

Protein Bioinformatics

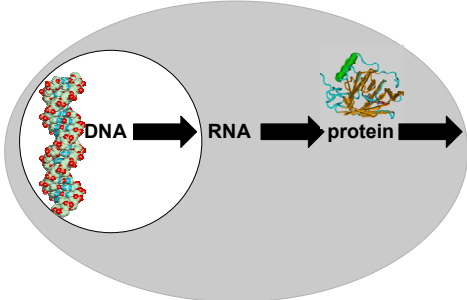
Part I: Access to information

260.655
April 6, 2006
Jonathan Pevsner, Ph.D.
pevsner@kennedykrieger.org

Outline

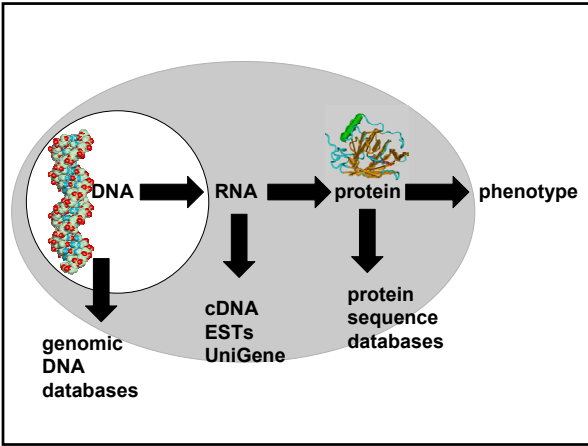
- [1] Proteins at NCBI
 - RefSeq accession numbers
 - Cn3D to visualize structures
- [2] The Protein Data Bank (PDB)
- [3] UniProt
- [4] ExpASy (Expert Protein Analysis System)
 - DeepView, the Swiss-Pdb Viewer.

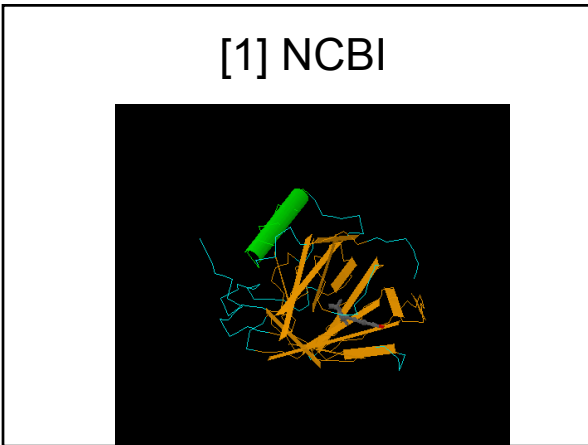
Central dogma of molecular biology

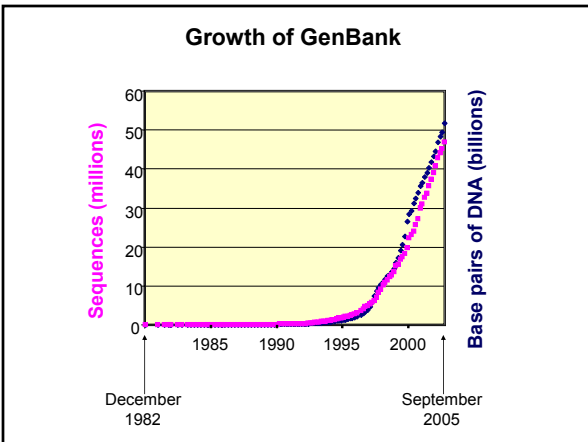


genome → transcriptome → proteome

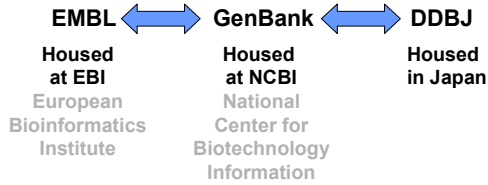
Central dogma of bioinformatics and genomics

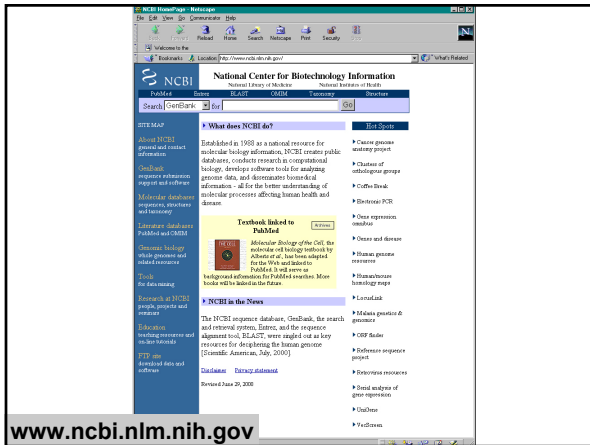






There are three major public DNA databases





Accession numbers are labels for sequences

NCBI includes databases (such as GenBank) that contain information on DNA, RNA, or protein sequences. You may want to acquire information beginning with a query such as the name of a protein of interest, or the raw nucleotides comprising a DNA sequence of interest.

DNA sequences and other molecular data are tagged with accession numbers that are used to identify a sequence or other record relevant to molecular data.

What is an accession number?

An accession number is label that used to identify a sequence. It is a string of letters and/or numbers that corresponds to a molecular sequence.

Examples (all for retinol-binding protein, RBP4):

X02775	GenBank genomic DNA sequence	DNA
NT_030059	Genomic contig	
Rs7079946	dbSNP (single nucleotide polymorphism)	
N91759.1	An expressed sequence tag (1 of 328)	RNA
NM_006744	RefSeq DNA sequence (from a transcript)	
NP_007635	RefSeq protein	protein
AAC02945	GenBank protein	
Q28369	SwissProt protein	
1KT7	Protein Data Bank structure record	

Accessing protein sequences via Entrez

Entrez Gene with RefSeq

Entrez Gene is a great starting point: it collects key information on each gene/protein from major databases. It covers all major organisms.

RefSeq provides a curated, optimal accession number for each DNA (NM_006744) or protein (NP_007635)

Example #1. Sean mentioned silk fibroin. How do you find its sequence?

From the NCBI home page, type "silk fibroin" and hit "Go"

The screenshot shows the NCBI homepage with a search bar containing the text "silk fibroin". The search results section is partially visible, showing "What does NCBI do?" and "Hot Spots". The "Hot Spots" section includes links to "Assembly Archive", "Clusters of orthologous groups", "Coffee Break, Genes & Disease, NCBI Handbook", "Electronic PCR", "Entrez Home", "Entrez Tools", "Gene expression omnibus (GEO)", and "Human genome resources".

You can try scrolling through the RefSeq list, or apply "Limits"

The screenshot shows the NCBI Protein search results for the query 'sperm whale'. The search results are displayed in a table with columns for 'Limits', 'Protein', 'Reports', and 'BLINK, Conserved Domains, Links'. The results are sorted by relevance and show the first 7 items out of 47. The first item is 'NP_976312' for the protein 'myoglobin [Bala capensis]'. The second item is 'NP_976311' for 'myoglobin [Bala capensis]'. The third item is 'NP_003359' for 'myoglobin [Bala capensis]'. The fourth item is 'NP_038621' for 'myoglobin [Bala macrorhynchus]'. The fifth item is 'NP_067299' for 'myoglobin [Bala macrorhynchus]'. The sixth item is 'NP_776306' for 'myoglobin [Bala tairua]'. The seventh item is 'NP_999401' for 'myoglobin [Bala tairua]'. The search bar at the top shows the query 'sperm whale' and the search button is labeled 'Go'.



As another approach, click "TaxBrowser"...

The screenshot shows the NCBI Taxonomy Browser search page. The search bar contains the text 'sperm whale' and the search button is labeled 'Go'. The page title is 'National Center for Biotechnology Information' and 'National Institutes of Health'. The search results are displayed in a table with columns for 'Taxonomy ID', 'Common name', and 'Scientific name'. The results are sorted by relevance and show the first 4 items out of 47. The first item is 'Kozia simus [genbank common name: dwarf sperm whale]'. The second item is 'Physeter catodon [genbank common name: sperm whale]'. The third item is 'Kozia brevicauda [genbank common name: pygmy sperm whale]'. The fourth item is 'Dwarf sperm whale zammaherperivivus'.

Enter the name of the organism you are interested in...

The screenshot shows the NCBI Taxonomy Browser search page. The search bar contains the text 'sperm whale' and the search button is labeled 'Go'. The page title is 'National Center for Biotechnology Information' and 'National Institutes of Health'. The search results are displayed in a table with columns for 'Taxonomy ID', 'Common name', and 'Scientific name'. The results are sorted by relevance and show the first 4 items out of 47. The first item is 'Kozia simus [genbank common name: dwarf sperm whale]'. The second item is 'Physeter catodon [genbank common name: sperm whale]'. The third item is 'Kozia brevicauda [genbank common name: pygmy sperm whale]'. The fourth item is 'Dwarf sperm whale zammaherperivivus'.

Follow a link of interest...

sperm whale

- [Kozia simus \[genbank common name: dwarf sperm whale\]](#)
- [Physeter catodon \[genbank common name: sperm whale\]](#)
- [Kozia brevicauda \[genbank common name: pygmy sperm whale\]](#)
- [Dwarf sperm whale zammaherperivivus](#)

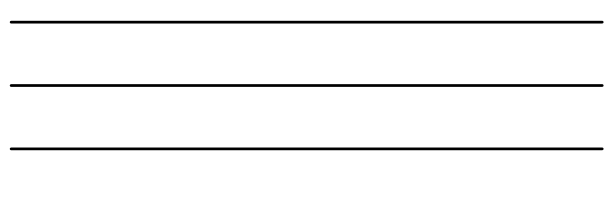


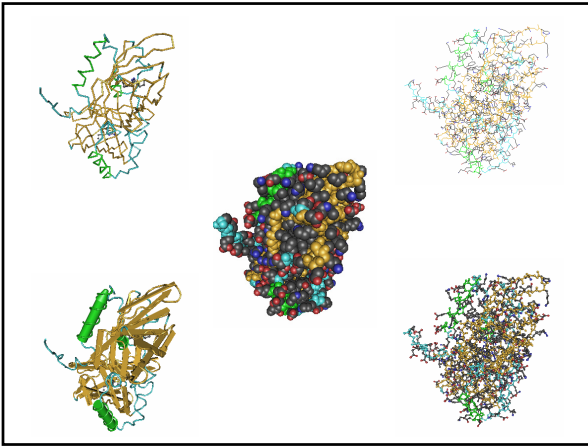
Now click protein...

The screenshot shows the NCBI Protein search results for the query 'Physeter catodon'. The search results are displayed in a table with columns for 'Limits', 'Protein', 'Reports', and 'BLINK, Conserved Domains, Links'. The results are sorted by relevance and show the first 1 item out of 1. The first item is 'NP_067299' for the protein 'myoglobin [Bala macrorhynchus]'. The search bar at the top shows the query 'Physeter catodon' and the search button is labeled 'Go'. The page title is 'National Center for Biotechnology Information' and 'National Institutes of Health'.

You now can view all sperm whale proteins...

The screenshot shows the NCBI Protein search results for the query 'sperm whale'. The search results are displayed in a table with columns for 'Limits', 'Protein', 'Reports', and 'BLINK, Conserved Domains, Links'. The results are sorted by relevance and show the first 1 item out of 1. The first item is 'NP_067299' for the protein 'myoglobin [Bala macrorhynchus]'. The search bar at the top shows the query 'sperm whale' and the search button is labeled 'Go'. The page title is 'National Center for Biotechnology Information' and 'National Institutes of Health'.





Overlay two or more structures with VAST at NCBI

The graphics below indicate the individual chains, 3D domains and ligands identified, if present, in the MMDB structure. You may view them in [Cn3D](#) by clicking "View 3D Structure" above. You may also click each icon to get more information.

Click "Chain"...

1390 neighbors found. 57 representatives from the Medium redundancy subset displayed.

Click one or more boxes then "View 3D Structure" ...

[2] PDB



The Protein Data Bank (PDB)

- PDB is the principal repository for protein structures
- Established in 1971
- Accessed at <http://www.rcsb.org/pdb> or simply <http://www.pdb.org>
- Currently contains over 35,000 structure entities

Updated 3/06

WELCOME TO THE RCSB PDB

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

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

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

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

A narrated tutorial illustrates how to search, navigate, browse, generate reports, and visualize structures using this new site. (The [summary homepage](#) has your [tutorial](#).)

Comments? info@rcsb.org

Molecule of the Month: Tissue Factor

Blood performs many essential jobs in your body: it transports oxygen and nutrients, it carries your cells from infection, and it carries hormones and other messages from place to place in your body. But one blood cell is a kind of what is pumped under pressure, we must protect ourselves from leaks. Fortunately, the blood has a built-in repair method that quickly seals up breaks in the blood circulatory system as soon as they happen. You see these repairs in action whenever you cut yourself: the blood clots and forms a gray *clot*, which then dries into a scab that seals and protects the cut until it can heal.

More ...

Previous Features

NEWS

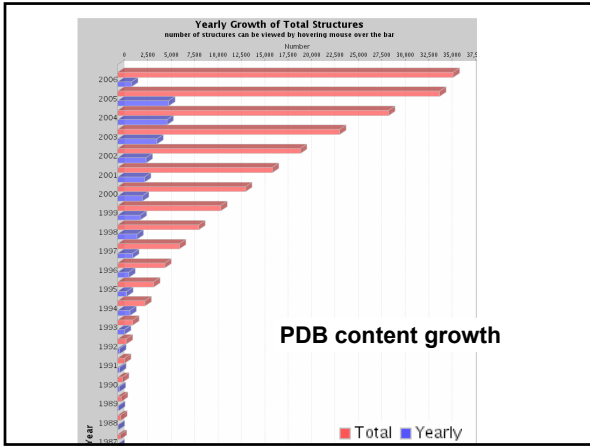
- Complete News
- Newsletter
- Discussion Forum

28-Mar-2006
Art of Science Exhibit and "PDB In-a-Clave" at Virginia Tech Structural Biology Symposium

Images from the RCSB PDB's "Art of Science" exhibit will be on display at Virginia Polytechnic Institute and State University, as part of their Structural Biology Symposium (March 24 - April 1, Blacksburg, VA).

Full Story ...

21-Mar-2006
RCSB PDB Exhibit Booth and Presentations at Experimental Biology



PDB holdings (September, 2005)

29,876	proteins, peptides
1,338	protein/nucl. complexes
1,500	nucleic acids
13	carbohydrates
32,727	total

Search for keyword DNAC yields mouse zinc finger binding proteins

PDB
PROTEIN DATA BANK

An Information Portal to Biological Macromolecular Structure
As of Tuesday Mar 28, 2006 there are 36823 Structures

Search Results for: DNAC

Structure ID	Resolution	Exp. Method	Author
1A1E	1.90 Å	X-Ray Diffraction	Elnaf-Eriksson, M., Benson, T.E., Pabo, C.O.
1A1G	1.90 Å	X-Ray Diffraction	Elnaf-Eriksson, M., Benson, T.E., Pabo, C.O.
1A1H	1.90 Å	X-Ray Diffraction	Elnaf-Eriksson, M., Benson, T.E., Pabo, C.O.

gateways to access PDB files

Swiss-Prot, NCBI, EMBL



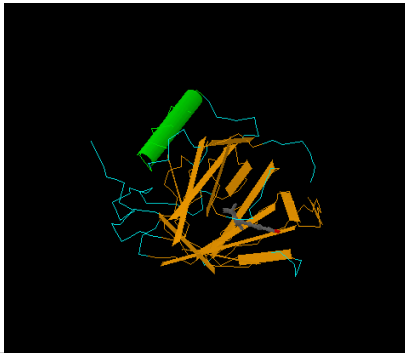
Protein Data Bank



CATH, Dali, SCOP, FSSP

databases that interpret PDB files

[3] UniProt



UniProt (Universal Protein Resource) at www.uniprot.org

UniProt combines information in Swiss-Prot, TrEMBL, and PIR. UniProt is comprised of three components

- The **UniProt Knowledgebase (UniProtKB)** is the central access point for extensive curated protein information.
- The **UniProt Reference Clusters (UniRef)** databases combine closely related sequences into a single record.
- The **UniProt Archive (UniParc)** is a comprehensive repository, reflecting the history of all protein sequences.

www.uniprot.org

UniProt (Universal Protein Resource) is the world's most comprehensive catalog of information on proteins. It is a central repository of protein sequence and function created by joining the information contained in Swiss-Prot, TrEMBL, and PIR.

UniProt is comprised of three components, each optimized for different uses. The **UniProt Knowledgebase (UniProtKB)** is the central access point for extensive curated protein information, including function, classification, and cross-reference. The **UniProt Reference Clusters (UniRef)** databases combine closely related sequences into a single record to speed searches. The **UniProt Archive (UniParc)** is a comprehensive repository, reflecting the history of all protein sequences.

The sequences and information in UniProt are accessible via **text search**, **BLAST similarity search**, and **FTP**.

Example: search for *E. coli* DnaC at NCBI

Search: Protein For [tax:562[Organism:ex] dnaC] Go Clear Save Search

Display: Summary Show 20 Sort by Relevance Send to

All: 67 Bacteria: 67 RefSeq: 34 x Related Structures: 55 X

Items 1 - 20 of 34

- 1: [NP_418781](#) Reports
DNA replication protein DnaC [Escherichia coli K12]
gi16132182|ref|NP_418781.1|16132182
- 2: [NP_757286](#) Reports
DNA replication protein DnaC [Escherichia coli CFTV73]
gi26251246|ref|NP_757286.1|26251246
- 3: [NP_290977](#) Reports
DNA replication protein DnaC [Escherichia coli O157:H7 EDL933]
gi15804935|ref|NP_290977.1|15804935

Approach: NCBI → TaxBrowser → *E. coli* → proteins → dnaC → RefSeq → three entries shown here.

Example: search for *E. coli* DnaC at UniProt

Text Search Result (UniProtKB)

Search: Any field For [tax:562] - Add input box

202 entries found. Page 1 of 1. 10 per page. Display: Summary Show 20 Sort by Relevance Send to

Accession	Protein Name	Length	Organism Name	Taxonomic Group	UniProt ID	Matched Fields
P31242 / P31242.1	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.2	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.3	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.4	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.5	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.6	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.7	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.8	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.9	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.10	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.11	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.12	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.13	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.14	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.15	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.16	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.17	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.18	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.19	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields
P31242 / P31242.20	Replicative DNA helicase	459	<i>Escherichia coli</i> str. O157:H7	Bacteria	NP_418781	Matched Fields

202 entries found

ExPASy to access protein and DNA sequences

ExPASy sequence retrieval system
(ExPASy = Expert Protein Analysis System)

Visit <http://www.expasy.ch/>

Site Map Search ExpAsy Contact us

Search | Swiss-Prot/EMBL for | Go | Clear

ExPASy Proteomics Server

The ExPASy (Expert Protein Analysis System) proteomics server of the Swiss Institute of Bioinformatics (SB) is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE (Disclaimer / References)

[Announcements](#) | [Job opening](#) | [Mirror Sites](#)

Databases	Tools and software packages
<ul style="list-style-type: none">• Swiss-Prot and TrEMBL - Protein knowledgebase• PROSITE - Protein families and domains• SWISS-2D-PAGE - Two-dimensional polyacrylamide gel electrophoresis• ENZYME - Enzyme nomenclature• SWISS-MODEL Repository - Automatically generated protein models• GermOnLine - Knowledgebase on germ cell differentiation• Ashbya Genome Database• Links to many other molecular biology databases	<ul style="list-style-type: none">• Proteomics and sequence analysis tools<ul style="list-style-type: none">◦ Proteomics Explorer (POE), Pipistem (MSMS), Phenyx (MSMD), Finkid, Pepsidatam, ...]◦ DNA → Protein (translat)◦ Similarity searches (BLAST)◦ Pattern and profile searches (ScanProsite)◦ Posttranslational modification and topology prediction◦ Primary structure analysis (Phospho, pMIR, MedCite)◦ Secondary and tertiary structure prediction (SWISS-MODEL, Swiss-PairView)◦ Image recognition (T-COFFEE, SM)◦ Phylogenetic analysis◦ Biological text analysis• ImageMaster / Milano - Software for 2-D PAGE analysis• MSight - Mass Spectrometry imager• Make2D-DB II - A package to build a web proteomics database• Roche Applied Science's Biochemical Pathways

ExPASy to access protein and DNA sequences

When you search the ExPASy database, you are now querying the UniProt Knowledgebase.

► UniProtKB/Swiss-Prot; a curated protein sequence database which strives to provide a high level of annotation (such as the description of the function of a protein, its domains structure, post-translational modifications, variants, etc.), a minimal level of redundancy and high level of integration with other databases.

UniProtKB/Swiss-Prot Release 49.3 of 21-Mar-2006:
212,425 entries (More statistics)

ExpASY to access protein and DNA sequences

When you search the ExpASY database, you are now querying the UniProt Knowledgebase.

► UniProtKB/TrEMBL; a computer-annotated supplement of Swiss-Prot that contains all the translations of EMBL nucleotide sequence entries not yet integrated in Swiss-Prot.

UniProtKB/TrEMBL Release 32.3 of 21-Mar-2006:
2,666,963 entries

Example: find human tyrosinase at ExpASY



Top Page Query Form Query Manager View Manager Databases Tick

Select one or more databases and continue

Continue Reset

Sequence SWISS_PROT TrEMBL

From the ExpASY home,
click Swiss-Prot → SRS
→ Start → Continue

Example: find human tyrosinase at ExpASY

Search: [SWISS_PROT](#) [TrEMBL](#)

Combine searches with: AND OR Append wildcard "*" to words

Info	AllText	tyrosinase
Info	Organism	human
Info	AllText	
Info	AllText	

Include fields in output: ID, AccNumber, Date, SubmissionDate, Description, GeneName, Keywords

Entry List in chunks of: 100

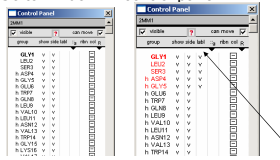
Sequence Format: *default*

Use view: ShortDescription

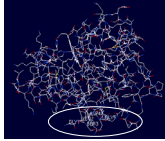
Retrieve set of: entry

Display in: list table

► Go to window → control panel.



- Shift/click to select the first five amino acid residues of myoglobin. They should appear red.
- Click "label" (i.e. label)(see arrow, above right). Those five residues now have a "v".
- Inspect the display panel; those five residues are labeled.



Download and practice using DeepView!
Try using myoglobin.

The ExpASy download site includes a helpful
web-based tutorial
<http://www.usm.maine.edu/~rhodes/SPVTut/>
