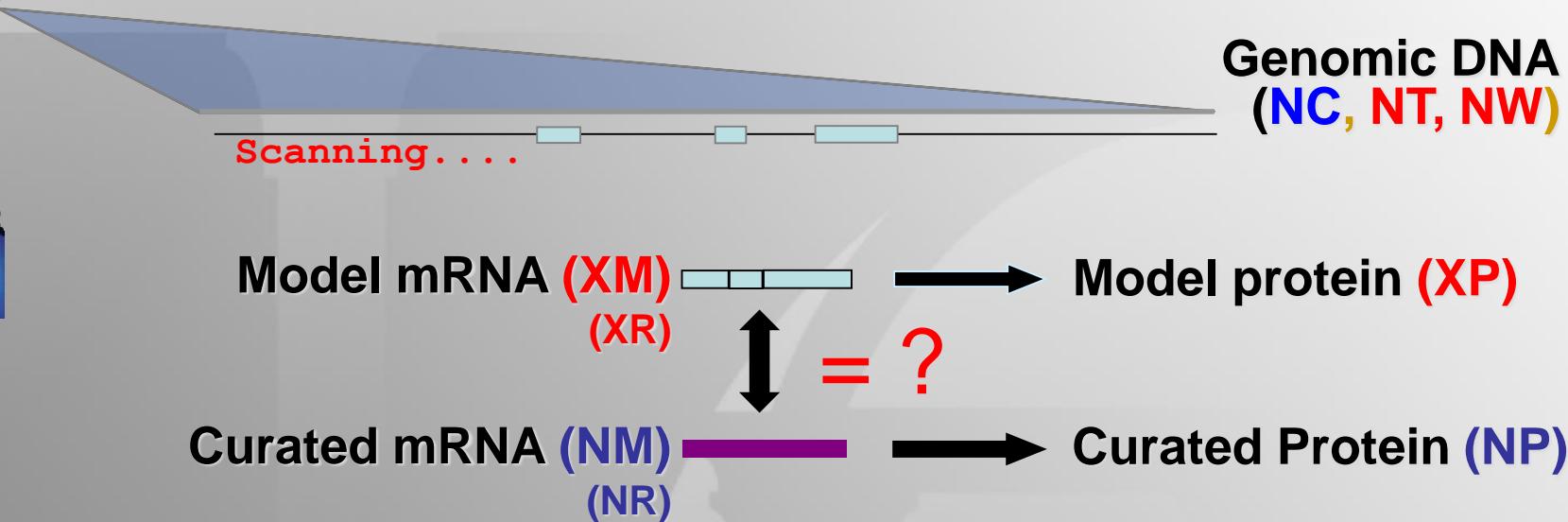
[Homo sapiens apolipoprotein E \(APOE\), mRNA](#)
1,223 bp linear mRNA
NM_000041.2 GI:48762938

[Human apolipoprotein E mRNA, complete cds](#)
5. 1,157 bp linear mRNA
M12529.1 GI:178848

[Homo sapiens preapolipoprotein E \(APOE\) mRNA, complete cds](#)
6. 1,156 bp linear mRNA
K00396.1 GI:178850

[Homo sapiens apolipoprotein E, mRNA \(cDNA clone MGC:1571 IMAGE:3355712\), complete cds](#)
7. 1,186 bp linear mRNA
BC003557.1 GI:13097698

RefSeqs: Annotation Reagents



RefSeq

GenBank
Sequences



RefSeq Benefits

- Non-redundancy
- Updates to reflect current sequence data and *b*
- Data *validation*
- Format *consistency*
- Distinct accession series

- **PubMed** is a free search engine accessing primarily the MEDLINE database of references and abstracts on life sciences and biomedical topics. The United States National Library of Medicine (NLM) at the National Institutes of Health maintains the database as part of the Entrez system of information retrieval.

Databases Architecture



Information system

Query system

Storage System

Data

Databases Architecture

Information system

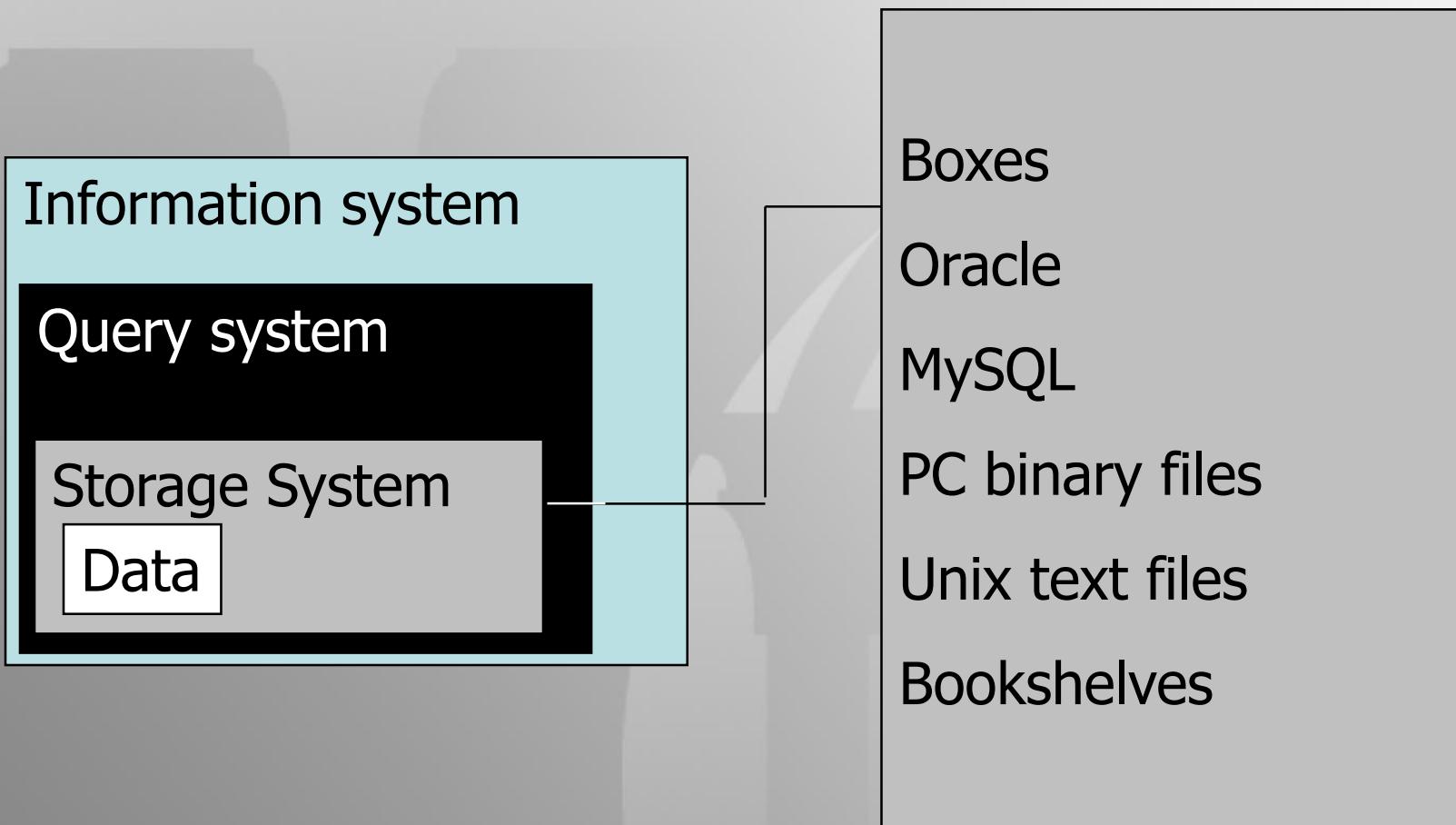
Query system

Storage System

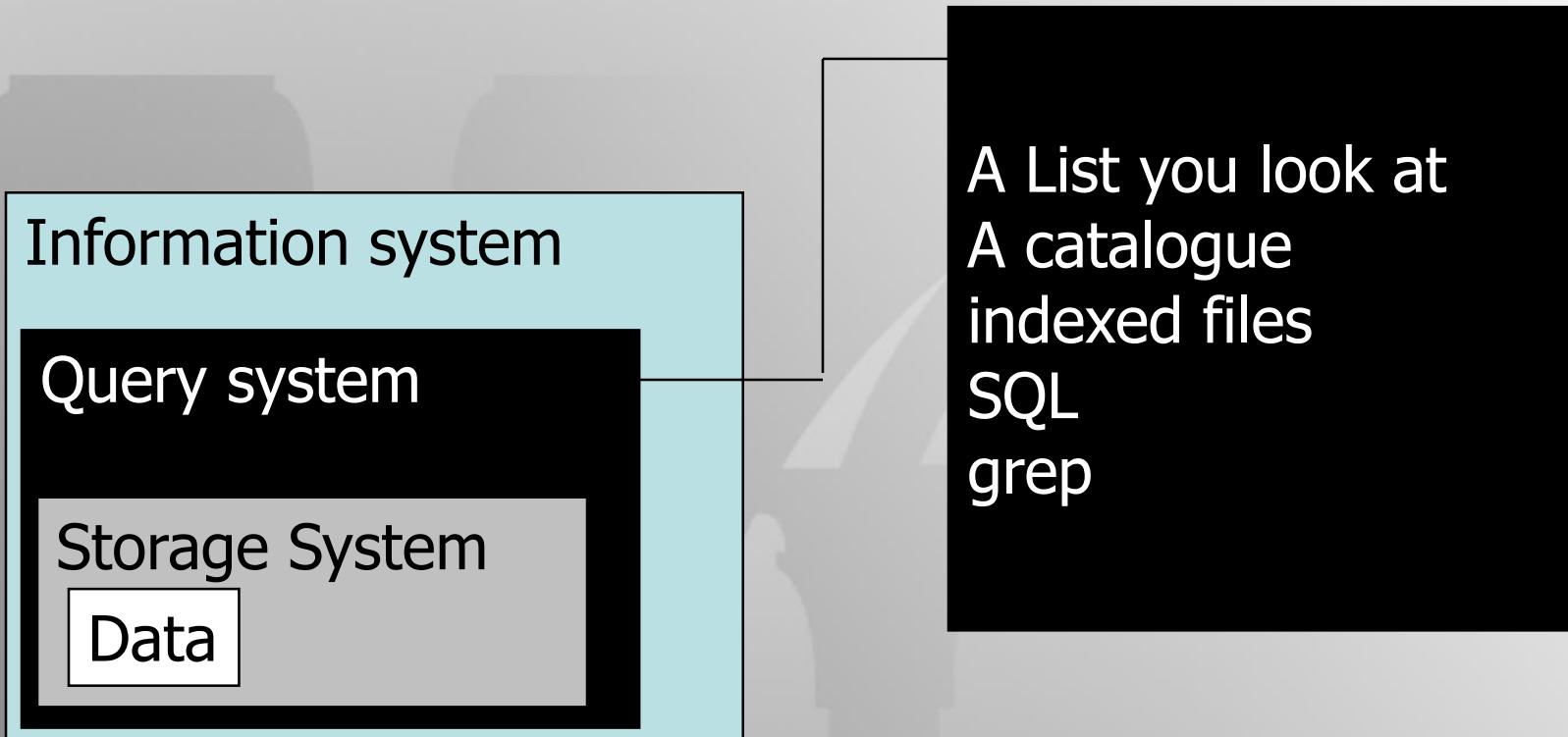
Data

GenBank flat file
PDB file
Interaction Record
Title of a book
Book

Databases Architecture



Databases Architecture



Databases Architecture



Information system

Query system

Storage System

Data

The Google
Entrez
SRS

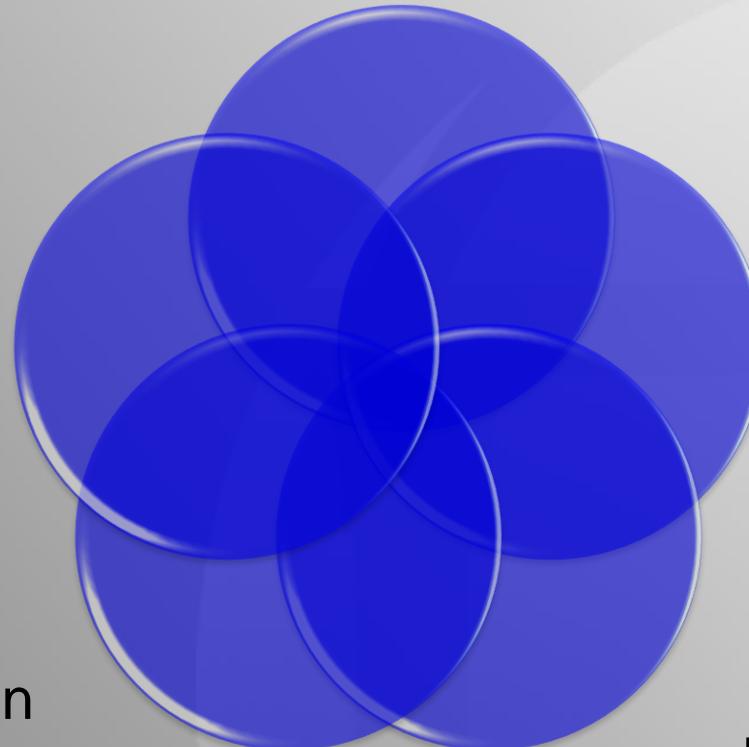
Aspects of genome analysis

Gene prediction
via comparison,
coding and
regulatory
regions

Gene prediction
via EST

Ab initio Gene
prediction

Locus



Gene
identification by
EST (expressed
sequence tags)

Features of biological databases

- 1) Data heterogeneity
- 2) High volume data
- 3) Uncertainty
- 4) Data Curation
- 5) Large scale data integration
- 6) Data sharing
- 7) Dynamic and subject to change

Classification scheme for biological databases

Data type

Maintenance status

Data access

Data source

Database design

Organism

THE ‘PERFECT’ DATABASE



- Comprehensive, but **easy to search**.
- Annotated, but not “too annotated”.
- A simple, easy to understand structure.
- **Cross-referenced**.
- Minimum redundancy.
- **Easy retrieval** of data.

Pitfalls of Biological Databases

- Errors in Sequence Databases
- Redundancy in the Primary Sequence Databases
- False or Incomplete Genes Annotations

Errors in Nucleotide Sequences

- sequencing errors
- frame-shifts
- Contaminated with sequences from cloning vectors
 - Exceptional Care for sequences produced before the 1990s

- repeated submission
 - identical or overlapping sequences by the same or different authors
 - revision of annotations
 - dumping of expressed sequence tags (EST) data
 - poor database management

NAR Articles



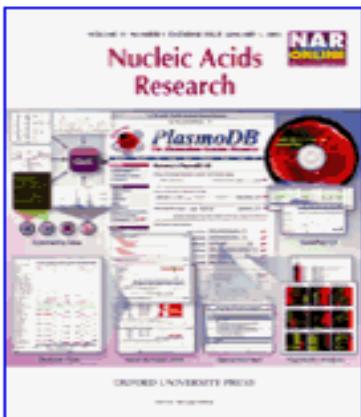
Nucleic Acids Research

OXFORD
Journals online

[HOME](#) [HELP](#) [FEEDBACK](#) [SUBSCRIPTIONS](#) [ARCHIVE](#) [SEARCH](#) [TABLE OF CONTENTS](#)

[B F Francis Ouelette](#) || [Change Password](#) || [View/Change User Information](#) || [CiteTrack Personal Alerts](#) || [Subscription HELP](#) || [Sign Out](#)

Receive this page by email each issue: [\[Sign up for eTOCs\]](#)



[\[Cover Caption\]](#)

Other Issues:



Contents: Volume 31, Number 1 January 1 2003 [\[Index by Author\]](#)

- [Editorial](#)
- [Articles](#)

Find articles in this issue containing these words:

Enter

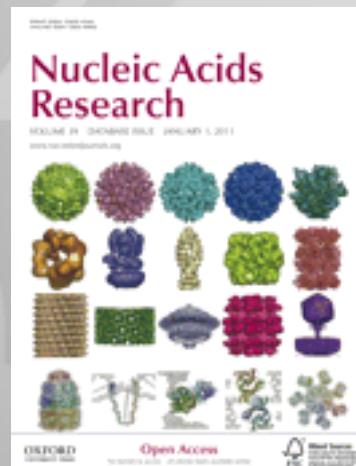
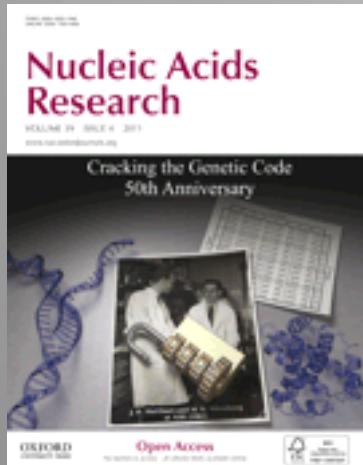
[\[Search ALL Issues\]](#)

<http://nar.oupjournals.org/content/vol31/issue1/>

The definitive source....



- More than 1300 DB
- http://nar.oxfordjournals.org/content/39/suppl_1.toc



DNA Sequence databases

- Main repositories:
 - GenBank (US)
 - (<http://www.ncbi.nlm.nih.gov/Genbank/index.html>)
 - EMBL (Europe)
 - (<http://www.ebi.ac.uk/embl/>)
 - DDBJ (Japan)
 - (<http://www.ddbj.nig.ac.jp/>)
- Primary databases
 - DNA sequences are identical



PubMed is...

- National Library of Medicine's search service
 - >20 million citations in MEDLINE
 - links to participating online journals
 - PubMed tutorial (via side bar)

All Databases

PubMed

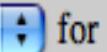
Nucleotide

Protein

Genome

Search

PubMed



for

Go

Clear

Limits

Preview/Index

History

Clipboard

Details

Entrez integrates...

- the scientific literature;
- DNA and protein sequence databases;
 - 3D protein structure data;
 - population study data sets;
- assemblies of complete genomes

Genome Databases

- Focus on one organism or group of organisms:
 - Colibase (*E. coli* and related species) <http://colibase.bham.ac.uk/>
 - GDB (human) <http://www.gdb.org/>
 - Flybase (*Drosophila*) <http://flybase.bio.indiana.edu/>
 - WormBase (*C. elegans*) <http://wormbase.org>
 - AtDB (*Arabidopsis*) <http://www.arabidopsis.org>
 - SGD (*S. cerevisiae*) <http://genome-www.stanford.edu/Saccharomyces/>



SWISSPROT

- European/Swiss Bioinformatics Institute 1986
- Highly accurate, hand curated resource
- Aims:
 - Have a high level of annotation
 - Often by the people who have been working with the gene
 - Have a low level of redundancy
 - Have a high level of integration with other databases

<http://www.ebi.ac.uk/trembl/>

- SWISSPROT's Big Brother
 - All genes which have been left out of SWISSPROT
 - Computer annotated rather than human annotated

- Families of proteins
- Can search using regular expressions
 - Similar to unix commands using wildcards, etc.
 - E.g., [AC]-x-V-x(4)-{ED}
 - Interpreted as:
 - [Ala or Cys]-any-Val-any-any-any-any-{any but Glu or Asp}
- Families exhibit these patterns
 - So we can search over families
- 1574 documents about 1308 different patterns

<http://ca.expasy.org/prosite/>



PFAM

The University
of Georgia

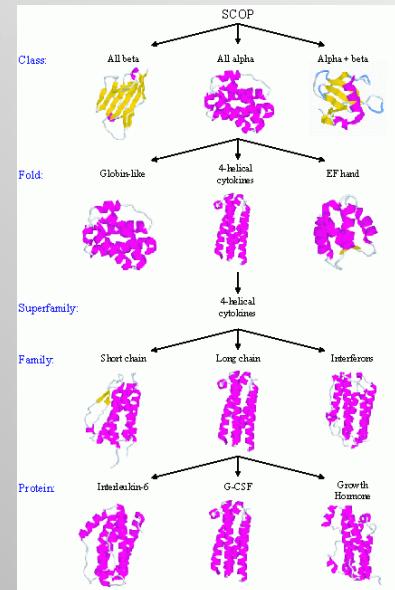
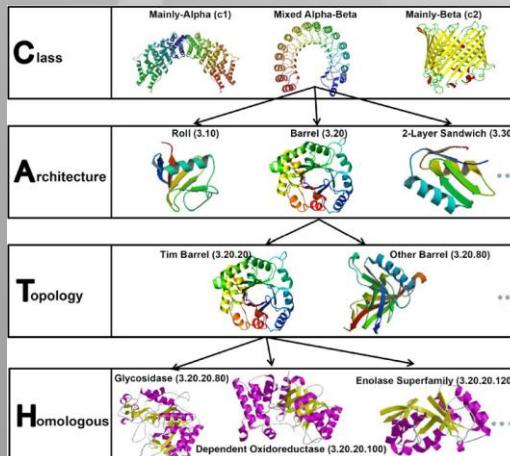
<http://pfam.sanger.ac.uk/>

- Maintained by the Sanger Centre (Cambridge)
- Protein families aligned using HMMs
 - Hidden Markov Models (see later lecture)
- Given a new sequence
 - Find families which the sequence might fit into
- Sequence Coverage
 - 11912 families
 - Split into Pfam-A (high quality) and Pfam-B (low quality)

The University SCOP and CATH of Georgia

<http://scop.mrc-lmb.cam.ac.uk/scop/>
<http://www.cathdb.info/>

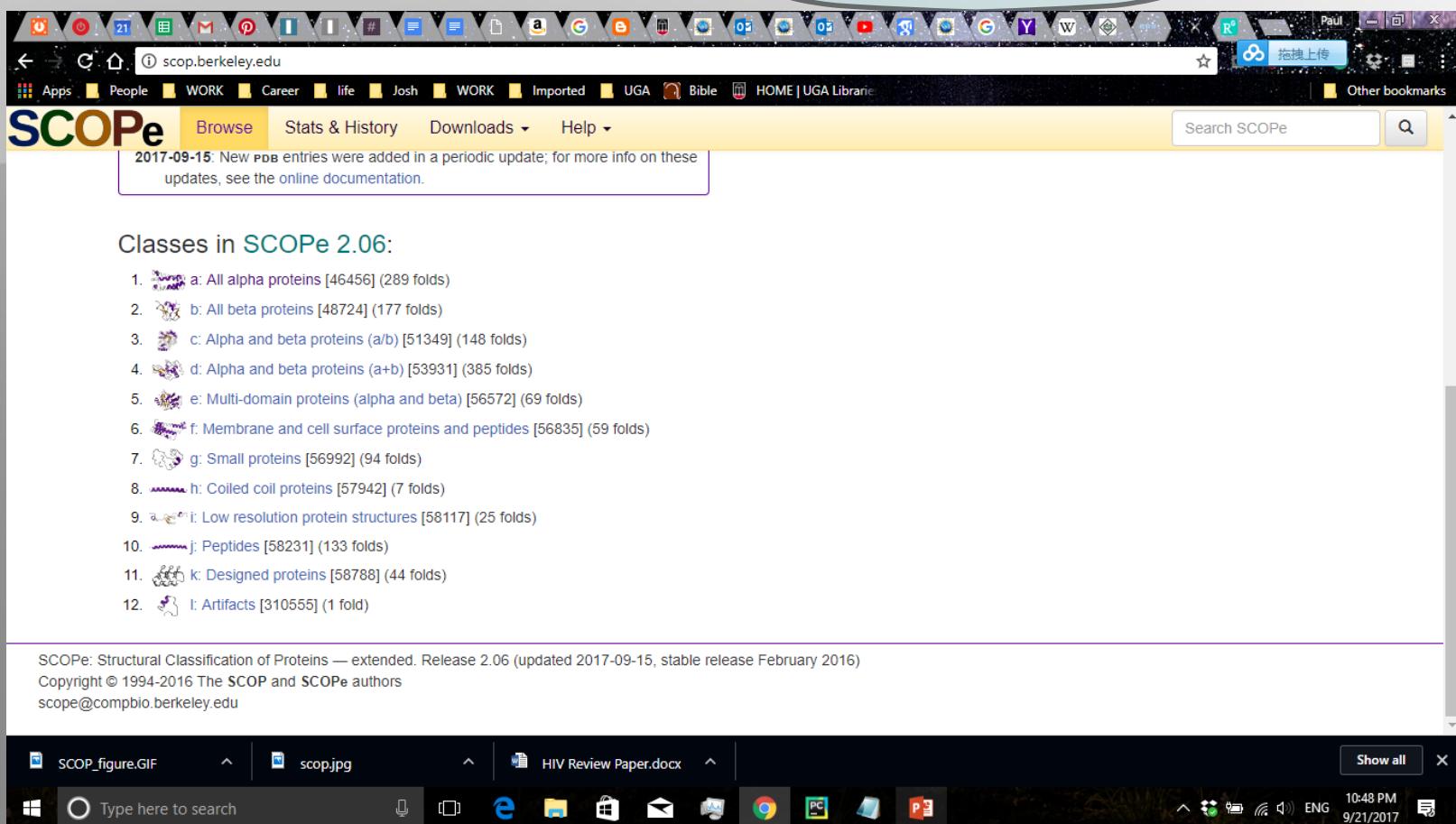
- SCOP
 - Structural Classification of Proteins
 - Hierarchically ordered and manually curated
 - 38221 PDB Entries
 - 110800 Domains
- CATH
 - Classification of protein domain structures
 - 124 folds
 - 226 Superfamily
 - 1148 Sequence family
 - 14473 Domain



The University SCOP and CATH of Georgia

- SCOP

<http://scop.mrc-lmb.cam.ac.uk/scop/>
<http://www.cathdb.info/>



Classes in SCOPe 2.06:

1. a: All alpha proteins [46456] (289 folds)
2. b: All beta proteins [48724] (177 folds)
3. c: Alpha and beta proteins (a/b) [51349] (148 folds)
4. d: Alpha and beta proteins (a+b) [53931] (385 folds)
5. e: Multi-domain proteins (alpha and beta) [56572] (69 folds)
6. f: Membrane and cell surface proteins and peptides [56835] (59 folds)
7. g: Small proteins [56992] (94 folds)
8. h: Coiled coil proteins [57942] (7 folds)
9. i: Low resolution protein structures [58117] (25 folds)
10. j: Peptides [58231] (133 folds)
11. k: Designed proteins [58788] (44 folds)
12. l: Artifacts [310555] (1 fold)

SCOPe: Structural Classification of Proteins — extended. Release 2.06 (updated 2017-09-15, stable release February 2016)
Copyright © 1994-2016 The SCOP and SCOPe authors
scope@compbio.berkeley.edu

PDB Format

- The PDB format consists of a collection of fixed format records that describe :
 - Atomic coordinates,
 - Chemical and biochemical features
 - Experimental details of the structure determination
 - Some structural features such as
 - Secondary structure assignments,
 - Hydrogen bonding
 - Biological assemblies
 - Active sites



PDB
PROTEIN DATA BANK

An Information Service of Tuesday Sep 11, 2006 at 2:49 PDT Version 21227 Structures | PDB Statistics

A Resource for Studying Biological Macromolecules

The PDB archive contains information about experimentally determined structures of proteins, nucleic acids, and macromolecular assemblies. As a member of the RCSB, the PDB stores and distributes PDB data according to agreed upon standards.

The RCSB PDB also provides a variety of tools and resources. Users can perform simple and advanced searches using annotations related to sequence, structure and function. These molecules are visualized, manipulated, and analyzed by users who range from students to specialist scientists.

Molecule of the Month
Enhancersome

Take a minute to ponder what form our year takes. What shape is your home, the color of your dress, the length of your fingernails, the perfect simulation of your blood and muscles. One may say your hand, spine or brain. How let your imagination travel instead, and think of the complex shapes and functions of your different cells, and the amazing processes they undergo each day. Macromolecules, the building blocks of life, are the physical embodiment of the information in the genes which encode it. Over 20,000-21,000 protein-encoding genes. One of like great power being placed together by scientists in the mechanism by which these genes, and the methods used to control their expression, specify all of these different elements of life.

Side Panels

- [Notified Proteins](#)
- [Chemical Family](#)
- [Sequence Search](#)
- [Protein Comparison](#)

Side Panels

- [News](#)
- [New PDB Proteins](#)
- [PDB Advances](#)
- [ASGT Resources](#)

Presente

- [CDS@PDB](#)

Close This Panel

[View PDB ITTB Advisory Information](#)

[View PDB-ICM-2004](#)

[Download PDB, PDBx, and PDBx-VR](#)

[View the services for model users](#)

[View the services for individual users](#)

[Read More...](#)

PDB Format

```

SHEET 3 L 3 THR H 212 VHL H 214 -1 O THR H 212 N HIS H 207
SHEET 1 M 3 THR H 158 TRP H 161 0
SHEET 2 M 3 TVR H 201 HIS H 207 -1 O ASN H 206 N THR H 158
SHEET 3 M 3 LVS H 217 VAL H 218 -1 O VAL H 218 N TYR H 201
SHEET 1 N 3 TRP C 28 LEU C 29 0
SHEET 2 N 3 VAL C 13 ALA C 18 -1 N VAL C 17 O LEU C 29
SHEET 3 N 3 LEU C 36 ALA C 38 -1 O LEU C 36 N HIS C 15
SHEET 1 O 5 TRP C 28 LEU C 29 0
SHEET 2 O 5 VAL C 13 ALA C 18 -1 N VAL C 17 O LEU C 29
SHEET 3 O 5 TVR C 151 ALA C 156 -1 O ILE C 154 N ALA C 14
SHEET 4 O 5 GLY C 54 GLN C 67 -1 N TYR C 59 O GLY C 153
SHEET 5 O 5 PRO C 113 LEU C 126 -1 O GLY C 122 N ILE C 58
SHEET 1 P 5 GLU C 42 ARG C 44 0
SHEET 2 P 5 GLN C 47 VAL C 49 -1 O VAL C 49 N GLU C 42
SHEET 3 P 5 ARG C 131 ILE C 136 -1 O LEU C 132 N LEU C 48
SHEET 4 P 5 LEU C 76 ILE C 83 -1 N THR C 79 O GLU C 135
SHEET 5 P 5 LYS C 90 LYS C 98 -1 O LEU C 94 N ILE C 80
SSBOND 1 CVS L 23 CVS L 88 1555 1555 2.05
SSBOND 2 CVS L 134 CVS L 194 1555 1555 2.03
SSBOND 3 CVS H 22 CVS H 98 1555 1555 2.05
SSBOND 4 CVS H 147 CVS H 203 1555 1555 2.03
SSBOND 5 CVS C 69 CVS C 101 1555 1555 2.06
CISPEP 1 SER L 7 PRO L 8 0 2.35
CISPEP 2 TRP L 94 PRO L 95 0 6.74
CISPEP 3 TVR L 140 PRO L 141 0 1.27
CISPEP 4 GLU H 1 VAL H 2 0 3.36
CISPEP 5 GLY H 140 GLY H 141 0 9.68
CISPEP 6 PHE H 153 PRO H 154 0 -3.35
CISPEP 7 GLU H 155 PRO H 156 0 0.91
CISPEP 8 CVS C 101 GLN C 102 0 -3.69
CRYST1 153.663 153.663 99.279 90.00 90.00 120.00 H 3 9
ORIGX1 1.000000 0.000000 0.000000 0.000000
ORIGX2 0.000000 1.000000 0.000000 0.000000
ORIGX3 0.000000 0.000000 1.000000 0.000000
SCALE1 0.006508 0.003757 0.000000 0.000000
SCALE2 0.000000 0.007514 0.000000 0.000000
SCALE3 0.000000 0.000000 0.010073 0.000000
ATOM 1 N ASP L 1 24.141 -2.203 -3.932 1.00 34.47 N
ATOM 2 CA ASP L 1 25.160 -2.458 -2.905 1.00 44.54 C
ATOM 3 C ASP I 1 25.272 -1.169 -2.301 1.00 43.71 C

```

how file was produced

HEADER	the 17 numbered atoms				coordinates		
REMARK Spartan ST exported M001	1	Cu	UNK	0001	0.000	0.000	0.000
	2	N	UNK	0001	-0.870	0.000	1.728
	3	H	UNK	0001	-1.884	0.000	1.654
	4	H	UNK	0001	-0.607	0.825	2.263
	5	H	UNK	0001	-0.607	-0.825	2.263
	6	N	UNK	0001	-1.728	0.000	-0.870
	7	H	UNK	0001	-1.654	0.000	-1.884
	8	H	UNK	0001	-2.263	0.825	-0.607
	9	H	UNK	0001	-2.263	-0.825	-0.607
	10	N	UNK	0001	0.870	0.000	-1.728
	11	H	UNK	0001	1.884	0.000	-1.654
	12	H	UNK	0001	0.607	0.825	-2.263
	13	H	UNK	0001	0.607	-0.825	-2.263
	14	N	UNK	0001	1.728	0.000	0.870
	15	H	UNK	0001	1.654	0.000	1.884
	16	H	UNK	0001	2.263	0.825	0.607
	17	H	UNK	0001	2.263	-0.825	0.607
CONECT	1	2	6	10	14		
CONECT	2	3	4	5	1		
CONECT	3	2					
CONECT	4	2					
CONECT	5	2					
CONECT	6	7	8	9	1		
CONECT	7	6					
CONECT	8	6					
CONECT	9	6					
CONECT	10	11	12	13	1		
CONECT	11	10					
CONECT	12	10					
CONECT	13	10					
CONECT	14	15	16	17	1		
CONECT	15	14					
CONECT	16	14					
CONECT	17	14					
END							

How the 17 atoms are connected

PDB Format

Atomic Coordinates: PDB Format

Element	Amino Acid	Chain name	Sequence Number	Coordinates-----					
				X	Y	Z	(etc.)		
ATOM	1	N	ASP	L	1	4.060	7.307	5.186	...
ATOM	2	CA	ASP	L	1	4.042	7.776	6.553	...
ATOM	3	C	ASP	L	1	2.668	8.426	6.644	...
ATOM	4	O	ASP	L	1	1.987	8.438	5.606	...
ATOM	5	CB	ASP	L	1	5.090	8.827	6.797	...
ATOM	6	CG	ASP	L	1	6.338	8.761	5.929	...
ATOM	7	OD1	ASP	L	1	6.576	9.758	5.241	...
ATOM	8	OD2	ASP	L	1	7.065	7.759	5.948	...
\\									
Element position within amino									

ATOM	1132	NH1	ARG	A	149	31.814	-31.597	16.995
ATOM	1133	NH2	ARG	A	149	32.203	-32.934	18.816
ATOM	1134	N	ASN	A	150	29.346	-24.359	18.812
ATOM	1135	CA	ASN	A	150	28.480	-23.190	18.933
ATOM	1136	C	ASN	A	150	28.606	-22.168	17.808
ATOM	1137	O	ASN	A	150	27.803	-21.276	17.678
ATOM	1138	CB	ASN	A	150	28.732	-22.524	20.282
ATOM	1139	CG	ASN	A	150	28.284	-23.389	21.447
ATOM	1140	OD1	ASN	A	150	27.205	-23.981	21.430
ATOM	1141	ND2	ASN	A	150	29.110	-23.463	22.466
ATOM	1142	N	LEU	A	151	29.629	-22.313	16.996
ATOM	1143	CA	LEU	A	151	29.868	-21.415	15.894
ATOM	1144	C	LEU	A	151	29.953	-22.205	14.597
ATOM	1145	O	LEU	A	151	30.149	-23.422	14.614
ATOM	1146	CB	LEU	A	151	31.208	-20.735	16.100
ATOM	1147	CG	LEU	A	151	31.436	-19.884	17.337
ATOM	1148	CD1	LEU	A	151	32.846	-19.333	17.256

Atom Number Atom Type Amino Acid Type Chain Residue Number X,Y,Z Coordinates

PDB Format

Atomic Coordinates: PDB Format

Element	Amino Acid	Chain name	Sequence Number	Coordinates-----					
				X	Y	Z	(etc.)		
ATOM	1	N	ASP	L	1	4.060	7.307	5.186	...
ATOM	2	CA	ASP	L	1	4.042	7.776	6.553	...
ATOM	3	C	ASP	L	1	2.668	8.426	6.644	...
ATOM	4	O	ASP	L	1	1.987	8.438	5.606	...
ATOM	5	CB	ASP	L	1	5.090	8.827	6.797	...
ATOM	6	CG	ASP	L	1	6.338	8.761	5.929	...
ATOM	7	OD1	ASP	L	1	6.576	9.758	5.241	...
ATOM	8	OD2	ASP	L	1	7.065	7.759	5.948	...
\\									
Element position within amino									

ATOM	1132	NH1	ARG	A	149	31.814	-31.597	16.995
ATOM	1133	NH2	ARG	A	149	32.203	-32.934	18.816
ATOM	1134	N	ASN	A	150	29.346	-24.359	18.812
ATOM	1135	CA	ASN	A	150	28.480	-23.190	18.933
ATOM	1136	C	ASN	A	150	28.606	-22.168	17.808
ATOM	1137	O	ASN	A	150	27.803	-21.276	17.678
ATOM	1138	CB	ASN	A	150	28.732	-22.524	20.282
ATOM	1139	CG	ASN	A	150	28.284	-23.389	21.447
ATOM	1140	OD1	ASN	A	150	27.205	-23.981	21.430
ATOM	1141	ND2	ASN	A	150	29.110	-23.463	22.466
ATOM	1142	N	LEU	A	151	29.629	-22.313	16.996
ATOM	1143	CA	LEU	A	151	29.868	-21.415	15.894
ATOM	1144	C	LEU	A	151	29.953	-22.205	14.597
ATOM	1145	O	LEU	A	151	30.149	-23.422	14.614
ATOM	1146	CB	LEU	A	151	31.208	-20.735	16.100
ATOM	1147	CG	LEU	A	151	31.436	-19.884	17.337
ATOM	1148	CD1	LEU	A	151	32.846	-19.333	17.256

Atom Number Atom Type Amino Acid Type Chain Residue Number X,Y,Z Coordinates

PDBsum

EMBL-EBI Enter Text Here Find Help | Feedback

Databases Tools Research Training Industry About Us Help Site Index

PDBsum

PDB code (4 chars) 1gb1 Find Example: "1kfv"

Contents
PDBsum contains 79,259 entries, including 1,620 superseded Last update: 8 October, 2011

Related databases

- EC-POB** Enzyme 3D structures organized by the EC numbering hierarchy.
- DrugPort** Structures of drugs and their target proteins in the PDB.
- SAS** Searches sequence against all PDB sequences and structurally annotates alignment.
- ProFunc** Prediction of protein function from 3D structure.
- Arch Schema** NEW Graphs of protein sequences having related Pfam domain architectures.

Text and sequence searches

Submit your PDB file for PDBsum analysis

Search by sequence

Perform FASTA search vs all sequences in the PDB to get a list of the closest matches.

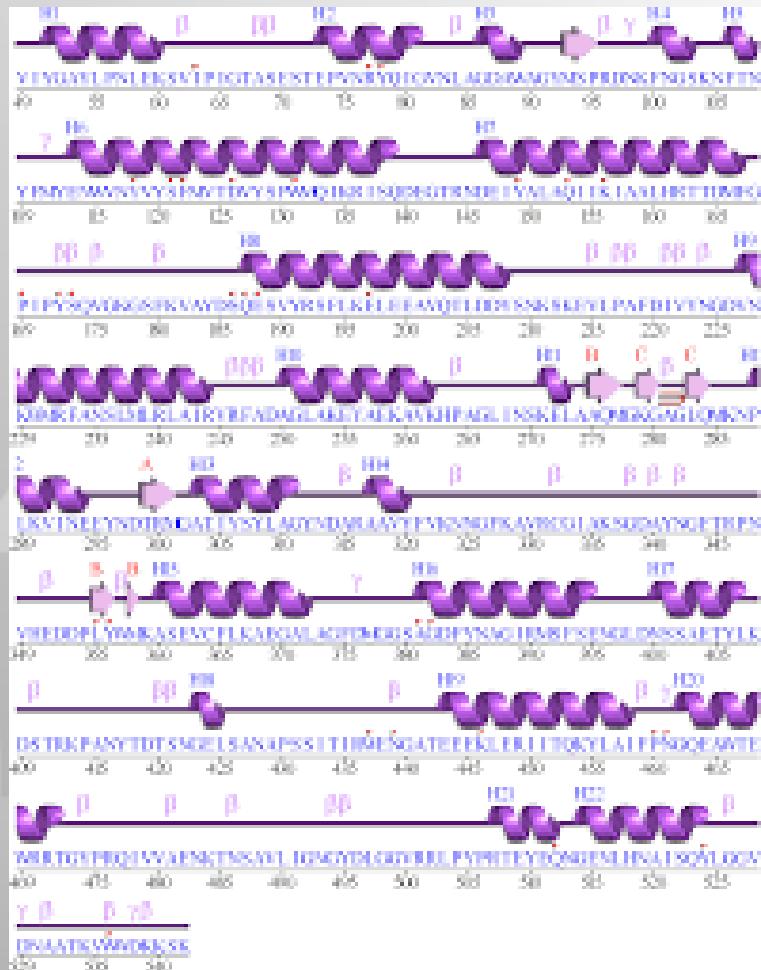
Notes

You can use the Generate option on the left to submit your own structure and get a password protected PDBsum analysis generated for it.

The development of PDBsum has been partly funded by the Wellcome Trust

1kv

Terms of Use EBI Funding Contact EBI © European Bioinformatics Institute 2011. EBI is an Outstation of the European Molecular Biology Laboratory.

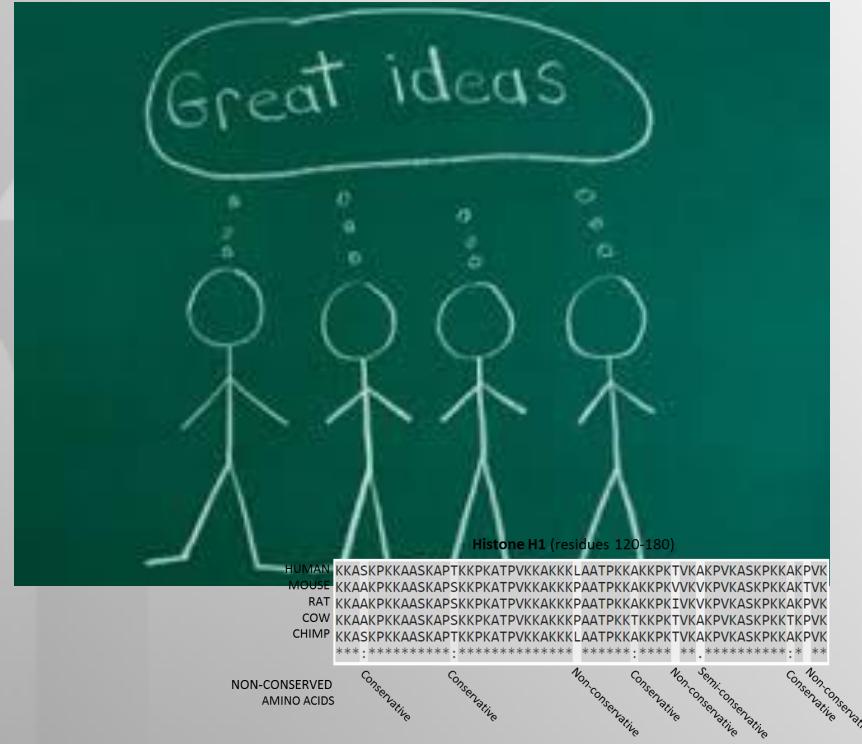


Python Time





Group Discussion





Sequence Databases

- Annotated sequence databases
 - SWISS-PROT, GenBank etc...
 - Usage: identifying function, retrieving information
- Low-annotation sequence databases
 - EST databases, high-throughput genome sequences
 - Usage: discovery of new genes

General Protein Databases

- SWISS-PROT
 - Manually curated
 - high-quality annotations, less data
- GenPept/TREMBL
 - Translated coding sequences from GenBank/EMBL
 - Few annotations, more up to date
- PIR
 - Phylogenetic-based annotations
- All 3 now combining efforts to form UniProt (<http://www.uniprot.org>)

Low-annotation Databases

- ESTs (Expressed Sequence Tags)
 - Low quality sequences generated by high - volume sequencing the 3' or 5' end of cDNAs
- High-throughput genome sequences
 - Produced by mass-sequencing of genomic DNA

Non-redundant Databases

- Sequence data only: cannot be browsed, can only be searched using a sequence
- Combine sequences from more than one database
- Examples:
 - NR Nucleic (genbank+EMBL+DDBJ+PDB DNA)
 - NR Protein (SWISS-PROT+TrEMBL+GenPept+PDB protein)

Sequence & Structure Databases

- PDB (Protein Databank)
 - Stores 3-dimensional atomic coordinates for biological molecules including protein and nucleic acids
 - Data obtained by X-ray crystallography, NMR, or computer modelling
 - <http://www.rcsb.org/pdb/>
- MMDB (Molecular Modelling database)
 - Over 28,000 3D macromolecular structures, including proteins and polynucleotides
 - (<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Structure>)
- SCOP (Structural Classification of Proteins)
 - Classification of proteins according to structural and evolutionary relationships

File Formats

- GenBank/GB, genbank flatfile format
- NBRF format
- EMBL, EMBL flatfile format
- Swissprot
- GCG, single sequence format of GCG software
- DNAStrider, for common Mac program
- Pearson/Fasta, a common format used by Fasta programs and others
- Phylip3.2, sequential format for Phylip programs
- Phylip, interleaved format for Phylip programs (v3.3, v3.4)
- Plain/Raw, sequence data only (no name, document, numbering)
- MSF multi sequence format used by GCG software
- PAUP's multiple sequence (NEXUS) format
- ASN.1 format used by NCBI

ID TRBG361 standard; mRNA; PLN; 1859 BP.

XX

AC X56734; S46826;

XX

SV X56734.1

XX

DT 12-SEP-1991 (Rel. 29, Created)

DT 15-MAR-1999 (Rel. 59, Last updated, Version 9)

XX

DE Trifolium repens mRNA for non-cyanogenic beta-glucosidase

XX

KW beta-glucosidase.

XX

OS Trifolium repens (white clover)

OC Eukaryota; Viridiplantae; Streptophyta; Embryophyta; Tracheophyta;

OC Spermatophyta; Magnoliophyta; eudicots; core eudicots; rosids;

OC eurosids I; Fabales; Fabaceae; Papilionoideae; Trifolieae; Trifolium.

XX

RN [5]

RP 1-1859

RX MEDLINE; 91322517.

RX PUBMED; 1907511.

RA Oxtoby E., Dunn M.A., Pancoro A., Hughes M.A.;

RT "Nucleotide and derived amino acid sequence of the cyanogenic

RT beta-glucosidase (linamarase) from white clover (Trifolium repens L.).";

RL Plant Mol. Biol. 17(2):209-219(1991).

XX

RN [6]

RP 1-1859

RA Hughes M.A.;

RT ;

RL Submitted (19-NOV-1990) to the EMBL/GenBank/DDBJ databases.

RL M.A. Hughes, UNIVERSITY OF NEWCASTLE UPON TYNE, MEDICAL SCHOOL, NEW CASTLE

RL UPON TYNE, NE2 4HH, UK

XX

DR GOA; P26204.

DR MENDEL; 11000; Trirp;1162;11000.

DR SWISS-PROT; P26204; BGLS_TRIRP.

XX

FH Key Location/Qualifiers

FH

FT source 1..1859

FT /db_xref="taxon:3899"

FT /mol_type="mRNA"

FT /organism="Trifolium repens"

FT /tissue_type="leaves"

FT /clone_lib="lambda gt10"

FT /clone="TRE361"

FT CDS 14..1495

FT /db_xref="GOA:P26204"

FT /db_xref="SWISS-PROT:P26204"

FT /note="non-cyanogenic"

FT /EC_number="3.2.1.21"

FT /product="beta-glucosidase"

FT /protein_id="CAA40058.1"

FT /translation="MDFIVAIHALFVISSFTITSTNAVEASTLLDIGNLSRSSFPFGFI

FT FGAGSSAYQFEGAVNEGGRGPSIWDTFTHKYPEKIRDGSNADITVDQYHRYKEDVGIMK

FT DQNMDSYRFSISWPRLPKGKLGGINHEGIKYNNLINELLANGIQPFVTLFHWDLPQ

FT VLEDEYGGFLNSGVINDFRDYTDLCFKEFGDRVRYWSTLNEPWVFSNSGYALGTNAPGR

FT CSASNVAKGPGDSGTGPYIVTHNQILAHAEAVHVYKTKYQAYQKGKIGITLVSNWLMP LD

FT DNSIPDIKAAERSLDFQFGLFMEQLTTGDYSKSMRRIVKNRLPKFSKFESSLVNGSFD F

FT IGINYYSSYYISNAPSHGNNAKPSYSTNPMTNISFEKHGIPLGPRASIWIYVYPYMF IQ

FT EDFEIFCYILKINITILQFSITENGMEFNNDATLPVEEALLNTYRIDYYRHLYYIRSA

FT IRAGSNVKGFYAWSFLDCNEWFAGFTVRFGLNFVD"

FT mRNA 1..1859

FT /evidence=EXPERIMENTAL

XX

SQ Sequence 1859 BP; 609 A; 314 C; 355 G; 581 T; 0 other;

aaacaaaacca aatatggatt ttattgtac catatttgc ctgtttgtta ttagctcatt 60

cacaattact tccacaaaatg cagttgaagc ttctactctt cttgacatag gtaaacctgag 120

tccggagcagt tttcctcggt gcttcatctt tggtgctgga tcttcagcat accaatttga 180

agtgccagta aacgaaggcg gtagaggacc aagtatttgg gataacctca cccataaaata 240

tccagaaaaaa ataaggggatg gaagcaatgc agacatcacg gttgaccaat atcaccgcta 300

caaggaagat gttgggattta tgaaggatca aaatatggat tcgtatagat tctcaatctc 360

ttggccaaga atactcccaa agggaaagtt gagcggaggc ataaatcacg aaggaa

Genbank Format

LOCUS SCU49845 **5028 bp** **DNA** **PLN**
21-JUN-1999
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.
REFERENCE 1 (bases 1 to 5028)
AUTHORS Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.
TITLE Cloning and sequence of REV7, a gene whose function
 is required for
 DNA damage-induced mutagenesis in *Saccharomyces*
cerevisiae
JOURNAL Yeast 10 (11), 1503-1509 (1994)
MEDLINE 95176709
PUBMED 7871890
REFERENCE 2 (bases 1 to 5028)
AUTHORS Roemer,T., Madden,K., Chang,J. and Snyder,M.
TITLE Selection of axial growth sites in yeast requires
 Axl2p, a novel
 plasma membrane glycoprotein
JOURNAL Genes Dev. 10 (7), 777-793 (1996)
MEDLINE 96194260
PUBMED 8846915
REFERENCE 3 (bases 1 to 5028)
AUTHORS Roemer,T.
TITLE Direct Submission
JOURNAL Submitted (22-FEB-1996) Terry Roemer, Biology, Yale
 University, New
 Haven, CT, USA

gene 687_3158
CDS 687_3158
 /gene="AXL2"
 /note="plasma membrane glycoprotein"
 /product="Ax12p"
 /protein_id="AAA98666.1"
 /db_xref="GI:1293615"

 /translation="MTQLQISLLLTTATISLLHLVVATPYEAYPIGKQYPPVARVNESF
 TFQISNDTYKSSVDKTAQITYNCFDLPSWLSFDSSRTFSGEPSSDILSDANTTLYFN
 VILEGTDSDASTSLNNNTYQFVVTNRPSISLSSDFNLALLKNYGYTNGKNALKLDPNE
 VFNVTFDRSMFTNEESIVSYYGRSQLYNAPLPNWLFDFSGELKFTGTAPVINSIAPE
 TSYSFVIIATDIEGFSAVEVEFELVIGAHQLTTSIQNSLIINVTDGNVSYDLPLNYV
 YLDDDPPISSDKLGSINLLADPDWALDNATISGSVPDELLGKNSNPANFSVSIYDTYG
 DVIYFNFEVVSTTDLFAISSLPNINATRGEWFSYYFLPSQFTDYVNTNVSLFTNSSQ
 DHDWVKFQSSNLTLAGEVPKNFDKLSLGLKANQGSQSQELYFNIIGMDSKITHSNHSA
 NATSTRSSHSTSTSSYTAKISSTSAAATSSAPAALPAANKTSSHNNKKAVAIA
 CGVAIPLGVILVALICFLIFWRRRRENPDDENLPHAIISGPDLNNPANKPNQENATPLN
 NPFDDDASSYDDTSIARRLAALNTLKLDNHSATESDISSVDEKRDSLSGMNTYNDQFQ
 SQSKEELLAKPPVQPPESPFFDPQRSSSVYMDSEPAVNKSWRYTGNLSPVSDIVRDS
 YGSQKTVDEKLFLEAPEKEKRTSRDVMTSSLDPWNSNISPSPVRKSVTPSPYNTK
 HRNRHLQNIQDSQSGKNGITPTTMSTSSSDFVPVKDGENFCWVHSMEPDRRPSKKRL
 VDFSNKSNVNVGQVKDIHGRIPEML
BASE COUNT 1510 **a** 1074 **c** 835 **g** 1609 **t**
ORIGIN
 1 gatcctccat atacaacggt atctccacct caggtttaga tctcaacaac
 ggaaccattg
 61 ccgacatgag acagtttagt atcgtcgaga gttacaagct aaaacgagca
 gtagtcagct
 121 ctgcacatcta agccgctgaa gttctactaa gggtggataa catcatccgt
 gcaagaccaa
 181 gaaccggccaa tagacaacat atgtaacata tttaggatat acctcgaaaa
 taataaaaccg

Swiss-Prot: Q823P0

NiceProt - a user-friendly view of this Swiss-Prot entry

ID DNA1_CHLCV STANDARD; PRT: 450 AA.
AC Q823P0;
DT 10-OCT-2003 (Rel. 42, Created)
DT 10-OCT-2003 (Rel. 42, Last sequence update)
DT 25-OCT-2004 (Rel. 45, Last annotation update)
DE Chromosomal replication initiator protein dnaA 1.
GN Name=dnaA1; Synonyms=dnaA-1; OrderedLocusNames=CCA00368;
OS Chlamydophila caviae.
OC Bacteria; Chlamydiae; Chlamydiales; Chlamydiaceae; Chlamydophila.
OX NCBI_TaxID=83557;
RN [1]
RP SEQUENCE FROM N.A.
RC STRAIN=GPIC;
RX MEDLINE=22569155; PubMed=12682364 [NCBI, ExPASY, EBI, Israel, Japan]; DOI=10.1093/nar/gkg100
RA Read T.D., Myers G.S.A., Brunham R.C., Nelson W.C., Paulsen I.T.,
RA Heidelberg J.F., Holtzapfle E.K., Khouri H.M., Federova N.B.,
RA Carty H.A., Umayam L.A., Haft D.H., Peterson J.D., Beanan M.J.,
RA White O., Salzberg S.L., Hsia R.-C., McClarty G., Rank R.G.,
RA Bavoil P.M., Fraser C.M.;
RT "Genome sequence of Chlamydophila caviae (Chlamydia psittaci GPIC): examining the role of niche-specific genes in the evolution of the Chlamydiaceae.";
RL Nucleic Acids Res. 31:2134-2147(2003).
CC --!- FUNCTION: Plays an important role in the initiation and regulation of chromosomal replication. Binds to the origin of replication; it binds specifically double-stranded DNA at a 9 bp consensus (dnaA box): 5'-TTATC(C/A)A(C/A)A-3'. DnaA binds to ATP and to acidic phospholipids (By similarity).
CC --!- SIMILARITY: Belongs to the dnaA family.
CC
CC This SWISS-PROT entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL outstation - the European Bioinformatics Institute. There are no restrictions on its use by non-profit institutions as long as its content is in no way modified and this statement is not removed. Usage by and for commercial entities requires a license agreement (See <http://www.isb-sib.ch/announce/>

CC or send an email to license@isb-sib.ch).

CC

DR EMBL; AEO16995; AAP05115.1; -. [EMBL / GenBank / DDBJ] [CoCodingSequence]
DR HSSP; P03004; 1J1V. [HSSP ENTRY / SWISS-3DIMAGE / PDB]
DR TIGR; CCA00368; -.
DR HAMAP; MF 00377; -. 1.
DR InterPro; IPR003593; AAA ATPase.
DR InterPro; IPR001957; Bac DnaA.
DR InterPro; IPR010921; Trp repress rep.
DR InterPro; Graphical view of domain structure.
DR Pfam; PF00308; Bac DnaA; 1.
DR Pfam; Graphical view of domain structure.
DR PRINTS; PRO0051; DNAA.
DR SMART; SMO0382; AAA; 1.
DR TIGRFAMS; TIGR00362; DnaA; 1.
DR PROSITE; PS01008; DNAA; 1.
DR ProDom [Domain structure / List of seq. sharing at least 1 domain]
DR HOBACGEN [Family / Alignment / Tree]
DR BLOCKS; Q823P0.
DR ProtoNet; Q823P0.
DR ProtoMap; Q823P0.
DR PRESAGE; Q823P0.
DR DIP; Q823P0.
DR ModBase; Q823P0.
DR SMR; Q823P0.
DR SWISS-2DPAGE; GET REGION ON 2D PAGE.
KU ATP-binding; Complete proteome; DNA replication; DNA-binding.
FT NP_BIND 156 163 ATP (Potential).
SQ SEQUENCE 450 AA; 51099 MW; CF440A7B300210D8 CRC64;
MLTCSDCSTW EQFVNYYVKTR CSKTAFENUI SPIQIIEETQ EKIRLEVFPNI FVQNYLLDNY
KQDLCFSVPL DAQGEPALEF VVAEIKKAPA QPIAPREPQE SPAETFEESK DFEKLKNAAAY
RFDMNIEGPS NQFVKSAAVG IAGRPGRSYN PLFIHGGVGL GKTHLLHAVG HYVREHHKNL
RVHCITTEAF INDLVQHLRL KSIDMKMFY RSLDLLLVDD IQFLQNQRNF EEEFCNTFET
LINLNKQIVI TSDKPPGQLK LSERIIARMF WGLVAHVGIP DLETRVAILQ HKAFQKGLHI
PNEIAFYIAD HIYGWVRQLE GAINKLTAYC RLFGKTLTES IVRDTLRELFSRSPSKQKVSV
ESILKSVATV FQVKLQLDLKG TSRSKELVLA RQVAMYLAKT LITDSLVAIG SAFGKTHSTV
LYACKTIEQK IEKDETLLTRQ ISLCKNHIVG

//

Specialized Sequence Databases

- Focus on a specific type of sequences
- Sequences are often modified or specially annotated
- Usage depends on the database
- Examples:
 - Ribosomal RNA databases
 - Immunology databases

Protein domain databases

- Pfam (<http://www.sanger.ac.uk/Software/Pfam/>)
 - Collection of multiple sequence alignments and hidden Markov models covering many common protein domains and families
- SMART (**a Simple Modular Architecture Research Tool**)
 - Identification and annotation of genetically mobile domains and the analysis of domain architectures
 - (http://smart.embl-heidelberg.de/help/smart_about.shtml)
- CDD (<http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>)
 - Combines SMART and Pfam databases
 - Easier and quicker search