

BIOCH300D Bioinformatics Project 2012

Dr. Silvia Vidal

McGill Life Sciences Complex room 367

Phone: 398-2362

E-mail: silvia.vidal@mcgill.ca

<http://teodorolab.mcgill.ca/300D/>

OUTLINE INTRODUCTION

1- **WHAT IS BIOINFORMATICS?** definition, applications, goal

2- **WHERE DOES IT OCCUR?** Some important institutions

3- **HOW DOES IT HAPPEN?**

3.1 DATABASES NCBI databases

3.2 TOOLS for protein sequence analysis

3.2.1 BLAST

3.2.2 CLUSTAL W

3.3.3 3-D VIEWER

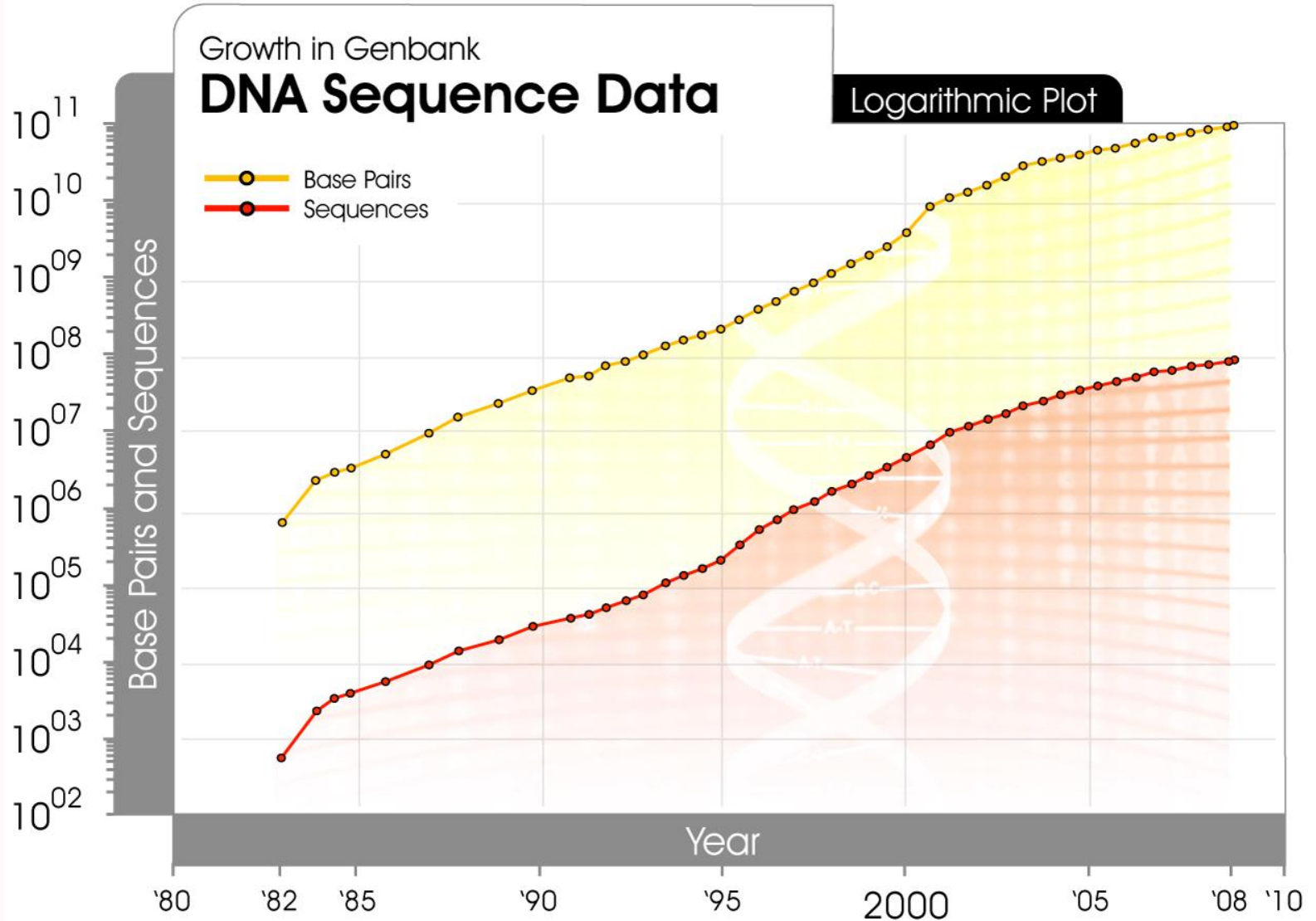
EXERCISE



What is Bioinformatics?

Scientific discipline resulting
from the merge of biology and
information sciences

Access and Storage of information is important



Access and Storage of information is important

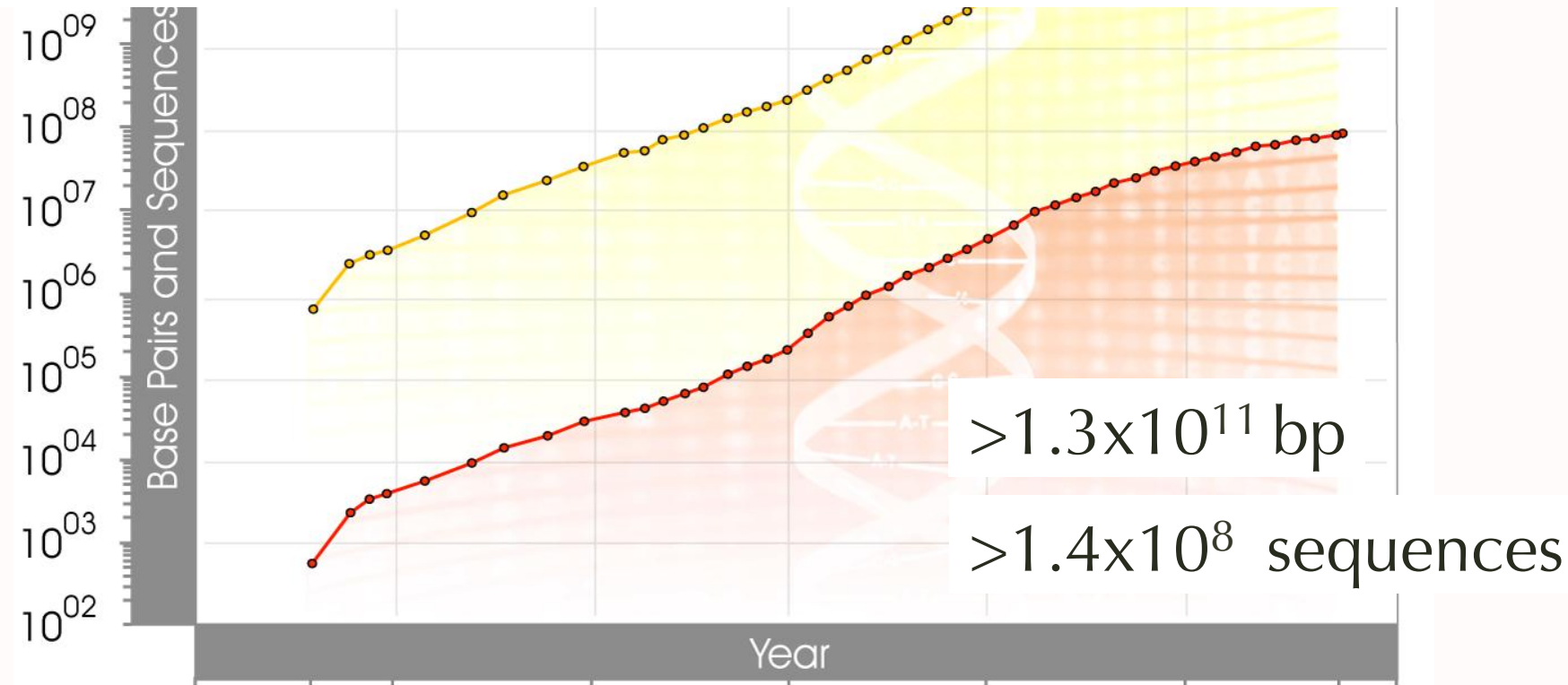
GBREL.TXT

Genetic Sequence Data Bank
February 15 2012

NCBI-GenBank Flat File Release 188.0

Distribution Release Notes

149819246 loci, 137384889783 bases, from 149819246 reported sequences



$>1.3 \times 10^{11}$ bp
 $>1.4 \times 10^8$ sequences

Bioinformatics comprises two main activities:

1- Computerized annotation of genomic
and biological information and data:

DATABASES

2- Transformation and manipulation of
the data: **SOFTWARE TOOLS**

The general goal of bioinformatics is:

- 1- Providing testable predictions from data and/or their manipulation
- 2- Needs to support data with statistics and probability calculations

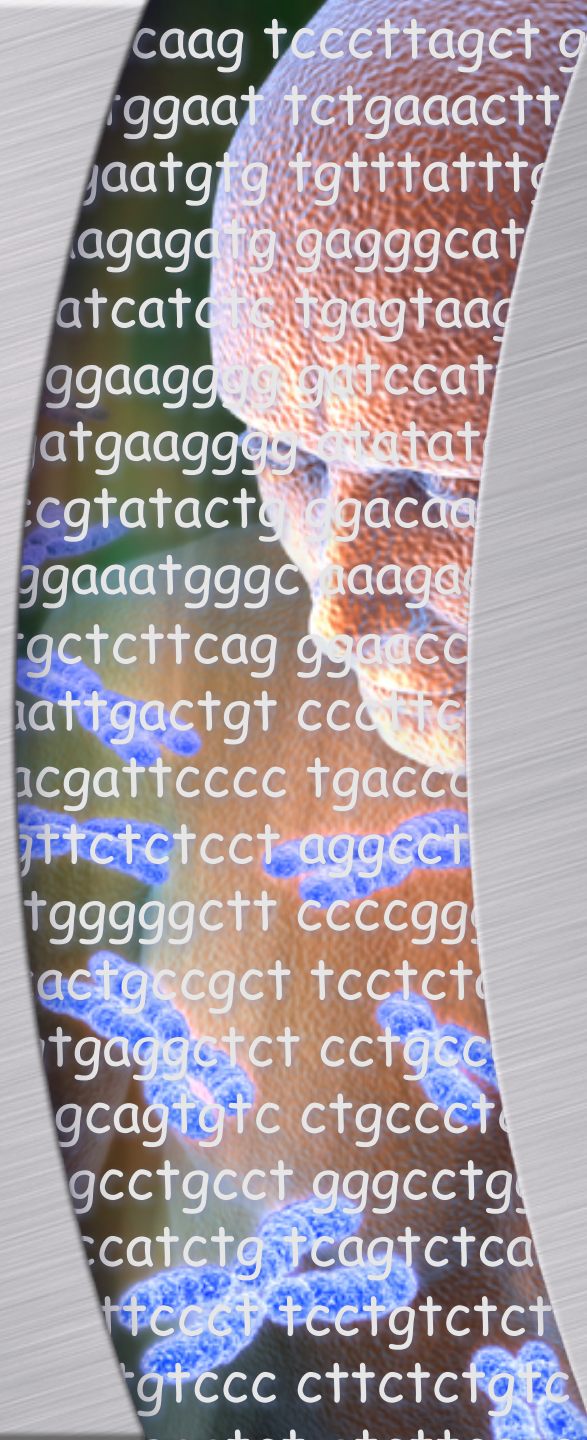
Bioinformatics is presently used in...

- Medical sciences
- Genomics (humans, plants, animals, pathogens)
- Ecology / population genetics
- Functional genomics
- Systems biology
- Pharmaceutical development
- Molecular biology

Bioinformatics and ... gene/protein sequence analysis

- ✓ identify new gene/protein
- ✓ predict protein function
- ✓ identify active sites, functional domains
- ✓ evolution
- ✓ 3D structure
- ✓ protein design
- ✓ identify inhibitors of protein function

Enter Bioinfo



International institutions centralize bioinformatic information, tools and processes

NCBI: National Center for Biotechnology Information, National library of Medicine and National Institutes of Health

EBI: European Bioinformatics Institute

PDB: Protein Data Bank.

The NCBI Homepage

- ✓ All Databases
- PubMed
- Protein
- Nucleotide
- GSS
- EST
- Structure
- Genome
- BioProject
- BioSample
- BioSystems
- Books
- Conserved Domains
- Clone
- dbGaP
- dbVar
- Epigenomics
- Gene
- GEO DataSets
- GEO Profiles
- HomoloGene
- MeSH
- NCBI Web Site
- NLM Catalog
- OMIA
- OMIM
- PMC
- PopSet
- Probe
- Protein Clusters
- PubChem BioAssay
- PubChem Compound
- PubChem Substance
- PubMed Health
- SNP
- SRA
- Taxonomy
- ToolKit
- ToolKitAll
- UniGene
- UniSTS



Search

- NCBI Home
- Resource List (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

Welcome to NCBI

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

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

- [Analyze data using NCBI software](#)
- [Downloads](#): Get NCBI data or software
- [Tutorials](#): Learn how to accomplish specific tasks at NCBI
- [Submissions](#): Submit data to GenBank or other NCBI databases

YouTube channel

How to get the most out of NCBI databases with video tutorials
Visit the NCBI YouTube Channel.

GO



1 2 3 4 5 6 7 8

Popular Resources

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



NCBI Announcements

BLAST+ v.2.2.26 is now available 08 Mar 2012

The new BLAST+ release contains a number of important changes and improvements, including the addition of

NIH/NCBI launches the Genetic Testing Registry (GTR) 29 Feb 2012

The GTR is a new resource for finding information about genetic tests and test

NCBI Discovery Workshop: A Practical Hands-On Course 24 Jan 2012

February 21-22, 2012 @ the NIH: Space is still available in the 2-day Discovery

http://www.ncbi.nlm.nih.gov/genome

NCBI Resources How To

Genome

Genome Information by organism

Overview [7012] Eukaryotes [1664] Prokaryotes

First Previous Shown: 1 - 100 out of 7012

Organism/Name	Kingdom
Abaca bunchy top virus	Viruses
Abalone herpesvirus Taiwan/2004	Viruses
Abalone shriveling syndrome-associated virus	Viruses
Abelson murine leukemia virus	Viruses
Abiotrophia defectiva	Bacteria
Abutilon Brazil virus	Viruses
Abutilon mosaic Bolivia virus	Viruses
Abutilon mosaic Brazil virus	Viruses
Abutilon mosaic virus	Viruses
Acacia mangium	Eukaryotes
Acanthamoeba castellanii	Eukaryotes
Acanthamoeba polyphaga mimivirus	Viruses
Acanthascus dawsoni	Eukaryotes
Acanthocheilonema viteae	Eukaryotes
Acanthocystis turfacea	Viruses
Chlorella virus 1	Viruses
Acaryochloris marina	Bacteria
Acaryochloris phage A-HIS1	Viruses
Acaryochloris sp. CCMEF	Bacteria

- ✓ All
- All Eukaryotes --
- Animals
- Fungi
- Other
- Plants
- Protists
- All Archaea --
- Crenarchaeota
- Euryarchaeota
- Korarchaeota
- Nanoarchaeota
- Thaumarchaeota
- unclassified Archaea
- All Bacteria --
- Actinobacteria
- Aquificae
- Armatimonadetes
- Bacteroidetes/Chlorobi group
- Caldiseica
- Chlamydiae/Verrucomicrobia group
- Chloroflexi
- Chrysiogenetes
- Cyanobacteria
- Deferribacteres
- Deinococcus-Thermus
- Dictyoglomi
- Elusimicrobia
- Fibrobacteres/Acidobacteria group
- Firmicutes
- Fusobacteria
- Gemmatimonadetes
- Nitrospirae
- Planctomycetes
- Proteobacteria
- Spirochaetes
- Synergistetes
- Tenericutes
- Thermodesulfobacteria
- Thermotogae
- unclassified Bacteria
- All Viruses --
- Avsunviridae
- Deltavirus
- Pospiviroidae
- Retro-transcribing viruses
- Satellites
- dsDNA viruses, no RNA stage
- dsRNA viruses
- ssDNA viruses
- ssRNA viruses
- unassigned viruses
- unclassified archaeal viruses
- unclassified phages
- unclassified viroids
- unclassified virophages
- unclassified viruses

access and storage databank ncbi

Sonic Radio

Top Sites

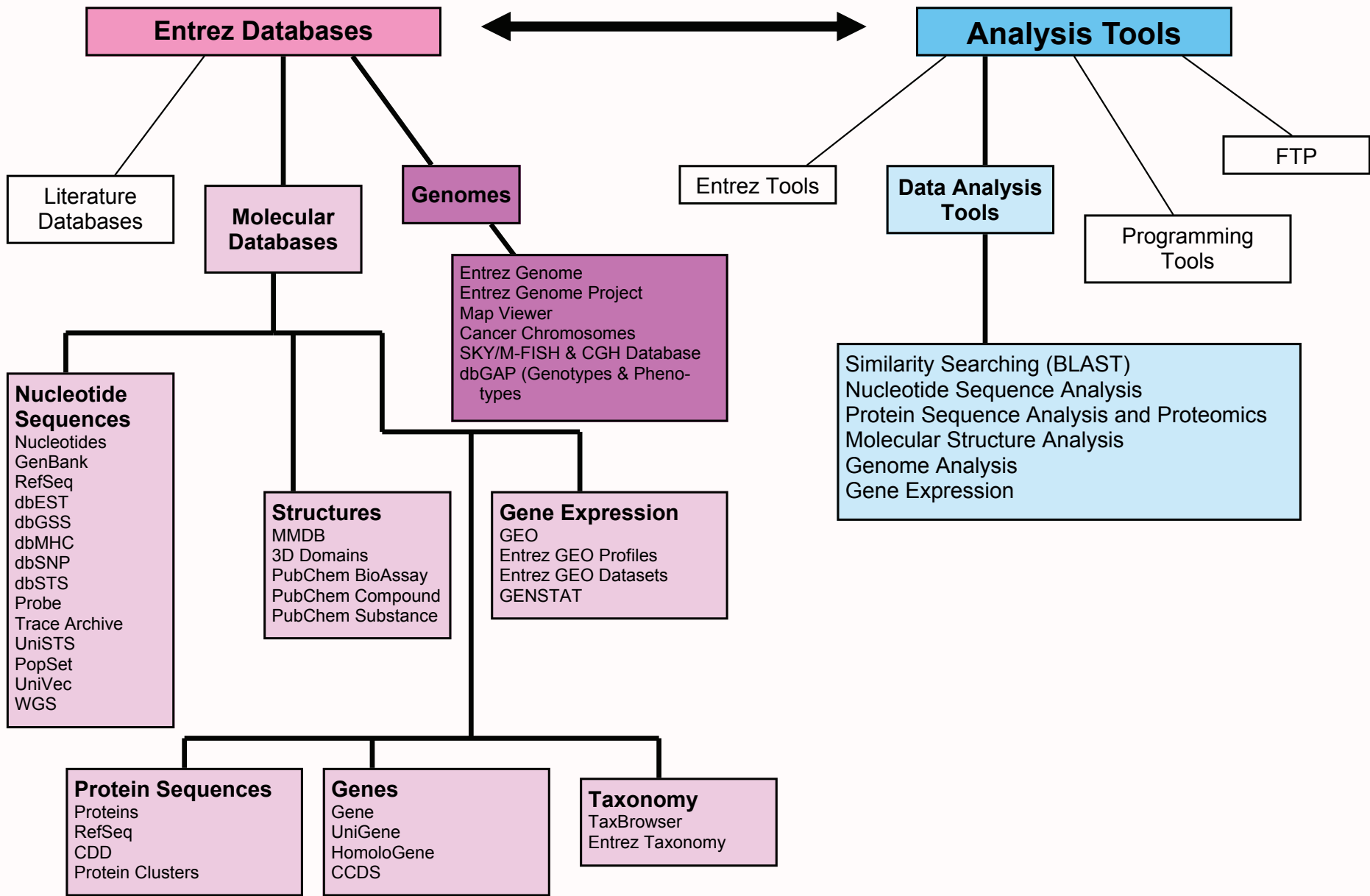
vidall_mcgill My NCBI Sign Out

Search

Download Reports from FTP site

Download selected records

SubGroup	Size (Mb)	Chr	Organelles	Plasmids	BioProjects
Nanoviridae	0.006	6	-	-	1
unclassified	0.044	1	-	-	1
unclassified	0.035	1	-	-	1
Retroviridae	0.006	1	-	-	1
Bacilli	3.48	-	-	-	1
Geminiviridae	0.005	2	-	-	1
Geminiviridae	0.005	2	-	-	1
Geminiviridae	0.005	2	-	-	1
Geminiviridae	0.005	2	-	-	1
Land Plants	0	13	-	-	1
Other Protists	46.43	-	1	-	1
Mimiviridae	1.18	1	-	-	1
Other Animals	0	-	-	-	1
Roundworms	0.014	-	1	-	1
Phycodnaviridae	0.29	1	-	-	1
Chroococcales	8.36	1	-	9	1
Siphoviridae	0	-	-	-	1
Chroococcales	7.88	-	-	-	1

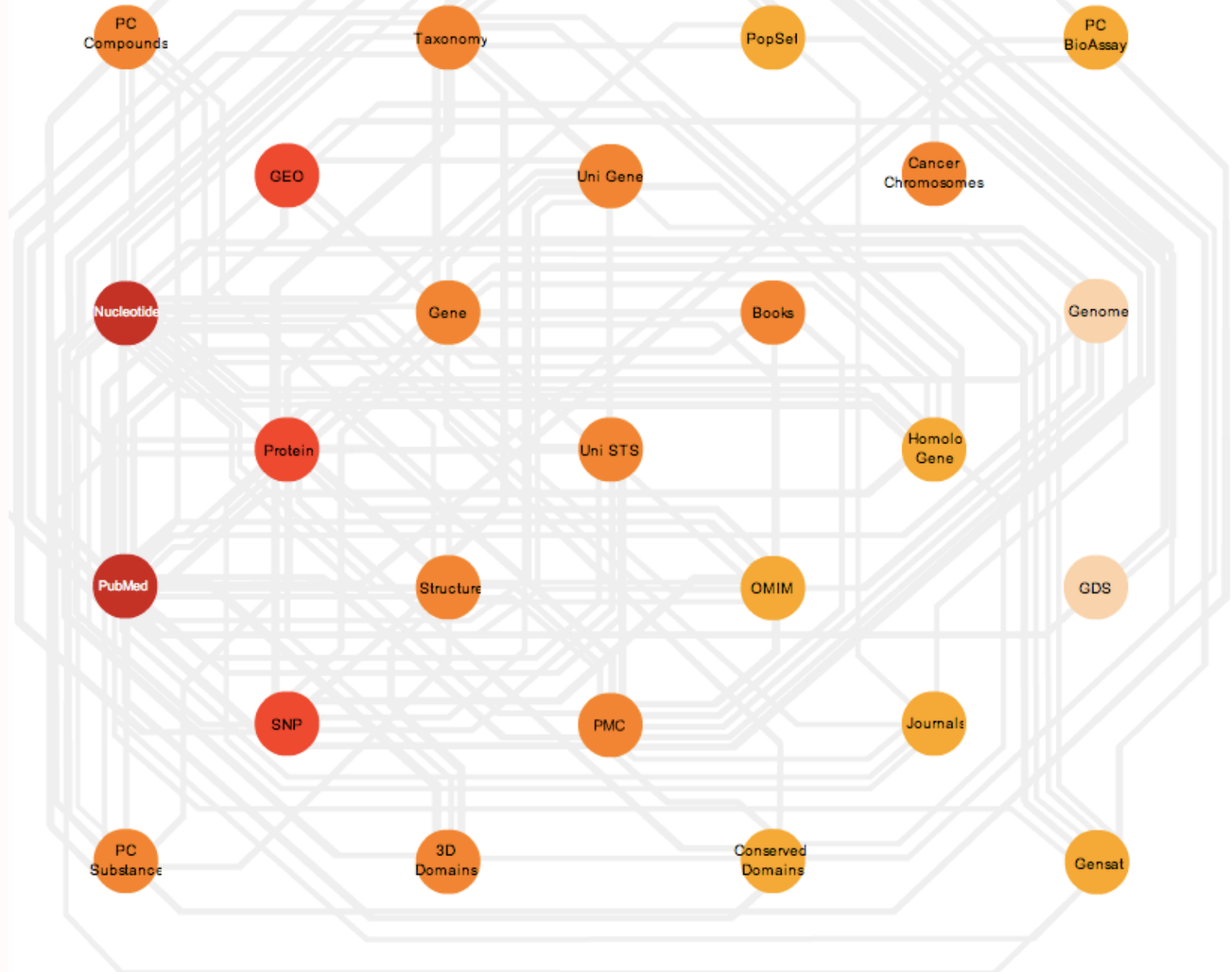


Created by Lyra Eugenio-Thomson.

For more complete information on NCBI Databases and Tools, go to:

<http://www.ncbi.nlm.nih.gov/Sitemap/index.html>

NCBI databases are integrated










Search across databases






















GO

CLEAR

Help

Welcome to the new Entrez cross-database search page

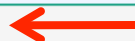
 PubMed: biomedical literature citations and abstracts ?	 Books: online books ?
 PubMed Central: free, full text journal articles ?	 OMIM: online Mendelian Inheritance in Man ?
	 Site Search: NCBI web and FTP sites ?

 Nucleotide: sequence database (GenBank) ?	 UniGene: gene-oriented clusters of transcript sequences ?
 Protein: sequence database ?	 CDD: conserved protein domain database ?
 Genome: whole genome sequences ?	 3D Domains: domains from Entrez Structure ?
 Structure: three-dimensional macromolecular structures ?	 UniSTS: markers and mapping data ?
 Taxonomy: organisms in GenBank ?	 PopSet: population study data sets ?
 SNP: single nucleotide polymorphism ?	 GEO Profiles: expression and molecular abundance profiles ?
 Gene: gene-centered information ?	 GEO DataSets: experimental sets of GEO data ?
 HomoloGene: eukaryotic homology groups ?	 Cancer Chromosomes: cytogenetic databases ?
 PubChem Compound: small molecule chemical structures ?	 PubChem BioAssay: bioactivity screens of chemical substances ?
 PubChem Substance: chemical substances screened for bioactivity ?	 GENSAT: gene expression atlas of mouse central nervous system ?
 Genome Project: genome project information ?	



Search across databases

mdr1
























































GO

Clear

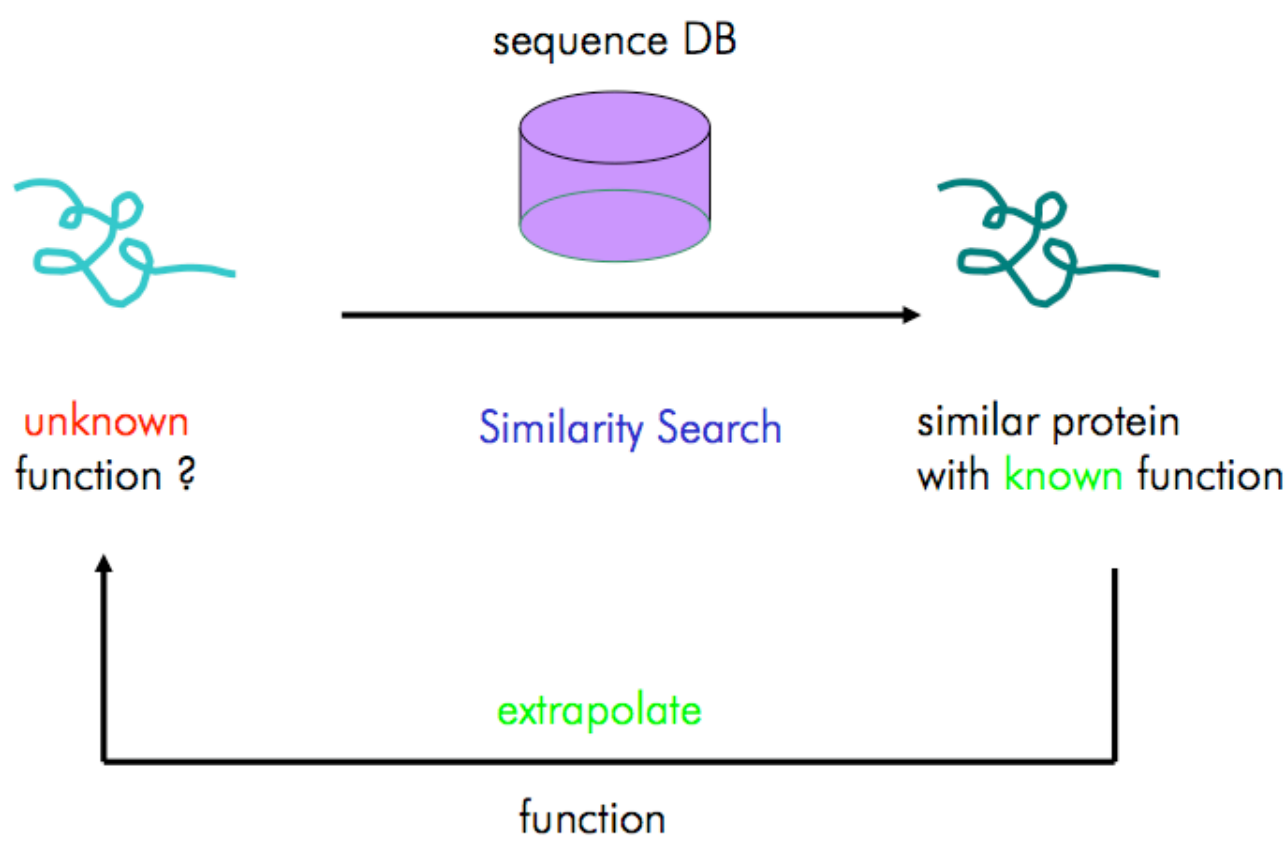
Help

- Result counts displayed in gray indicate one or more terms not found

5026		PubMed: biomedical literature citations and abstracts		59		Books: online books	
4095		PubMed Central: free, full text journal articles		15		OMIM: online Mendelian Inheritance in Man	
1		Site Search: NCBI web and FTP sites					
2132		Nucleotide: Core subset of nucleotide sequence records		none		dbGaP: genotype and phenotype	
2		EST: Expressed Sequence Tag records		12		UniGene: gene-oriented clusters of transcript sequences	
11		GSS: Genome Survey Sequence records		3		CDD: conserved protein domain database	
2779		Protein: sequence database		4389		Clone: integrated data for clone resources	
24		Genome: whole genome sequences		1		UniSTS: markers and mapping data	
2		Structure: three-dimensional macromolecular structures		17		PopSet: population study data sets	
none		Taxonomy: organisms in GenBank		5370		GEO Profiles: expression and molecular abundance profiles	
7956		SNP: short genetic variations		23		GEO DataSets: experimental sets of GEO data	
168		dbVar: Genomic structural variation		none		Epigenomics: Epigenetic maps and data sets	
749		Gene: gene-centered information		1570		PubChem BioAssay: bioactivity screens of chemical substances	
none		SRA: Sequence Read Archive		none		PubChem Compound: unique small molecule chemical structures	
81		BioSystems: Pathways and systems of interacting molecules		none		PubChem Substance: deposited chemical substance records	

1. **BLAST**: similarity search through sequence alignment

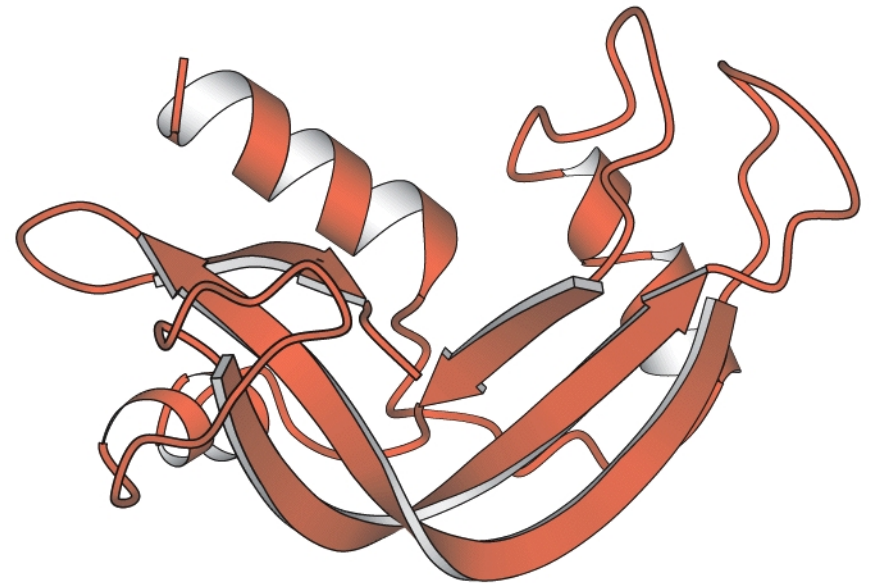
Value of sequence alignment



sequence/structure/function homology in
proteins from a common ancestor
(homologous proteins)



Bovine ribonuclease



Human ribonuclease

Structures of Ribonucleases from Cows and Human Beings. Structural similarity often follows functional similarity.

HOMOLOGY: The presence of a similar feature because of descent from a common ancestor

SIMILARITY: quantitative measurement of identities, gaps, conservative substitutions between protein sequences.

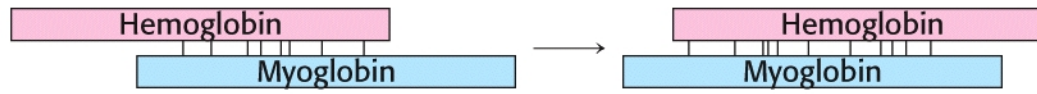
Human hemoglobin (α chain)

V L S P A D K T N V K A A W G K V G A H A G E Y G A E A L E R M F L S F P T T K T Y F P H F D L S H G
S A Q V K G H G K K V A D A L T N A V A H V D D M P N A L S A L S D L H A H K L R V D P V N F K L L S
H C L L V T L A A H L P A E F T P A V H A S L D K F L A S V S T V L T S K Y R

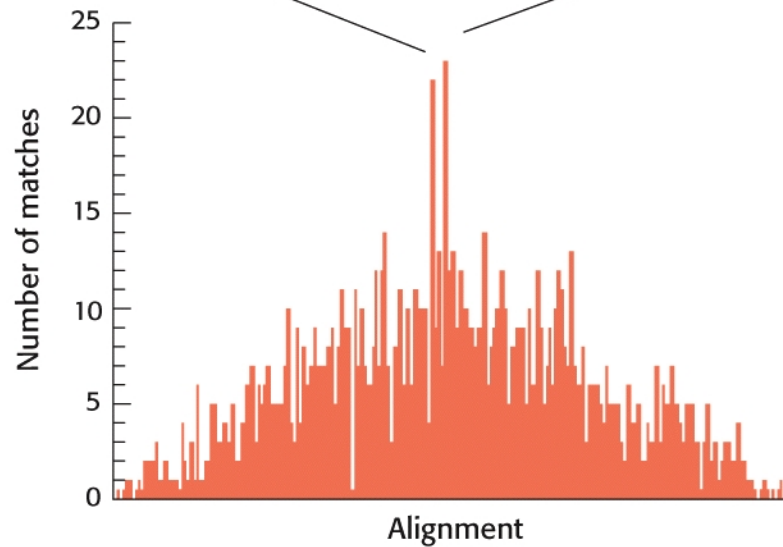
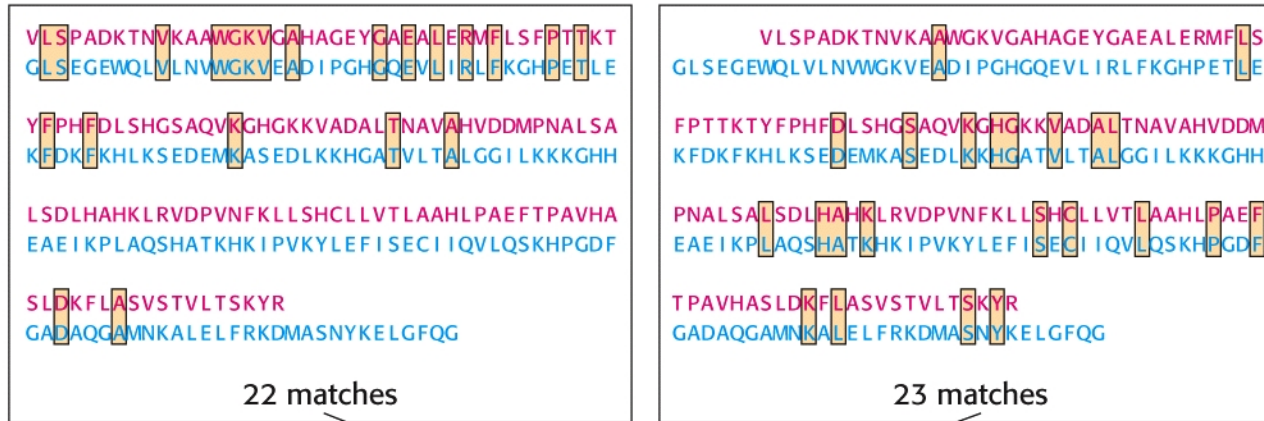
Human myoglobin

G L S D G E W Q L V L N W W G K V E A D I P G H G Q E V L I R L F K G H P E T L E K F D K F K H L K S
E D E M K A S E D L K K H G A T V L T A L G G I L K K K G H H E A E I K P L A Q S H A T K H K I P V K
Y L E F I S E C I I Q V L Q S K H P G D F G A D A Q G A M N K A L E L F R K D M A S N Y K E L G F Q G

(A)



(B)



Hemoglobin α V L S P A D K T N V K A A W G K V G A H A G E Y G A E A L E R M F L S F P T T K T Y F P H F Gap D

Myoglobin G L S E G E W Q L V L N W G K V E A D I P G H G Q E V L I R L F K G H P E T L E K F D K F K H L K S E D

L S H G S A Q V K G H G K K V A D A L T N A V A H V D D M P N A L S A L S D L H A H K L R V D P V N K K L

E M K A S E D L K K H G A T V L T A L G G I L K K K G H H E A E I K P L A Q S H A T K H K I P V K Y L E F

L S H C L L V T L A A H L P A E F T P A V H A S L D K F L A S V S T V L T S K Y R

I S E C I I Q V L Q S K H P G D F G A D A Q G A M N K A L E L F R K D M A S N Y K E L G F Q G

Alignment with gap insertion

Hemoglobin α
Myoglobin

V	L	S	P	A	D	K	T	N	V	K	A	A	W	G	K	V	G	A	H	A	G	E	Y	G	A	E	A	L	E	R	M	F	L	S	F	P	T	T	K	T	Y	F	P	H	F	-----				
G	L	S	E	G	E	W	Q	L	V	L	N	W	G	K	V	E	A	D	I	P	G	H	G	Q	E	V	L	I	R	L	F	K	G	H	P	E	T	L	E	K	F	D	K	F	K	H	L	K	S	
-	D	L	S	H	G	S	A	Q	V	K	G	H	G	K	K	V	A	D	A	L	T	N	A	V	A	H	V	D	D	M	P	N	A	L	S	A	L	S	D	L	H	A	H	K	L	R	V	D	P	V
E	D	E	M	K	A	S	E	D	L	K	K	H	G	A	T	V	L	T	A	L	G	G	I	L	K	K	K	G	H	E	A	E	I	K	P	L	A	Q	S	H	A	T	K	H	K	I	P	V	K	
N	F	K	L	S	H	C	L	L	V	T	L	A	A	H	L	P	A	E	F	T	P	A	V	H	A	S	L	D	K	F	L	A	S	V	S	T	V	L	T	S	K	Y	R							
Y	L	E	F	I	S	E	C	I	I	Q	V	L	Q	S	K	H	P	G	D	F	G	A	D	A	Q	G	A	M	N	K	A	L	E	L	F	R	K	D	M	A	S	N	Y	K	E	L	G	F	Q	

Alignment with conservative substitutions

BLAST

A method to ascertain sequence **similarity**

- 1- The program takes a *query* sequence and searches it against the **database** selected by user.
- 2- It aligns a *query* sequence against each of all *subject* sequence in the database.
- 3- The results are reported in a form of a ranked list followed by a series of individual sequence alignments, plus various **statistics and scores**.

Sequence Analysis:

Basic Local Alignment Search Tool

```
--T--CC-C-AGT--TATGT-CAGGGGACACG--A-GCATGCAGA-GAC
|  ||| |  ||  |  |  |  |||  |  |  |  |  ||||  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG--T-CAGAT--C

                tccCAGTTATGTcAGgggacacgagcatgcagagac
                |||||
aattgcccgcgctcgttttcagCAGTTATGTcAGatc
```

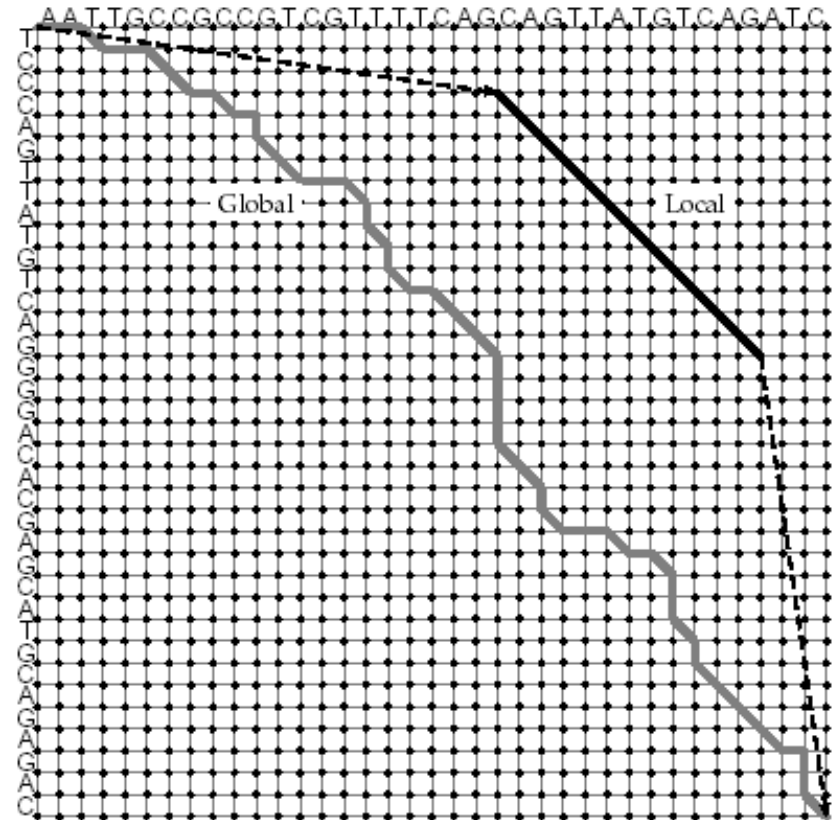


Figure 6.16 (a) Global and (b) local alignments of two hypothetical genes that each have a conserved domain. The local alignment has a much worse score according to the global scoring scheme, but it correctly locates the conserved domain.



► NCBI/ BLAST/ blastp suite

[blastn](#) | **[blastp](#)** | [blastx](#) | [tblastn](#) | [tblastx](#)

BLASTP programs search protein subjects using a protein query. [more...](#)

Enter Query Sequence

Enter accession number, gi, or FASTA sequence

[Clear](#)

Query subrange

From

To

Or, upload file

no file selected

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Enter Subject Sequence

Enter accession number, gi, or FASTA sequence

[Clear](#)

Subject subrange

From

To

Or, upload file

no file selected

Program Selection

Algorithm

blastp (protein-protein BLAST)

[Choose a BLAST algorithm](#)

BLAST

Search [protein sequence](#) using [Blastp \(protein-protein BLAST\)](#)

Show results in a new window

► [Algorithm parameters](#)



BLAST

Search **protein sequence** using **Blastp (protein-protein BLAST)**

Show results in a new window

Algorithm parameters

General Parameters

Max target sequences

100

Select the maximum number of aligned sequences to display

Short queries

Automatically adjust parameters for short input sequences

Expect threshold

10

Word size

3

Scoring Parameters

Matrix

BLOSUM62

Assigns a score for aligning pairs of residues, and determines overall alignment score. [more...](#)

Gap Costs

Existence: 11 Extension: 1

Compositional adjustments

Conditional compositional score matrix adjustment

Filters and Masking

Filter

Low complexity regions

Mask

Mask for lookup table only

Mask lower case letters

BLOSUM-62

BLAST substitution matrices

A key element in evaluating the quality of a pairwise sequence alignment is the "substitution matrix", which assigns a score for aligning any possible pair of residues. The theory of amino acid substitution matrices is described in [1], and applied to DNA sequence comparison in [2]. In general, different substitution matrices are tailored to detecting similarities among sequences that are diverged by differing degrees [1-3]. A single matrix may nevertheless be reasonably efficient over a relatively broad range of evolutionary change [1-3]. Experimentation has shown that the BLOSUM-62 matrix [4] is among the best for detecting most weak protein similarities. For particularly long and weak alignments, the BLOSUM-45 matrix may prove superior. A detailed statistical theory for gapped alignments has not been developed, and the best gap costs to use with a given substitution matrix are determined empirically. Short alignments need to be relatively strong (i.e. have a higher percentage of matching residues) to rise above background noise. Such short but strong alignments are more easily detected using a matrix with a higher "relative entropy" [1] than that of BLOSUM-62. In particular, short query sequences can only produce short alignments, and therefore database searches with short queries should use an appropriately tailored matrix. The BLOSUM series does not include any matrices with relative entropies suitable for the shortest queries, so the older PAM matrices [5,6] may be used instead. For proteins, a provisional table of recommended substitution matrices and gap costs for various query lengths is:

Query Length	Substitution Matrix	Gap Costs
<35	PAM-30	(9,1)
35-50	PAM-70	(10,1)
50-85	BLOSUM-80	(10,1)
85	BLOSUM-62	(10,1)

Gap Costs



results of BLAST

BLASTP 2.2.10 [Oct-19-2004]

Reference:
Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs", Nucleic Acids Res. 25:3389-3402.

RID: 1111270526-12062-108720234474.BLASTQ2

Query=
(153 letters)

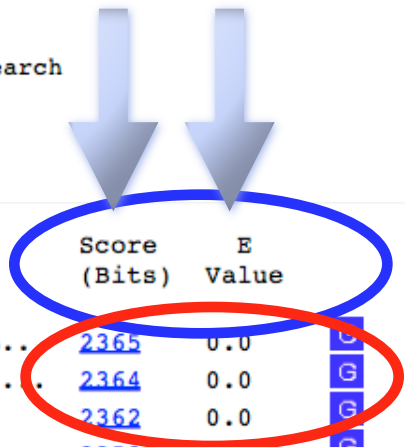
Database: All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples
2,367,365 sequences; 802,797,248 total letters

If you have any problems or questions with the results of this search please refer to the [BLAST FAQs](#)

[Taxonomy reports](#)

Sequences providing significant alignments:

gi 58802449 gb AAW82430.1	ATP-binding cassette, sub-family B..	2365	0.0	G
gi 42741632 gb AAW82430.2	ATP-binding cassette sub-family . .	2364	0.0	G
gi 307180 gb AAA59575.1	P-glycoprotein [Homo sapiens]	2362	0.0	G
gi 2353264 gb AAB69423.1	P-glycoprotein [Homo sapiens]	2358	0.0	G
gi 60326712 gb AAX18881.1	P-glycoprotein [Cercopithecus aethiop	2285	0.0	G
gi 31442763 gb AAN07780.2	multidrug resistance p-glycoprotein [2285	0.0	G
gi 74136329 ref NP_001028059.1	multidrug resistance p-glycop...	2280	0.0	G
gi 46394984 gb AAS91648.1	multidrug resistance protein; P-glyco	2279	0.0	G
gi 46394982 gb AAS91647.1	multidrug resistance protein 1; P-...	2175	0.0	G
gi 67462127 gb AAV67840.1	multidrug resistance protein 1 [Canis	2174	0.0	G



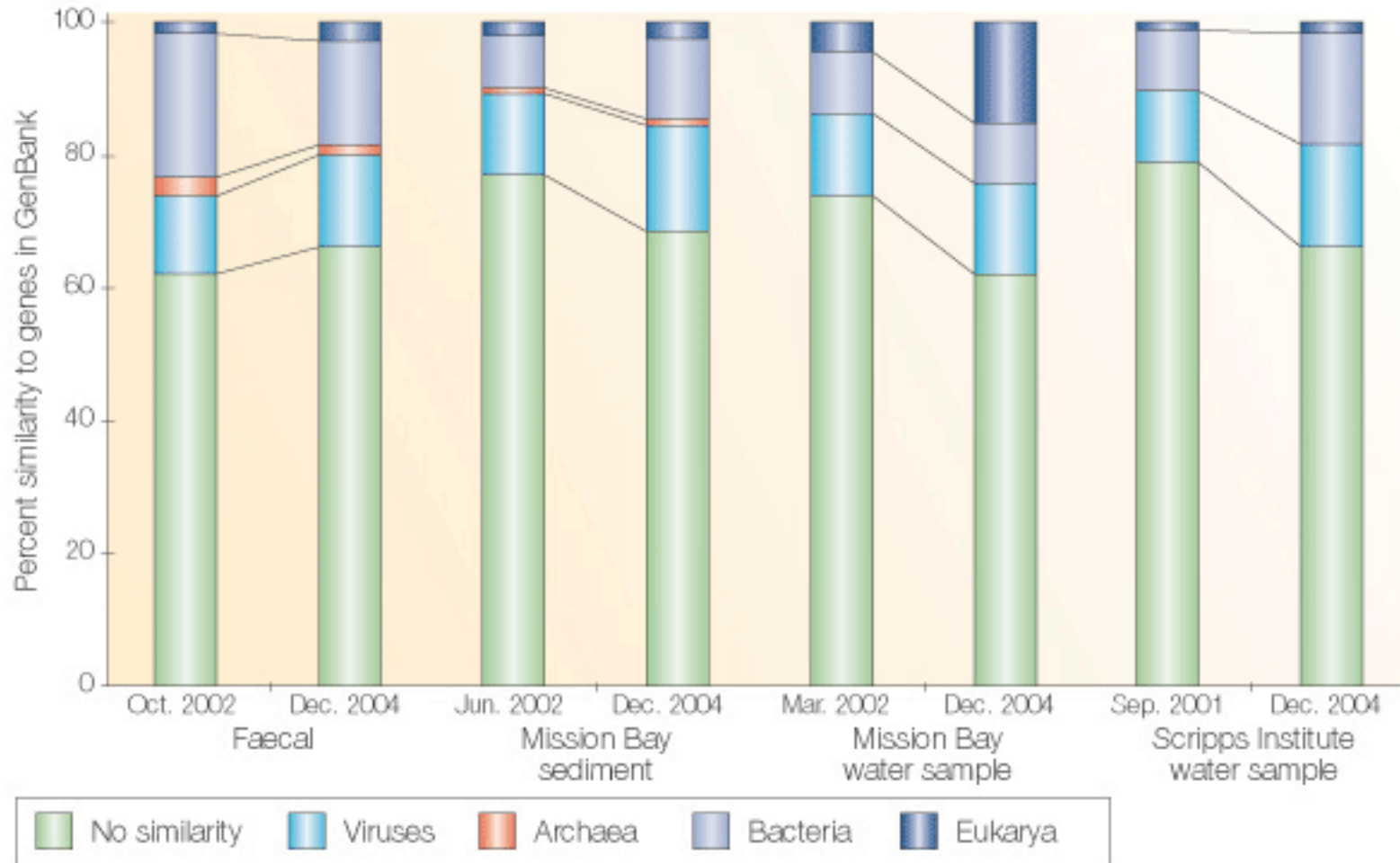
Identity Score: it is the extent of similarity between subject and query sequences. Takes into account similar and identical residues, as well as gaps introduced. The higher the identity score the better the alignment.

The E-Value: expected number of random comparisons in the database scoring above a given threshold; the E value indicates statistical significance of a given alignment. The lower the E value, the most significant is the hit.

USES of similarity searches with BLAST

- ✓ to identify homology between a new and a known sequence
- ✓ to identify new genes in any organism
- ✓ to predict gene structure of a new sequence (exon boundaries, regulatory regions)
- ✓ to identify members of a gene or protein family
- ✓ to infer biochemical function of a new protein
- ✓ to infer domain architecture, secondary and tertiary structure of a protein
- ✓ to infer evolutionary history of a sequence

Viral metagenomics: most viral genomes are novel



Nature Reviews | Microbiology
Vol 3, page 504, 2005

Viral metagenomic sequences from human faeces¹⁰, a marine sediment sample⁹ and two seawater samples² were compared to the GenBank non-redundant database at the date of publication and in December 2004. The percentage of each library that could be classified as Eukarya, Bacteria, Archaea, viruses or showed no similarities (E-value >0.001) is shown.

New hypotheses:


Were DNA viruses involved in the origin of the eukaryotic cell nucleus?

Does evolution tinker with the wide genetic diversity offered by virus sequences, for example, through gene transfer?

2. ClustalW2

multiple sequence alignment



- [Help Index](#)
 - [General Help](#)
 - [Formats](#)
 - [Gaps](#)
 - [Matrix](#)
 - [References](#)
 - [ClustalW2 Help](#)
 - [ClustalW2 FAQ](#)
 - [Jalview Help](#)
 - [Scores Table](#)
 - [Alignment](#)
 - [Guide Tree](#)
 - [Colours](#)
-
- [Similar Applications](#)
 - [Align](#)
 - [Kalign](#)
 - [MAFFT](#)
 - [MUSCLE](#)
 - [T-Coffee](#)
-
- [ClustalW Programmatic Access](#)
-
- www.clustal.org
-
- Clustal Related Literature** 

Search for Clustal related literature in Medline... [more](#)

EBI > Tools > Sequence Analysis > ClustalW2

ClustalW2

ClustalW2 is a general purpose multiple sequence alignment program for DNA or proteins. It produces biologically meaningful multiple sequence alignments of divergent sequences. It calculates the best match for the selected sequences, and lines them up so that the identities, similarities and differences can be seen. Evolutionary relationships can be seen via viewing Cladograms or Phylograms.

[New users, please read the FAQ.](#)

>> Download Software



<p>YOUR EMAIL</p> <input type="text"/>	<p>ALIGNMENT TITLE</p> <input type="text" value="Sequence"/>	<p>RESULTS</p> <input type="button" value="interactive"/>	<p>ALIGNMENT</p> <input type="button" value="full"/>	
<p>KTUP (WORD SIZE)</p> <input type="button" value="def"/>	<p>WINDOW LENGTH</p> <input type="button" value="def"/>	<p>SCORE TYPE</p> <input type="button" value="percent"/>	<p>TOPDIAG</p> <input type="button" value="def"/>	<p>PAIRGAP</p> <input type="button" value="def"/>
<p>MATRIX</p> <input type="button" value="def"/>	<p>GAP OPEN</p> <input type="button" value="def"/>	<p>NO END GAPS</p> <input type="button" value="yes"/>	<p>GAP EXTENSION</p> <input type="button" value="def"/>	<p>GAP DISTANCES</p> <input type="button" value="def"/>
<p>ITERATION</p> <input type="button" value="none"/>			<p>NUMITER</p> <input type="button" value="1"/>	
<p>OUTPUT</p> <p>OUTPUT FORMAT</p> <input type="button" value="aln w/numbers"/>		<p>PHYLOGENETIC TREE</p> <p>OUTPUT ORDER</p> <input type="button" value="aligned"/> <p>TREE TYPE</p> <input type="button" value="none"/> <p>CORRECT DIST.</p> <input type="button" value="off"/> <p>IGNORE GAPS</p> <input type="button" value="off"/> <p>CLUSTERING</p> <input type="button" value="NJ"/>		

Enter or paste a set of sequences in any supported format:

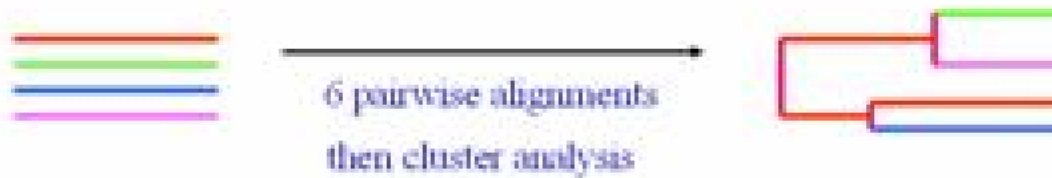
Upload a file: no file selected

ClustalW is a heuristic approach

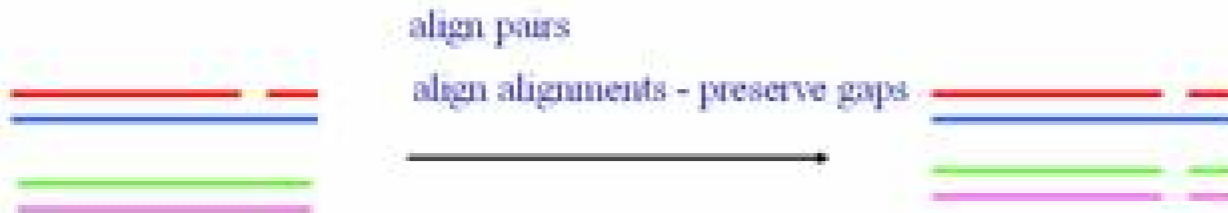
- ✓ It exploits the fact that homologous sequences are evolutionary related
- ✓ Closely related sequences are aligned first gradually adding in the more distant ones.
- ✓ Penalties are increased as more sequences are added
- ✓ Different weight matrices (PAM and BLOSUM series) are used as the alignment proceeds because different matrices are used for highly-related sequences or more divergent sequences

Steps in a multiple alignment

(1) Pairwise alignment



(2) Multiple alignment following the tree from (A)

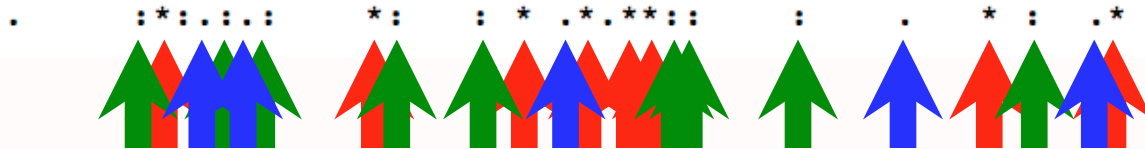


USES of multiple sequence alignments

- ✓ to find diagnostic patterns of protein families
- ✓ to detect homology between new sequences and existing families of sequences
- ✓ to help predict secondary and tertiary structure on new sequences
- ✓ to determine evolutionary relationships: phylogenetics

alignment file

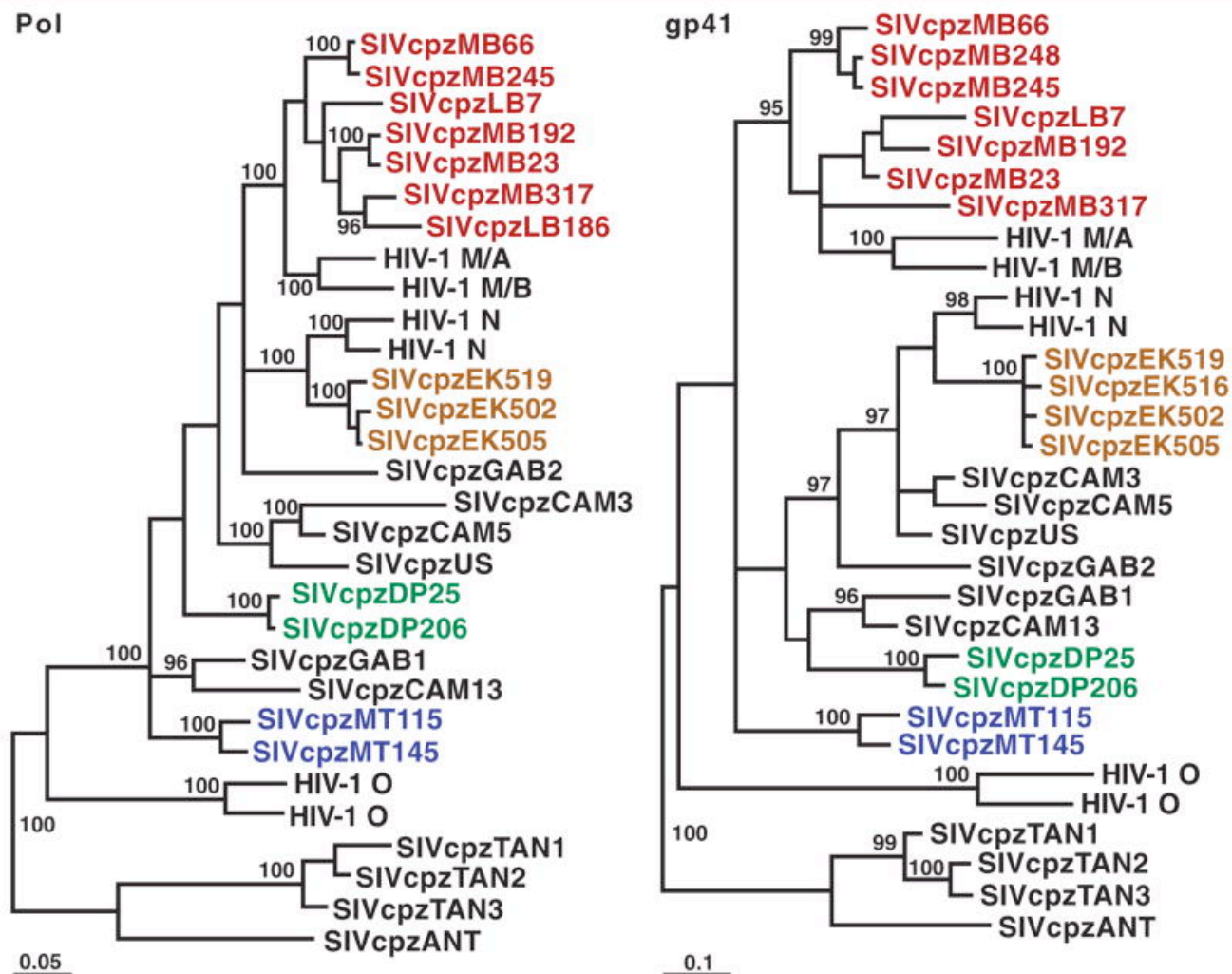
```
MDR1_HUMAN      YPSRKEVKILKGLNLKVQSGQTVALVGNSGCGKSTTVQLMQRLYDPTEGMVSVDGQDIRT 460
MDR3_HUMAN      YPSRANVKILKGLNLKVQSGQTVALVGSSGCGKSTTVQLIQRLYDPDEGTINIDGQDIRN 462
1B0U_A|PDBID    YGG---HEVLKGVSLQARAGDVISIIGSSGSGKSTFLRCINFLEKPSGAIIVNGQNINL  72
STE6_YEAST      YPSRPSEAVLKNVSLNFSAGQFTFIVGKSGSGKSTLSNLLLRFYDGYNGSISINGHNIQT 425
TAP1_HUMAN      YPNRPDVLVLQGLTFTLRPGEVTALVGPNGSGKSTVAALLQONLYQPTGGQLLLDGKPLPQ 571
CFTR_HUMAN      NFSLLGTPVLKDINFKIERGQLLAVAGSTGAGKTSLLMMIMGELEPSEGKIKHSG----- 486
```



- "*" means that the residues or nucleotides in that column are identical in all sequences in the alignment.
- ":" means that conserved substitutions have been observed, according to the CPLOR table at <http://www.ebi.ac.uk/clustalw>.
- "." means that semi-conserved substitutions are observed.

PRESENCE OF STRICTLY CONSERVED AMINO ACID POSITIONS AMONG RELATED PROTEINS INDICATE THE LOCALIZATION OF KEY RESIDUES, WHICH MIGHT BE IMPORTANT FOR FUNCTION

Chimpanzee Reservoirs of Pandemic and Nonpandemic HIV-1



3. Viewer 3D-protein structure

Protein Database

PDB ID or keyword Author **SEARCH**

Home Search

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

Welcome to the RCSB PDB

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

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

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

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

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

Comments? info@rcsb.org

Molecule of the Month: Tissue Factor



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

- More ...
- Previous Features

The RCSB PDB is supported by funds from the National Science Foundation (NSF), the National Institute of General Medical Sciences (NIGMS), the Office of Science, Department of Energy (DOE), the National Library of Medicine (NLM), the National Cancer Institute (NCI), the National Center for Research Resources (NCRR), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), and the National Institute of Neurological Disorders and Stroke (NINDS).

NEWS

- Complete News
- Newsletter
- Discussion Forum

14-Mar-2006

RCSB PDB at Science Expo for NJ Students

On March 21, the RCSB PDB will take part in a Science Expo held at Princeton University for middle school students from New Jersey.

- Full Story ...

07-Mar-2006

RCSB PDB Focus: Frequently Asked Questions

28-Feb-2006

RCSB PDB Exhibit News

21-Feb-2006

Virtual Reality Environment Highlights PDB Structures

14-Feb-2006

PDB Statistics: Structures Solved by Multiple Methods

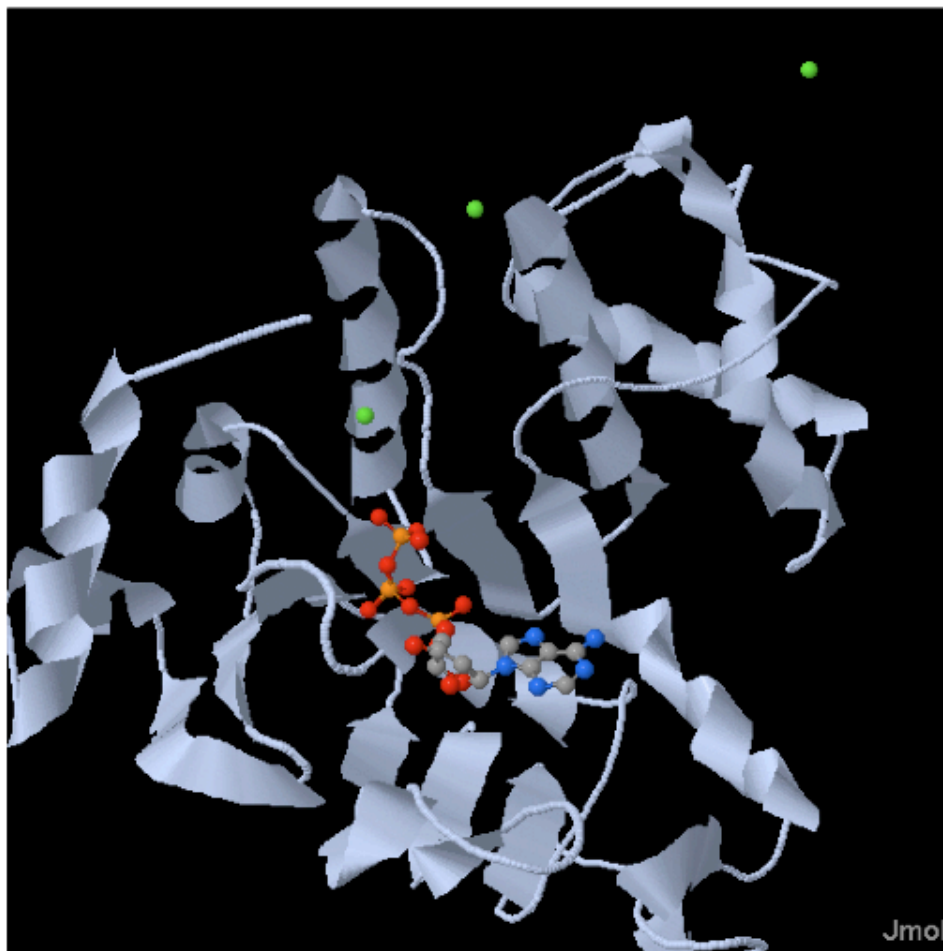
The RCSB PDB is supported by funds from the National Science Foundation (NSF), the National Institute of General Medical Sciences (NIGMS), the Office of Science, Department of Energy (DOE), the National Library of Medicine (NLM), the National Cancer Institute (NCI), the National Center for Research Resources (NCRR), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), and the National Institute of Neurological Disorders and Stroke (NINDS).

- ▼
- 1B0U
- ▶ Download Files
 - FASTA Sequence
- ▼ Display Files
 - Custom Structure Summary
 - PDB File
 - PDB File (Header)
 - mmCIF File
 - mmCIF File (Header)
 - PDBML/XML File
 - PDBML/XML (Header)
- ▼ Display Molecule
 - Image Gallery
 - KING Viewer
 - Jmol Viewer
 - WebMol Viewer
 - Rasmol Viewer (Plugin required)
 - Swiss-PDB Viewer (Plugin required)
 - 🔗 KING Help
 - 🔗 Jmol Help
 - 🔗 WebMol Help
 - 🔗 Protein Workshop Help
 - 🔗 QuickPDB
 - Asymmetric Unit
 - Assumed Biological Molecule 1
- Structural Reports
- ▶ Structure Analysis
- ▶ Help

For help select one of the options below:

- 🔗 [Help interacting with Jmol](#)
- 🔗 [Simple Interaction Guide \(requires flash\)](#)
- 🔗 [Advanced Jmol Help](#)

1B0U



USES of 3D-viewers

- ✓ to annotate 3D-structure
- ✓ to predict functional consequences of mutations
- ✓ to identify drug targets
- ✓ to reveal evolutionary relationships not detected by sequence analysis

Unexpected similarities in coat protein topology suggest a common origin for viruses that infect bacteria and eukarya).

Curr Opin Struct Biol. 2005 Dec;15(6):655-63. Epub 2005 Nov 3

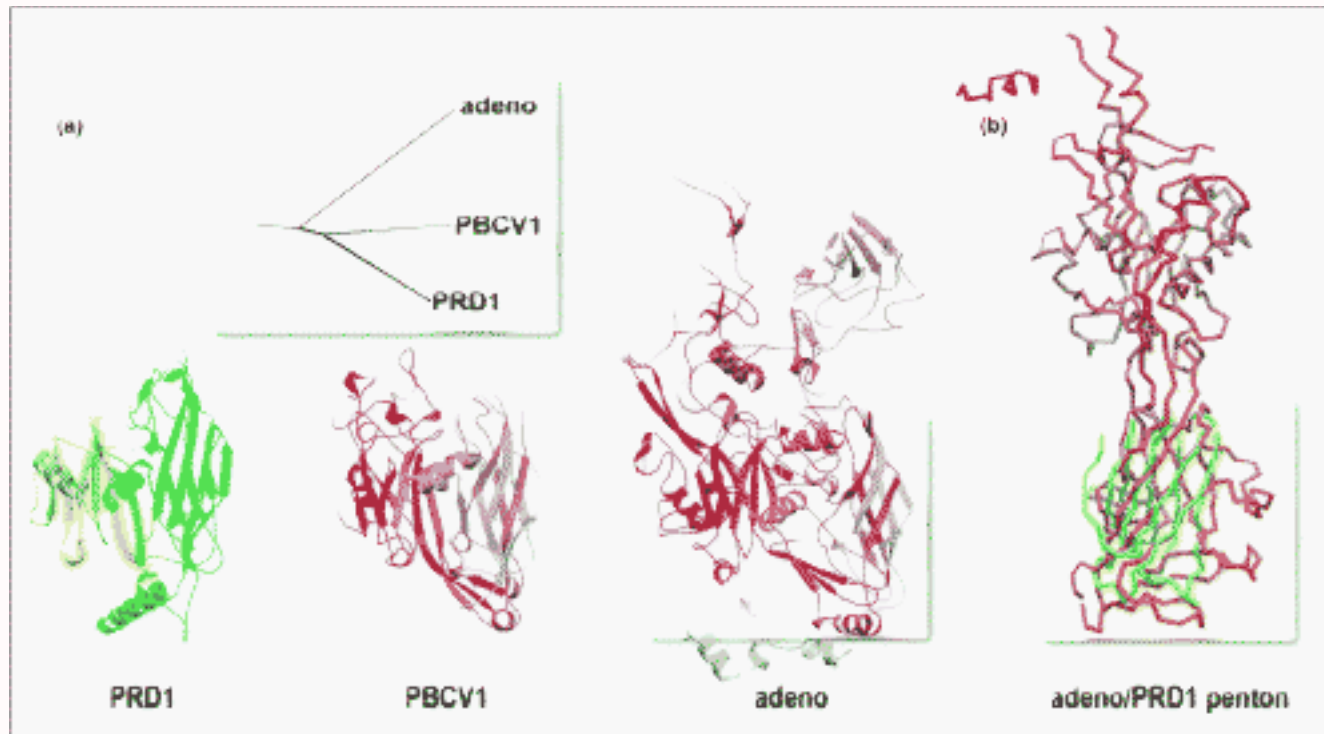


Figure 2. Comparison of the available X-ray-based coat and penton protein structures of PRD₁-like viruses. Colouring is according to the domain of life to which the virus host belongs (green, bacteria; red, eukarya). The structures are presented in side view, so that the lower part is towards the virus interior. (a) The hexon-like proteins of PRD₁, PBCV-1 and adenovirus. For adenovirus, the portion of the molecule comprising the two jelly-rolls is boxed (in performing the structure comparison, 728 of the 887 residues of adenovirus hexon were used, those deemed to comprise the jelly-rolls). Inset above is the relevant portion of the structure-based phylogenetic tree, which was calculated based on these structures. Note that the PBCV-1 protein, although somewhat closer to the equivalent protein of PRD₁ than of adenovirus, also has some structural features (e.g. the lower β extensions) that are closer to adenovirus. (b) The adenovirus penton compared to its stripped-down counterpart from PRD₁. Note that the two β sheets that form the jelly-roll are packed at different angles in the two molecules; nevertheless, the superposition matches 87 out of <100 residues of the PRD₁ jelly-roll.

1 aaataatgta ttggctctgc aaatgcagct tcagaacaag tcccttago
cacc
61 caccctaagt caccaccctt aagcctcacc catgtggaat tctgaaact
gtaga
121 **The Exercise ...** cacattgatc ctggaatgtg tgtttattt
atata
181 aatctgttct gtggaagcca cctgaagtca ggaagagatg gagggcatc
gagtg
241 agatgagacc tcatcactact tgactgtcca gcatcatctc tgagtaagg
aaaa
301 tttatcttcc aaactaggac actttcaaga gtggaagggg gatccatta
cacc
361 tggacaagag gcaaacacca gaatgtcccc gatgaagggg atatataat
cttg
421 atgtgaaacc tgccagatgg gctggaaagt ccgtatactg ggacaagta
gagtt
481 gtttgggaca aggacagggg tacaagagaa ggaaatgggc aaagagaga
actc
541 agccaagggg gcagagatgt tatatatgat tgctcttcag ggaaccggg
gctca
601 caccacagct gctcaaccac ctctctctg aattgactgt cccttcttt

go to the website

<http://teodorolab.mcgill.ca/300D/>

BIOC-300 Bioinformatics Mini-Project (BIMP)

Monday April 2nd, 2012

Contact

Coordinator

Dr Silvia Vidal
silvia.vidal@mcgill.ca [url]
514-398-2362

Meta-Instructions

- **Write a flow sheet.** Not required, but you will find that it is much easier to perform the exercises when you organize yourself - just like in the wet lab.
 - **Use the references/tutorials to guide you.** After attending the lecture, you will want to consolidate what you have seen by looking at the step-by-step tutorials written on each tool used in these exercise. The tutorials are in Chapter 3 of the "Introduction to Bioinformatics", a manual originally written for the Bioinformatics Project in the MicroImm lab course.
 - **Write the report.** A good report should be concise and to the point. Show what you have understood, and present it in an orderly fashion, as suggested below.
-

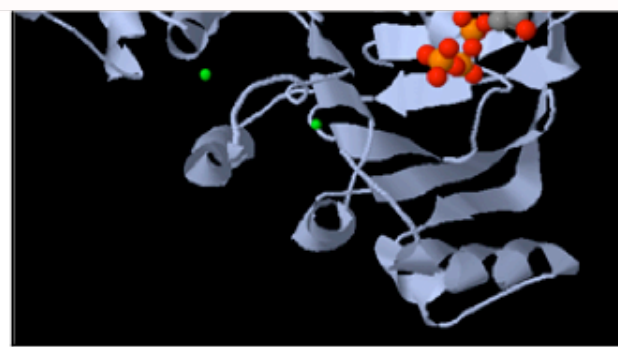
Instructions

This exercise intends to show you how to use some of the most basic tools used in applied bioinformatics. This year, we have chosen to use the ATP-binding cassette (ABC) transporter superfamily to illustrate the use of these tools. The NCBI has a free [online book on ABC transporters](#) on their website.



Part One: BLAST

The [Basic Local Alignment Search Tool](#) (BLAST) finds regions of local similarity between sequences. It is routinely used by millions of people to find sequences in extremely large databases. Typically, we say that we "BLAST" a sequence to say we "search" a sequence, without specifying where we look for it. The pool of sequences to search from is generally the NCBI sequence database, Genbank (which contained in 2005 about 60 billion bp, for 55 million sequence entries!).



The ABC transporter Histidine permease with ATP, seen with Jmol.

Your task is to use the human MDR1 sequence and perform a BLAST to discover similar sequences stored in public databases. We want you to be able to read the BLAST output and parse information from it.

Download: [[Human MDR1 sequence](#)]

Part Two: ClustalW

[ClustalW](#) is a multiple sequence alignment program. On top of aligning sequences, it may also be used to infer phylogenetic trees (while other programs are usually preferred).

Using ClustalW, we compare different ATP-binding cassette (ABC) transporters and try to find conserved motifs. ABC transporters are known to contain sequences that are involved in the binding of the ATP and the catalysis of the ATP hydrolysis reaction, which are present in all members of the superfamily. (What are ABC transporters? See the References section at the bottom of this page.)

Part Three: PDB

The [Protein Data Bank](#) is a repository of structural data, consisting overwhelmingly of protein structures. Much of the data can be appreciated best using standard molecular visualization tools. Classically, [RasMol](#) was used, but several new tools have been developed recently, many of which are web-based and just require [Java](#) to run.

Follow the instructions to:

- 1) identify human MDR1 protein homologs;
- 2) provide multiple sequence alignment of the NBD loop in ABC transporters;
- 3) visualize the structural localization of the ATP binding site in Histidine permease: highlight walker A, linker LSGGQ and walker B domains.

The Report

- ✓ Once you have done the sequence manipulation in each of 3 parts, observe the results. Try to answer question for each point of the exercise. You may need to look at your results again. **Your written answer is required for all questions** in the exercise, either in point form, or in short-essay form.
- ✓ Probably running the programs has taken you 15 min. Observation and reporting of the results might take you much longer but it is absolutely worth it.
- ✓ The report should be at least 2 pages. However, if you are inspired go ahead. You may add figures, tables, references as you need. Please, no BLAST printouts out of the web.



**Thank you and Good
Luck!**

