



# Bioinformatics: The Management and Analysis of Biological Data



Elliot J. Lefkowitz  
‘OMICS – MIC 741  
January 29, 2016

# Contact Information: Elliot Lefkowitz

## Professor, Microbiology



- ◆ Email
  - ◆ ElliotL@uab.edu
  
- ◆ Web Site
  - ◆ <http://bioinformatics.uab.edu>
  
- ◆ Office
  - ◆ BBRB 277A
  
- ◆ Phone
  - ◆ 934-1946

# Objectives

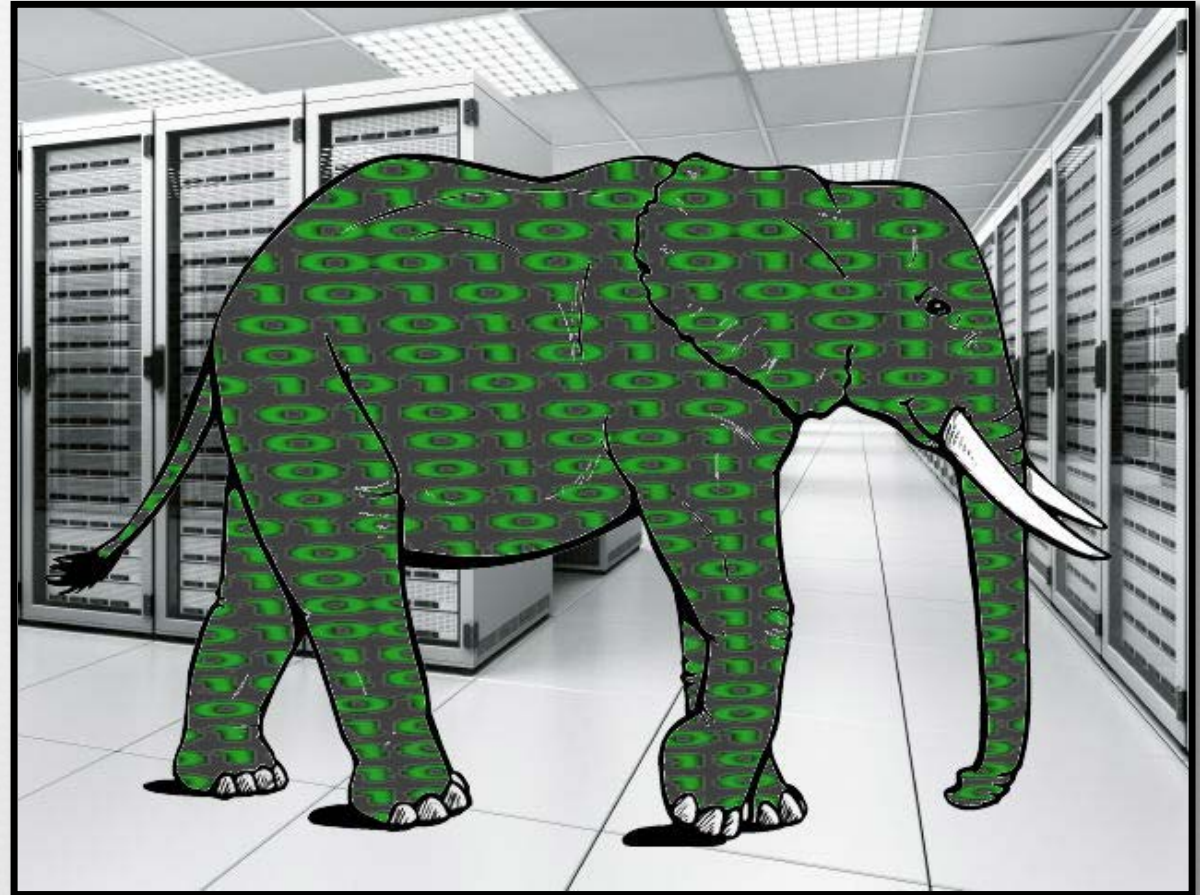


- ◆ Defining and understanding the role of Bioinformatics in biological sciences
- ◆ Becoming familiar with the basic bioinformatic vocabulary
- ◆ Providing an overview of biomedical data and databases
- ◆ Providing an overview of biomedical analytical tools
- ◆ Learning how to discover, access, and utilize information resources

# Biological Big Data

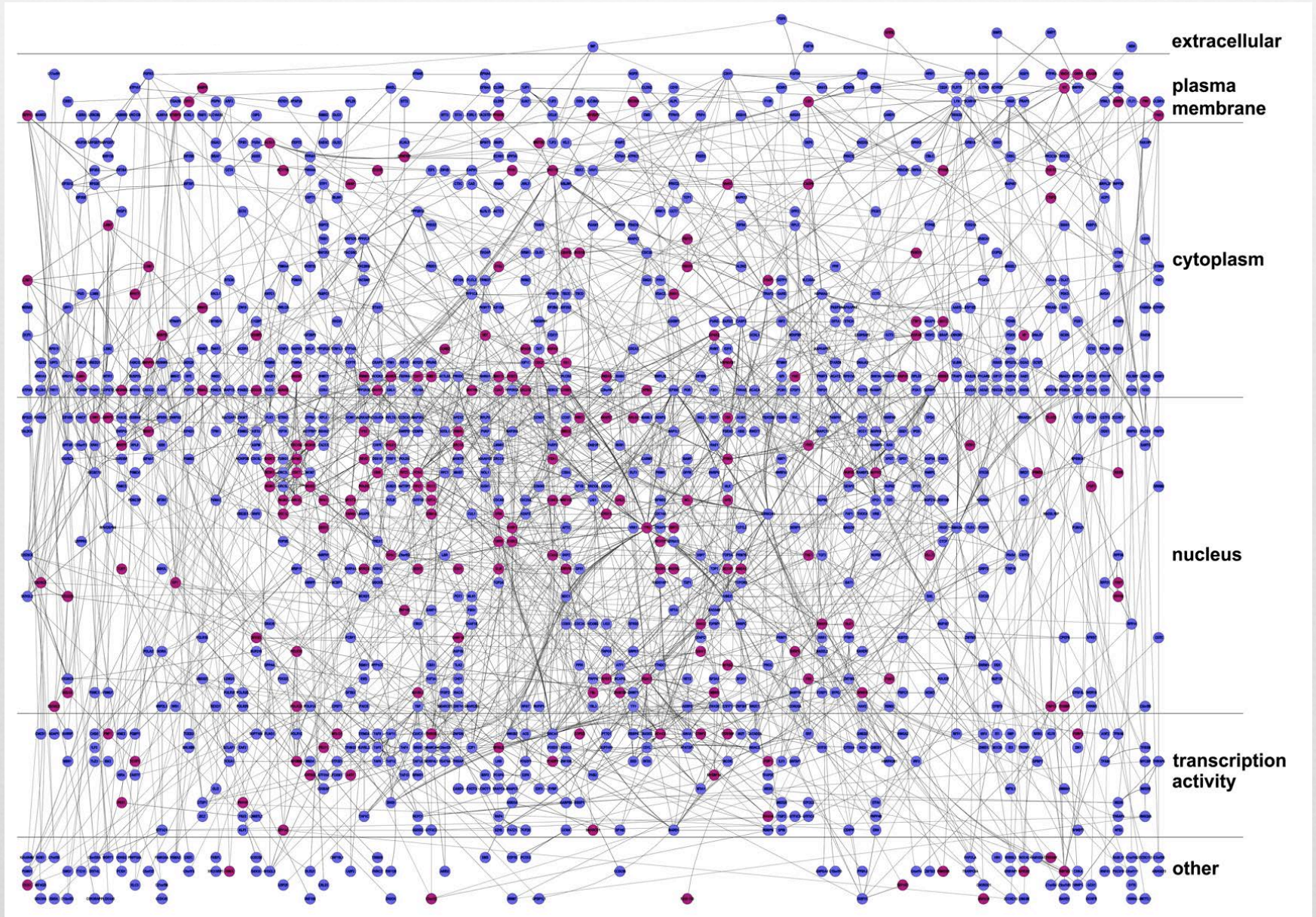


- ◆ Genomic
- ◆ Transcriptomic
- ◆ Epigenetic
- ◆ Proteomic
- ◆ Metabolomic
- ◆ Glycomic
- ◆ Lipidomic
- ◆ Imaging
- ◆ Health
- ◆ Integrated data sets



<http://www.bigdatabytes.com/managing-big-data-starts-here/>

# Systems Biology





# Bioinformatics



# What is Bioinformatics?



- ◆ Computer-aided analysis of biological information

# Bioinformatics = Pattern Discovery



- ◆ Discerning the characteristic (repeatable) patterns in biological information that help to explain the properties and interactions of biological systems.





# Biological Patterns



- ◆ A pattern is a property of a biological system that can be defined such that the definition can be used to detect other occurrences of the same property
  - ◆ Sequence
  - ◆ Expression
  - ◆ Structure
  - ◆ Function
  - ◆ ...

# Bibliography: Web Sites



- ◆ Databases:
  - ◆ Nucleotide: [www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)
  - ◆ Protein: [www.uniprot.org/](http://www.uniprot.org/)
  - ◆ Protein structure: [www.rcsb.org/pdb/](http://www.rcsb.org/pdb/)
  
- ◆ Analysis Tools:
  - ◆ Protein analysis: [us.expasy.org/](http://us.expasy.org/)
  - ◆ UAB Galaxy: [www.uab.edu/galaxy](http://www.uab.edu/galaxy)
  - ◆ Penn State Galaxy: [main.g2.bx.psu.edu/](http://main.g2.bx.psu.edu/)
  - ◆ Comparative genomics: [ecrbrowser.dcode.org](http://ecrbrowser.dcode.org)
  
- ◆ Genome Browsers:
  - ◆ NCBI: ([www.ncbi.nlm.nih.gov/genomes](http://www.ncbi.nlm.nih.gov/genomes))
  - ◆ Ensembl: ([www.ensembl.org](http://www.ensembl.org) )
  - ◆ UCSC Genome Browser: ([genome.ucsc.edu](http://genome.ucsc.edu))

# Bibliography: Journals



- ◆ Nucleic Acids Research
  - ◆ Database issue
    - ◆ Every January
  - ◆ Web Server issue
    - ◆ Every July
  
- ◆ Bioinformatics
  
- ◆ BMC Bioinformatics
  
- ◆ PLoS Computational Biology

# Bioinformatics at UAB

- ◆ Center for Clinical and Translational Sciences (CCTS)
  - ◆ Biomedical Informatics
  - ◆ <http://www.uab.edu/informatics>
  
- ◆ Center for AIDS Research
  - ◆ Molecular and Genetic Bioinformatics Facility
  - ◆ <http://www.genome.uab.edu/>
  
- ◆ Heflin Center for Genetics
  - ◆ <http://www.heflingenetics.uab.edu/>
  
- ◆ Comprehensive Cancer Center's Bioinformatics Shared Facility
  - ◆ <http://www.uab.edu/ccc/biostatistics/>
  
- ◆ Section on Statistical Genetics
  - ◆ <http://www.soph.uab.edu/ssg/>
  
- ◆ UAB-IT Research Computing
  - ◆ <http://www.uab.edu/cores/uab-it-research-computing-services/>

# Bioinformatics Services Provided by CCTS Informatics

- ◆ Provision of shared hardware supporting UAB's bioinformatics needs
- ◆ Access to software
  - ◆ Galaxy
  - ◆ Other supplemental software as needed
- ◆ Access to databases
  - ◆ Sequence, structure, motif, proteomic, expression...
- ◆ Education
  - ◆ Introductory lectures in core graduate courses
  - ◆ Bioinformatics module in upper-level courses
- ◆ Training
  - ◆ Assistance in the use of software and databases
  - ◆ Assistance in planning, organizing, and carrying out bioinformatics analyses
- ◆ Analytical assistance
- ◆ Research Collaborations

# Major Information Resources



- ◆ NCBI
  - ◆ National Center for Biotechnology Information
  - ◆ Databases, tools, links
  - ◆ <http://www.ncbi.nih.gov/>
  
- ◆ EBI
  - ◆ European Bioinformatics Institute
  - ◆ Research Services Training
  - ◆ <http://www.ebi.ac.uk/>

# National Centers for Biomedical Computing



- ◆ National Centers for Biomedical Computing
  - ◆ NIH program to enhance the computational infrastructure for biomedical computing
  - ◆ <http://www.ncbcs.org/>

**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**

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

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

**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

**Genetic Testing Registry**

A portal to clinical genetics resources with detailed information about genetic tests and laboratories.

**GO**

**Popular Resources**

**NCBI Announcements**

NCBI Eukaryotic Genome Annotation Pipeline breaks record; over 100 organisms annotated this year

Nov 20, 2014

[The NCBI Eukaryotic Genome](#)

NCBI BankIt webinar on December 17th

Nov 20, 2014

On December 17th, NCBI will have a webinar entitled "A Submitter's Guide to GenBank: Using BankIt for Small Scale

NCBI E-Utilities webinar video now on YouTube

Nov 13, 2014

October's webinar, "An Introduction to NCBI's E-Utilities on NCBI ADI" is now

[More...](#)

You are here: [NCBI](#) > [National Center for Biotechnology Information](#)

[Write to the Help Desk](#)

**GETTING STARTED**

- [NCBI Education](#)
- [NCBI Help Manual](#)
- [NCBI Handbook](#)
- [Training & Tutorials](#)

**RESOURCES**

- [Chemicals & Bioassays](#)
- [Data & Software](#)
- [DNA & RNA](#)
- [Domains & Structures](#)
- [Genes & Expression](#)
- [Genetics & Medicine](#)
- [Genomes & Maps](#)
- [Homology](#)
- [Literature](#)
- [Proteins](#)
- [Sequence Analysis](#)
- [Taxonomy](#)
- [Training & Tutorials](#)
- [Variation](#)

**POPULAR**

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

**FEATURED**

- [Genetic Testing Registry](#)
- [PubMed Health](#)
- [GenBank](#)
- [Reference Sequences](#)
- [Gene Expression Omnibus](#)
- [Map Viewer](#)
- [Human Genome](#)
- [Mouse Genome](#)
- [Influenza Virus](#)
- [Primer-BLAST](#)
- [Sequence Read Archive](#)

**NCBI INFORMATION**

- [About NCBI](#)
- [Research at NCBI](#)
- [NCBI News](#)
- [NCBI FTP Site](#)
- [NCBI on Facebook](#)
- [NCBI on Twitter](#)
- [NCBI on YouTube](#)



# Analytical Resources



- ◆ Bioinformatics Packages
- ◆ Web-based tools
  - ◆ General Protein Analysis Tools
    - ◆ <http://us.expasy.org/>
  - ◆ Galaxy
    - ◆ <https://usegalaxy.org/>

# Local Bioinformatics Packages



- ◆ Geneious
  - ◆ <http://www.geneious.com>
  
- ◆ MacVector
  - ◆ <http://www.macvector.com>
  
- ◆ Vector NTI
  - ◆ <http://www.invitrogen.com/site/us/en/home/Products-and-Services/Applications/Cloning/vector-nti-software.html?CID=fl-vectornti>
  
- ◆ DNASTar
  - ◆ <http://www.dnastar.com>
  
- ◆ Sequencher
  - ◆ <http://genecodes.com>

# Galaxy



- ◆ Comprehensive set of bioinformatics tools
  - ◆ Basic sequence analysis (EMBOSS)
  - ◆ Next generation Sequencing (NGS) analysis
  - ◆ Statistical analysis
  
- ◆ <http://www.uab.edu/galaxy>
  - ◆ Requires BlazerID login
  
- ◆ <http://main.g2.bx.psu.edu/>

Tools



search tools

[Get Data](#)[Send Data](#)[Demo Tools](#)[ENCODE Tools](#)[Lift-Over](#)[Text Manipulation](#)[Filter and Sort](#)[Join, Subtract and Group](#)[Convert Formats](#)[Extract Features](#)[Fetch Sequences](#)[Get Genomic Scores](#)[Operate on Genomic Intervals](#)[Statistics](#)[Wavelet Analysis](#)[Graph/Display Data](#)[Regional Variation](#)[Multiple regression](#)[Multivariate Analysis](#)[Evolution](#)[Motif Tools](#)[Multiple Alignments](#)[Metagenomic analyses](#)[FASTA manipulation](#)[NCBI BLAST+](#)[NGS TOOLBOX BETA](#)[NGS: QC and manipulation](#)[NGS: Assembly](#)[NGS: Mapping](#)[NGS: Indel Analysis](#)[NGS: RNA Analysis](#)[NGS: SAM Tools](#)[NGS: HA GSL Tools](#)[NGS: Peak Calling](#)[SNP/WGA: Data; Filters](#)[SNP/WGA: QC; LD; Plots](#)[SNP/WGA: Statistical Models](#)[Human Genome Variation](#)[Picard Tools](#)[SnpEff tools](#)[VCF Tools](#)[DebugTools](#)[EMBOSS](#)[BEDTools](#)[NGS: CATK](#)**Galaxy is Implementing a New Deletion Policy**

To combat the growing problem of lack of disk space on UAB Galaxy, we will now be implementing an automated deletion of datasets that have not been updated in more than 6 months.

Deletion will happen on the third Monday of every month, with warning emails being issued the previous two Mondays.

Histories will not be affected by the auto-deletion only datasets.

For more information on UAB Galaxy's deletion policy please see: [Deletion Policy](#)

**Welcome to UAB Galaxy!**

Where all you need is a Blazerid and a web browser to run NGS analyses on the UAB Cheaha Cluster!

**Local Resources**

UAB Galaxy Wiki: [Overview](#), [Data Import](#)

UAB Mailing Lists

UAB Galaxy-users ([search archive](#); [subscribe](#)) discuss with other UAB users

UAB Galaxy-help ask the UAB admins for help!

UAB Cheaha Computing Cluster

Cluster Hardware ([wiki](#))

Request a [cheaha account](#) (needed only for command-line access and bulk data upload)

**Internet Resources**

[Learn Galaxy](#) - tutorials

[Galaxy Project](#) user mailing list ([searchable archives](#); [subscribe](#); [post](#))

[Galaxy Toolshed](#) plug-ins for additional tools that you can request for installation at UAB

Public Galaxy Server at Penn State (PSU): [UseGalaxy.org](#) (more tools, but small disk quotas)

**Brought to you by**

UAB IT [Research Computing](#) under the Office of the Vice President for Information Technology at UAB

UAB CCTS (Center for Clinical and Translational Science under grant UL1 RR025777 from the NIH National Center for Research Resources)

The [Galaxy Platform](#) is developed by Penn State and Emory University

History



search datasets

**Unnamed history**

0 bytes



**i** This history is empty. You can [load your own data](#) or [get data from an external source](#)



Query all databases   [help](#)

### Visual Guidance

### Categories

- proteomics
- genomics
- structural bioinformatics
- systems biology
- phylogeny/evolution
- population genetics
- transcriptomics
- biophysics
- imaging
- IT infrastructure
- drug design

### Resources A..Z

### Links/Documentation

ExPASy is the **SIB Bioinformatics Resource Portal** which provides access to scientific databases and software tools (i.e., *resources*) in different areas of life sciences including proteomics, genomics, phylogeny, systems biology, population genetics, transcriptomics etc. (see **Categories** in the left menu). On this portal you find resources from many different SIB groups as well as external institutions.

### Featuring today

#### World-2DPAGE Repository

Public repository for gel-based proteomics data published in the literature  
[\[details\]](#)



### How to use this portal?



- Features and updates
- New to ExPASy
- Experienced ExPASy users: what is different

### Popular resources

- UniProtKB
- SWISS-MODEL
- STRING
- PROSITE

### Latest News

#### Protein Spotlight: The hidden things

- 2015-01-08

Nature has its secret ways. During the course of the 19th century, the Augustinian friar Gregor Mendel worked out the basics of genetic inheritance as he crossbred pea plants.

[More](#)

#### UniProt Knowledgebase release

2015\_01 - 2015-01-07

#### Release notes

547,357 UniProtKB/Swiss-Prot entries  
[\(More.\)](#)

89,451,166 UniProtKB/TrEMBL entries  
[\(More.\)](#)

[\[More news\]](#) [\[SIB news\]](#)



# Bioinformatic Databases



Something to compare against

# Major Sequence Databases



- ◆ International Nucleotide Sequence Database Collaboration
  - ◆ Genbank
  - ◆ EMBL
  - ◆ DDBJ
  
- ◆ Protein
  - ◆ UniProt
    - ◆ PIR
    - ◆ Swiss-Prot
    - ◆ Swiss-Prot TrEMBL
  
- ◆ NCBI

# Other Databases



- ◆ Structural
  - ◆ Protein Data Bank (PDB): <http://www.rcsb.org/pdb/>
  
- ◆ Expression
  - ◆ Microarray Gene Expression Data Society (MGED): <http://www.mged.org/>
  - ◆ Gene Expression Omnibus (GEO – NCBI): <http://www.ncbi.nlm.nih.gov/geo/>
  
- ◆ Proteomic
  - ◆ Mascot: <http://www.matrixscience.com/>
  
- ◆ Metabolism
  - ◆ BioCyc: <http://biocyc.org/>
  - ◆ Reactome: <http://reactome.org/>
  
- ◆ Ontology
  - ◆ Gene Ontology (GO) Consortium: <http://www.geneontology.org/>
  - ◆ Controlled vocabulary for the description of biological processes



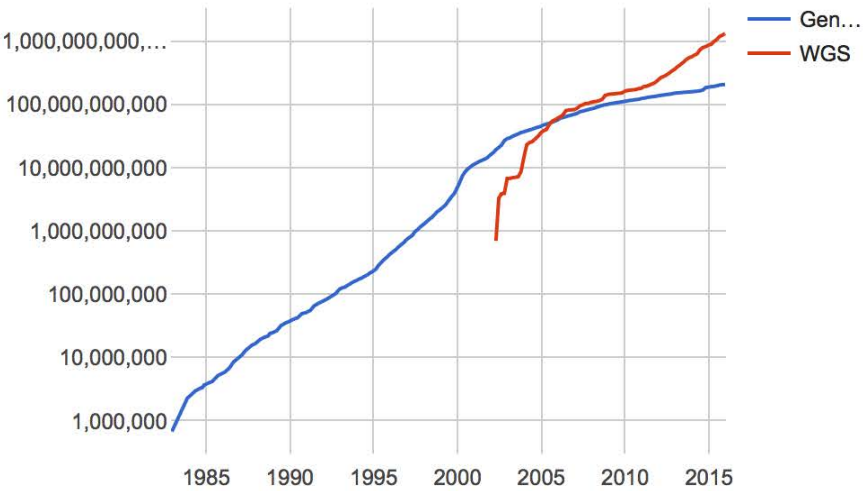
# GenBank



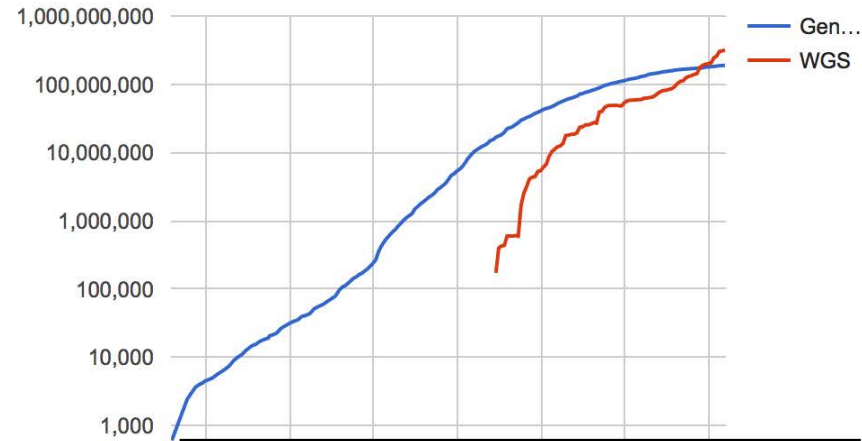
- ◆ Primary nucleic acid sequence database
- ◆ Maintained by NCBI
  - ◆ National Center for Biotechnology Information
  - ◆ <http://www.ncbi.nlm.nih.gov/genbank/>
  - ◆ December, 2016; Release 211
  - ◆ 203,939,111,071 bases
  - ◆ 189,232,925 sequences
- ◆ WGS (Whole Genome Shotgun) Sequences
  - ◆ Separate from GenBank
    - ◆ 1,297,865,618,365 Bases
    - ◆ 317,122,157 Sequences

# GenBank Growth

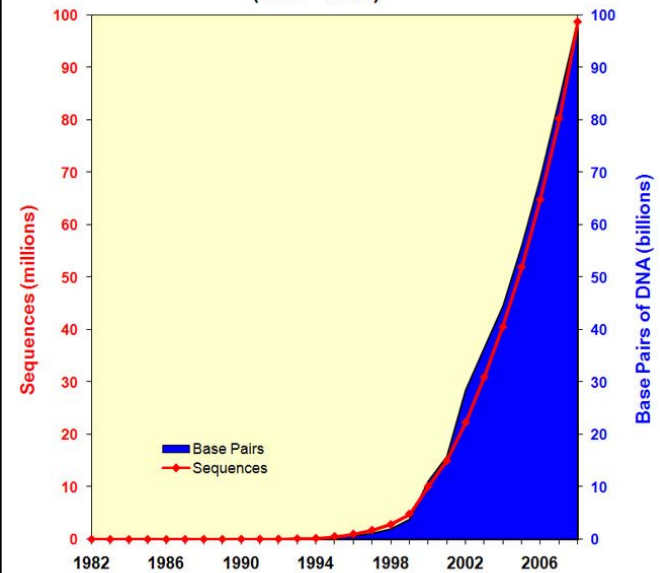
Bases



Sequences



Growth of GenBank  
(1982 - 2008)



# NCBI Databases for NGS Data



- ◆ Bioproject
  - ◆ Genomics, functional genomics, and genetics studies and links to their datasets
  
- ◆ Biosample
  - ◆ Descriptions of biological source materials used in experimental assays.
  
- ◆ Sequence Read Archive (SRA)
  - ◆ Raw sequence data
  
- ◆ Assembly archive

# Other NCBI Databases



- ◆ RefSeq
- ◆ HomoloGene
- ◆ Genomic
- ◆ SNPs
- ◆ ...

# RefSeq



- ◆ NCBI Reference Sequence project
- ◆ Provides reference sequence standards for the naturally occurring molecules from chromosomes to mRNAs to proteins
- ◆ Stable reference point for:
  - ◆ mutation analysis
  - ◆ gene expression studies
  - ◆ polymorphism discovery
- ◆ Accession numbers have two letters, an underscore, and six numbers
  - ◆ NM\_123456

# HomoloGene



- ◆ Database of homologs among the annotated genes of several completely sequenced eukaryotic genomes

## Release 68 Statistics

Release date:2014-04-09

Total species:21

Total HomoloGene groups:44233

[Download](#)

Scientific name	Common name	Genome release	Annotation release	Input genes	Grouped genes	Groups
<i>Homo sapiens</i>	human	<a href="#">GRCh38</a>	106	19806	19129	<a href="#">18732</a>
<i>Pan troglodytes</i>	chimpanzee	<a href="#">Pan_troglodytes-2.1.4</a>	102	22004	18730	<a href="#">18385</a>
<i>Macaca mulatta</i>	Rhesus monkey	<a href="#">Mmul_051212</a>	Build 1.2	22231	16843	<a href="#">16517</a>
<i>Canis lupus familiaris</i>	dog	<a href="#">CanFam3.1</a>	103	19872	18117	<a href="#">17434</a>
<i>Bos taurus</i>	cow	<a href="#">Bos_taurus_UMD_3.1</a>	103	21025	18798	<a href="#">17530</a>
<i>Mus musculus</i>	mouse	<a href="#">GRCm38.p2</a>	104	22627	21207	<a href="#">19033</a>
<i>Rattus norvegicus</i>	rat	<a href="#">Rnor_5.0</a>	104	22892	20616	<a href="#">18871</a>
<i>Gallus gallus</i>	chicken	<a href="#">Gallus_gallus-4.0</a>	102	17133	14600	<a href="#">13352</a>
<i>Xenopus tropicalis</i>	western clawed frog	<a href="#">Xtropicalis_v7</a>	101	22156	18447	<a href="#">14412</a>
<i>Danio rerio</i>	zebrafish	<a href="#">Zv9</a>	103	27270	20897	<a href="#">14559</a>
<i>Drosophila melanogaster</i>	fruit fly	<a href="#">Release 5</a>	Release 5.48	13944	8438	<a href="#">7016</a>
<i>Anopheles gambiae</i>	malaria mosquito	<a href="#">AgamP3</a>	AgamP3.3	12600	8428	<a href="#">6957</a>
<i>Caenorhabditis elegans</i>	nematode	<a href="#">WS195</a>	WS195	20519	7575	<a href="#">4308</a>
<i>Saccharomyces cerevisiae</i>	budding yeast	<a href="#">R64-1-1</a>	NA	5903	4579	<a href="#">4113</a>
<i>Kluyveromyces lactis</i>	ascmycetes	<a href="#">ASM251v1</a>	NA	5084	4283	<a href="#">4211</a>
<i>Eremothecium gossypii</i>	ascmycetes	<a href="#">ASM9102v4</a>	NA	4768	3874	<a href="#">3820</a>
<i>Schizosaccharomyces pombe</i>	fission yeast	<a href="#">ASM294v2</a>	NA	5132	3018	<a href="#">2741</a>
<i>Magnaporthe oryzae</i>	rice blast fungus	<a href="#">MG8</a>	NA	12673	6598	<a href="#">6061</a>
<i>Neurospora crassa</i>	ascmycetes	<a href="#">ASM18292v1</a>	NA	9821	5807	<a href="#">5701</a>
<i>Arabidopsis thaliana</i>	thale cress	<a href="#">TAIR10</a>	NA	27393	19143	<a href="#">10463</a>
<i>Oryza sativa</i>	rice	<a href="#">Build 4.0</a>	NA	28521	16112	<a href="#">9787</a>

# Multi-Genome Comparison

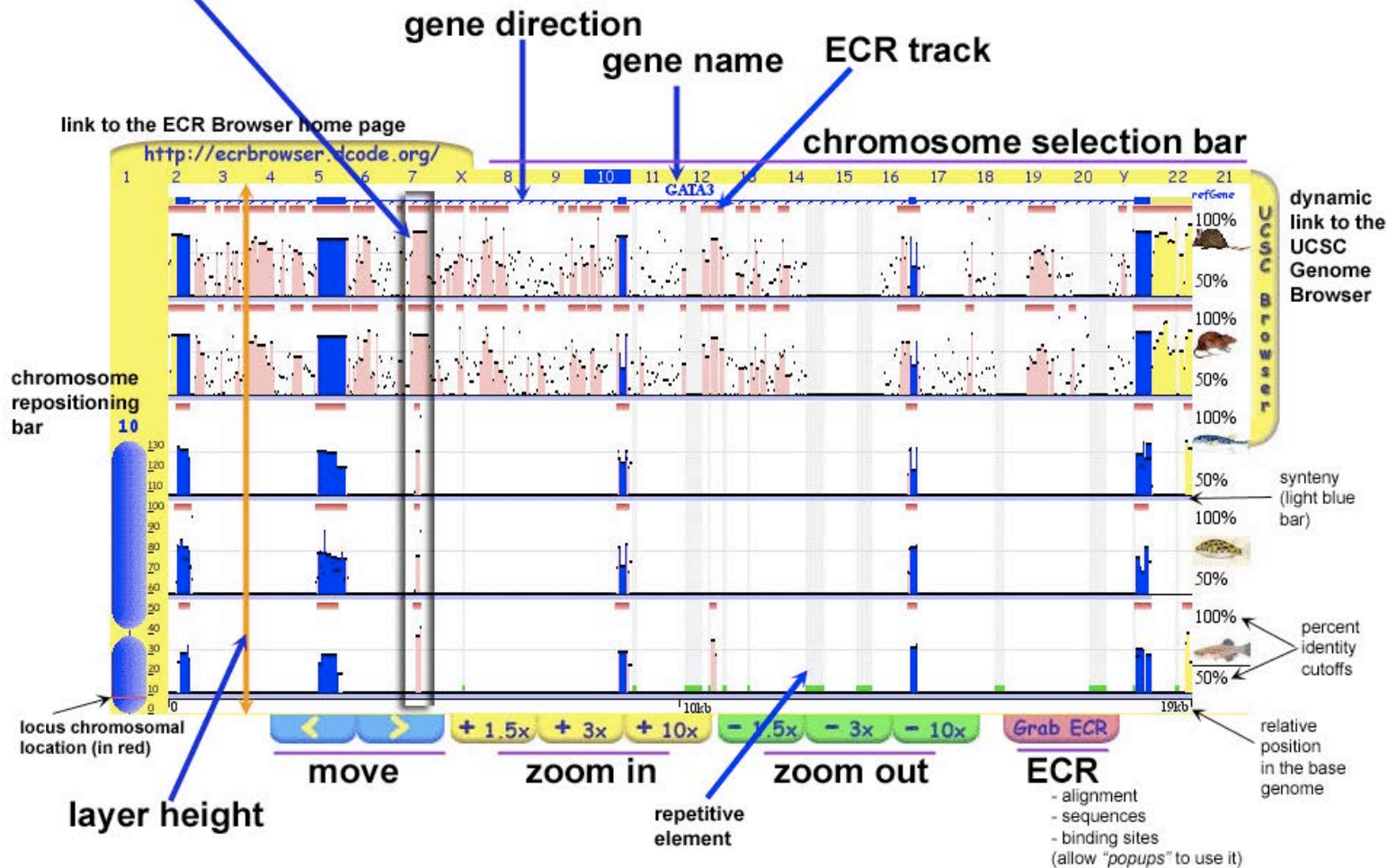


- ◆ Evolutionarily Conserved Regions
  - ◆ Coding Regions
  - ◆ Regulatory Regions
  
- ◆ ECR Browser
  - ◆ <http://ecrbrowser.dcode.org>



# ECR Browser (<http://ecrbrowser.dcode.org/>) legend

**putative regulatory ECR**  
(ECR that is conserved in all vertebrate species)



coding exon	UTR	intron	intergenic element

**Color codes**

**Species**

mouse	rat	fugu	tetraodon	zebrafish

ECR Browser on Human (hg18)

<http://ecrbrowser.dcode.org>

Parameters:

Graph	ECR length	ECR similarity	Layer height	Coordinate system
smooth	100	70	55	relative

[\[change\]](#)

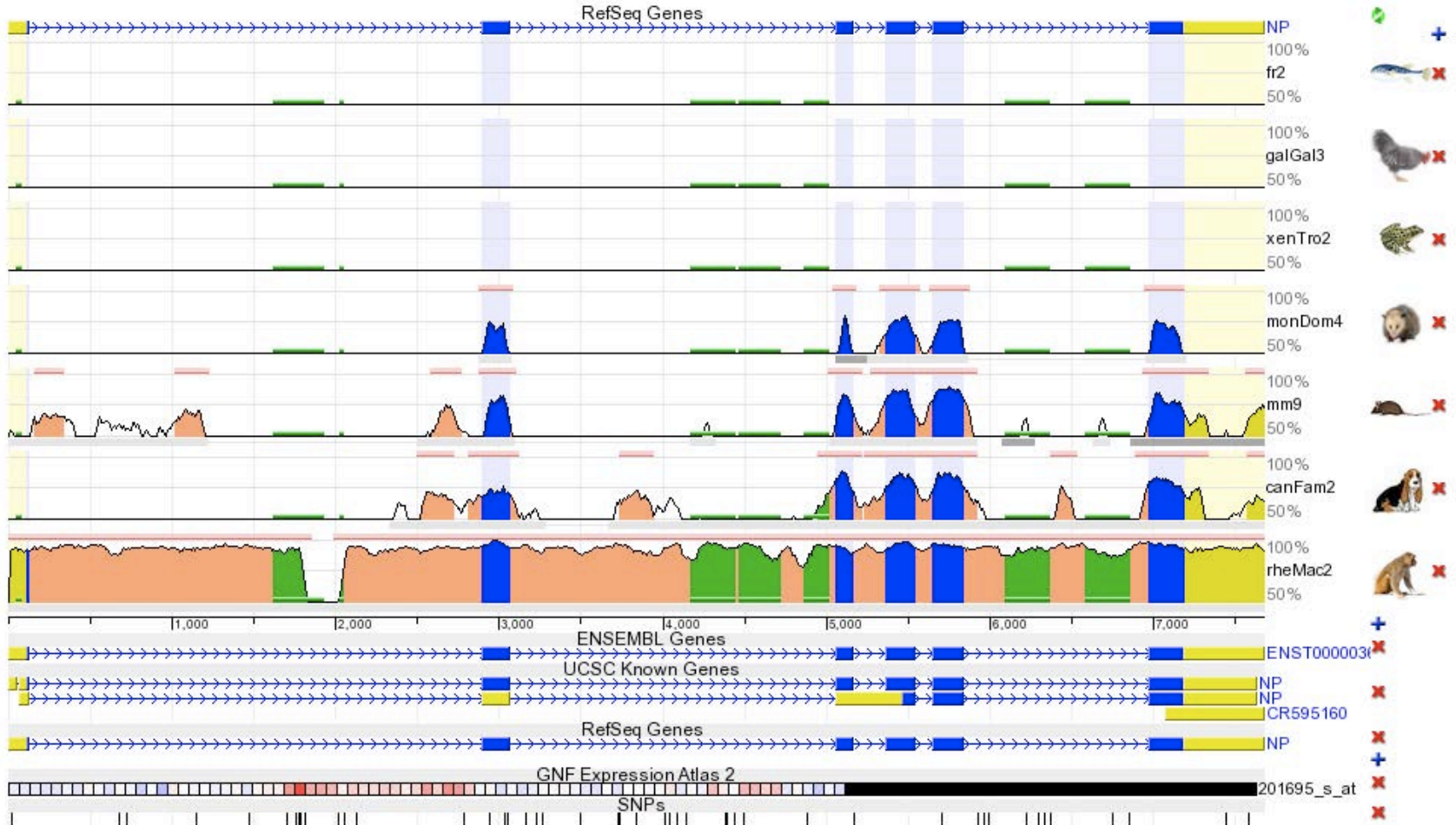
7,674 bps

gene or position (chrN:from-to)

chr14:20007409-20015082

**Submit**

**GENOME ALIGNMENT:** Align your sequence to a genome



# Sequences



# The Sequence Record



- ◆ Different for each database
- ◆ Locus (Name)
- ◆ Accession Number
- ◆ Keywords
- ◆ Description
- ◆ Properties
- ◆ References
- ◆ The Sequence

# GenBank Sample Record

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

**GenBank Sample Record - Mozilla Firefox**

File Edit View Go Bookmarks Tools Help

http://www.ncbi.nlm.nih.gov/Sitemap/samplerecord.html

Go uniprot

**NCBI** Sample GenBank Record

PubMed Entrez BLAST OMIM Taxonomy Structure

### GenBank Flat File Format

Click on any link in this sample record to see a detailed description of that data element or field. All of the descriptions are included on this page, so it can be printed as a single document. You can also return to the [Alphabetical Quicklinks Table](#) or [Resource Guide](#)

**LOCUS** SCU49845 5028 bp DNA PLN 21-JUN-1999

**DEFINITION** Saccharomyces cerevisiae TCP1-beta gene, partial ods, 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

**FEATURES**

source Location/Qualifiers

1..5028

/organism="Saccharomyces cerevisiae"

/db\_xref="taxon:4932"

/chromosome="IX"

/map="9"

CDS

<1..206

/codon\_start=3

/product="TCP1-beta"

/protein\_id="AAA98665.1"

/db\_xref="GI:1293614"

/translation="SSIYNGISTSGLDLNGTIADMRQLGIVESYKLRKRAVSSASEA AEVLLRVDNIIIRARPRTANRQHM"

gene

687..3158

/gene="AXL2"

# GenBank Record Human CFTR

LOCUS HUMCFTRM 6129 bp mRNA linear PRI 27-APR-1993  
DEFINITION Human cystic fibrosis mRNA, encoding a presumed transmembrane conductance regulator (CFTR).  
ACCESSION M28668  
VERSION M28668.1 GI:180331  
KEYWORDS cystic fibrosis; transmembrane conductance regulator.  
SOURCE Homo sapiens (human)  
ORGANISM Homo sapiens  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.  
REFERENCE 1 (bases 1 to 6129)  
AUTHORS Riordan,J.R., Rommens,J.M., Kerem,B., Alon,N., Rozmahel,R., Grzelczak,Z., Zielenski,J., Lok,S., Plavsic,N., Chou,J.-L., Drumm,M.L., Iannuzzi,M.C., Collins,F.S. and Tsui,L.-C.  
TITLE Identification of the cystic fibrosis gene: cloning and characterization of complementary DNA  
JOURNAL Science 245 (4922), 1066-1073 (1989)  
PUBMED 2475911

# Accession Numbers



- ◆ Each sequence submitted to a database is assigned a unique primary accession number
- ◆ Accession numbers do not change
- ◆ If a sequence is merged with another, a new accession number is assigned, and the original number becomes a secondary accession number
- ◆ Accession numbers may include version numbers
  - ◆ AO2428.2

# GI numbers



- ◆ GenInfo Identifier
- ◆ Series of digits that are assigned consecutively to each sequence record processed by NCBI
- ◆ The GI number bears no resemblance to the Accession number of the sequence record
- ◆ Used to track sequence histories in GenBank and the other sequence databases
  - ◆ If a sequence is changed/updated, a new gi number is assigned



COMMENT A three base-pair deletion spanning positions 1654-1656 is observed in cDNAs from cystic fibrosis patients.

FEATURES Location/Qualifiers

source 1. .6129  
/organism="Homo sapiens"  
/db\_xref="taxon:9606"

CDS 133. .4575  
/note="cystic fibrosis transmembrane conductance regulator"  
/codon\_start=1  
/db\_xref="PID:g180332"

/translation="MQRSPLEKASVVSKLFFSWTRPILRKG YRQRLELSDIYQI PSVD  
SADNLSEKLEREWDRELASKKNPKLINALRRCFFWRFMFYGIFLYLGEVTKAVQPLLL  
LNRFSKDIAILDLLPLTIFDFIQLLLIVIGAI AVVAVLQPYIFVATVPVIVAFIMLR  
AYFLQTSQQLKQLESEGRSPIFTHLVTS LKGLWTLRAFGRQPYFETLFHKALNLHTAN  
WFLYLS TLRWFQMRIEMIFVIF FIAVTFIS ILTTGEGEGR VGIIL TLAMNIMSTLQWA  
VNSSIDVDSL MRSVSRVFKFIDMPTEGKPTKSTKPYKNGQLSKVMI IENSHVKKDDIW  
PSGGQMTVKDLTAKYTEGGNAILEN ISFSISPGQRVGLLGRTGSGKSTLLSAFLRLN  
TEGEIQIDGVSWDSITLQQWRKAFGVIPQKVFI FSGTFRKNLDPYEQWSDQE IWKVAD  
EVGLRSVIEQFP GKLDVFLVDGGCVLSHG HKQLMCLARSVLSKAKILLLDEPSAHLDP  
VTYQII RR TLKQAFADCTVILCEHRIEAMLECQQFLVIEENKVRQYDSIQKLLNERSL  
FRQAISPSDRVKL FPHRNSSKCKSKPQIAALKEETEEEVQDTRL"

BASE COUNT 1886 a 1181 c 1330 g 1732 t

ORIGIN

HUMCFTRM Length: 6129 April 13, 1998 13:00 Type: N Check: 6781 ..

```
1  AATTGGAAGC AAATGACATC ACAGCAGGTC AGAGAAAAAG GGTTGAGCGG
51  CAGGCACCCA GAGTAGTAGG TCTTTGGCAT TAGGAGCTTG AGCCCAGACG
101 GCCCTAGCAG GGACCCCAGC GCCCGAGAGA CCATGCAGAG GTCGCCTCTG
151 GAAAAGGCCA GCGTTGTCTC CAAACTTTTT TTCAGCTGGA CCAGACCAAT
201 TTTGAGGAAA GGATACAGAC AGCGCCTGGA ATTGTCAGAC ATATACCAAA
251 TCCCTTCTGT TGATTCTGCT GACAATCTAT CTGAAAAATT GGAAAGAGAA
301 TGGGATAGAG AGCTGGCTTC AAAGAAAAT CCTAAACTCA TTAATGCCCT
351 TCGGCGATGT TTTTTCTGGA GATTTATGTT CTATGGAATC TTTTTATATT
401 TAGGGGAAGT CACCAAAGCA GTACAGCCTC TCTTACTGGG AAGAATCATA
451 GCTTCCTATG ACCCGGATAA CAAGGAGGAA CGCTCTATCG CGATTTATCT
```

analyze% typedata -ref GB\_PR:HUMIFNRF1A

!!NA\_SEQUENCE 1.0

LOCUS HUMIFNRF1A 7721 bp DNA PRI 10-NOV-1992

DEFINITION Homo sapiens interferon regulatory factor 1 gene, complete cds.

ACCESSION L05072

NID g184648

KEYWORDS interferon regulatory factor 1.

SOURCE Homo sapiens Placenta DNA.

ORGANISM Homo sapiens

Eukaryotae; mitochondrial eukaryotes; Metazoa; Chordata; Vertebrata; Eutheria; Primates; Catarrhini; Hominidae; Homo.

REFERENCE 1 (bases 1 to 7721)

AUTHORS Cha,Y., Sims,S.H., Romine,M.F., Kaufmann,M. and Deisseroth,A.B.

TITLE Human interferon regulatory factor 1: intron/exon organization

JOURNAL DNA Cell Biol. 11, 605-611 (1992)

MEDLINE 93000481

```

FEATURES                                 Location/Qualifiers
     source                               1. .7721
                                           /organism="Homo sapiens"
                                           /db_xref="taxon:9606"
                                           /tissue_type="Placenta"
                                           /map="5q23-q31"
     exon                                  1. .219
                                           /gene="IRF1"
                                           /note="putative"
                                           /number=1
     5'UTR                                 join(1. .219,1279. .1287)
                                           /gene="IRF1"
     gene                                   join(1. .219,1279. .1287)
                                           /gene="IRF1"
     intron                                 220. .1278
                                           /gene="IRF1"
                                           /number=1
     exon                                  1279. .1374
                                           /gene="IRF1"
                                           /number=2
     CDS                                    join(1288. .1374,2738. .2837,3630. .3806,3916. .3965,
4073. .4202,4386. .4508,5040. .5089,6248. .6383,6670.
.6794)
                                           /gene="IRF1"
                                           /codon_start=1
                                           /product="interferon regulatory factor 1"
                                           /db_xref="PID:g184649"
                                           /translation="MPI TRMRMRPWLEMQINSNQIPGLIWINK EEMIFQIPWKHAAKH
GWDINKDACLF RSWAIHTGRYKAGEKEPDPKTKANFRCAMNSLPDIEEVKDQSRNKG
SSAVRVYRMLPPLTKNQRKERKSKSRDAKSKAKRKSCGDS SPDTFSDGLSSSTLPDD
HSSYTVPGYMQDLEVEQALTPALSPCAVSS TLPDWHIPVEVVPDSTSDLYNFQVSPMP

```

intron	1375. .2737
	/gene="IRF1"
	/number=2
exon	2738. .2837
	/gene="IRF1"
	/number=3
intron	2838. .3629
	/gene="IRF1"
	/number=3
exon	3630. .3806
	/gene="IRF1"
	/number=4
intron	3807. .3915
	/gene="IRF1"
	/number=4
exon	3916. .3965
	/gene="IRF1"
	/number=5
intron	3966. .4072
	/gene="IRF1"
	/number=5

...

exon	5040. .5089
	/gene="IRF1"
	/number=8
intron	5090. .6247
	/gene="IRF1"
	/number=8
exon	6248. .6383
	/gene="IRF1"
	/number=9
intron	6384. .6669
	/gene="IRF1"
	/number=9
exon	6670. .7656
	/gene="IRF1"
	/number=10
3'UTR	6795. .7656

BASE COUNT 1750 a 1946 c 2253 g 1772 t  
ORIGIN

# Protein Sequence Databases



# UniProt



- ◆ Single site for access to protein sequences
- ◆ <http://www.uniprot.org/>
- ◆ Combines
  - ◆ Swiss-Prot
  - ◆ TrEMBL
  - ◆ PIR



UniProtKB

Advanced



BLAST Align Retrieve/ID Mapping

Help Contact

The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information.

**UniProtKB**

Swiss-Prot (547,357)  
Manually annotated and reviewed.

TrEMBL (89,451,166)  
Automatically annotated and not reviewed.

**UniRef**  
Sequence clusters

**UniParc**  
Sequence archive

**Proteomes**

**Supporting data**

Literature citations	Taxonomy	Subcellular locations
Cross-ref. databases	Diseases	Keywords

**News**

Thalidomide, the pharmacological version of yin and yang | Cross-references to DEPOD, MoonProt and Proteomes  
[UniProt release 2015\\_01](#)

Higher and higher | New mouse and zebrafish variation files | Structuring of 'cofactor' annotations  
[UniProt release 2014\\_11](#)

[News archive](#)

### Getting started

- Text search**  
Our basic text search allows you to search all the resources available
- BLAST**  
Find regions of similarity between your sequences
- Sequence alignments**  
Align two or more protein sequences using the Clustal Omega program

### UniProt data

- Download latest release**  
Get the UniProt data
- Statistics**  
View Swiss-Prot and TrEMBL statistics
- Forthcoming changes**  
Planned changes for the UniProt knowledgebase
- Submit your data**  
Submit your sequences and annotation

### Protein spotlight

**The Hidden Things**  
December 2014

Nature has its secret ways. During the course of the 19th century, the Augustinian friar Gregor Mendel worked out the basics of genetic inheritance as he crossbred pea plants. About a century later, it has become obvious that the inheritance of a given trait is in fact not so straightforward...





# Databases of Biological Knowledge



Categorizing, Describing, and Referencing Biological  
Information

# Biological Ontologies



- ◆ Dictionary of terms describing biological processes
- ◆ Uses a “controlled vocabulary” to precisely define the terms to be used along with their definitions
  - ◆ Also provides a table of synonyms
- ◆ Provides the ability to construct reliable systems for computer-searchable databases and comparative analyses

# A Few Ontologies



- ◆ Sequence Ontology (SO)
  - ◆ Ontology suitable for describing biological sequences
    - ◆ Sequence features with coordinates
      - ◆ exon, promoter, binding site...
    - ◆ Properties of features
  - ◆ <http://song.sourceforge.net/>
  
- ◆ Gene Ontology (GO)
  - ◆ Provides a controlled vocabulary to describe gene and gene product attributes in any organism
    - ◆ biological processes
    - ◆ cellular components
    - ◆ molecular functions
  - ◆ [www.geneontology.org](http://www.geneontology.org)



Quick search

Search

# Search Ontology

## Information about Ontology search

Free-text filtering

X

Your search is pinned to these filters

+ document\_category: ontology\_class

No current user filters.

### Ontology source

biological_process	(27810)	+	-
mouse_ontology	(19446)	+	-
molecular_function	(10805)	+	-
ncbi_taxonomy	(6146)	+	-
cellular_component	(3846)	+	-
adult_mouse_anatomy.gxd	(3229)	+	-
cell	(2128)	+	-
eco	(583)	+	-
go/extensions/gorel	(79)	+	-
cl	(22)	+	-
external	(6)	+	-

Subset

Ancestor

## Found entities

Total: 74298; showing 1-10

Results count 10

Navigation buttons: first, previous, next, last

Refresh and Filter buttons

<input type="checkbox"/>	Term	Definition	Ontology source	Synonyms
<input type="checkbox"/>	<a href="#">nucleic acid binding evidence</a>		eco	
<input type="checkbox"/>	<a href="#">transport assay evidence</a>		eco	IDA: transport assay
<input type="checkbox"/>	<a href="#">cobinamide phosphate guanylyltransferase activity</a>	Catalysis of the reaction: adenosylcobinamide phosphate + GTP + 2 H(+) = adenosylcobinamide-GDP + diphosphate.	molecular_function	GTP:adenosylcobinamide-phosphate guanylyltransferase activity GTP:cobinamide phosphate guanylyltransferase activity <a href="#">more...</a>
<input type="checkbox"/>	<a href="#">crossover junction endodeoxyribonuclease activity</a>	Catalysis of the endonucleolytic cleavage at a junction such as a reciprocal single-stranded crossover <a href="#">more...</a>	molecular_function	Hje endonuclease activity Holliday junction nuclease activity <a href="#">more...</a>



# Searching for Information



# Entrez Searching



- ◆ <http://www.ncbi.nlm.nih.gov/entrez/>
- ◆ Search via text patterns
- ◆ Cross-database search interface
  - ◆ Sequence
  - ◆ PubMed
  - ◆ OMIM
  - ◆ Linkage information
  - ◆ ...

## Search NCBI databases

Help

human cfr

Search

### Results found in 35 databases for "human cfr"

#### Literature

<b>Books</b>	141	books and reports
<b>MeSH</b>	14	ontology used for PubMed indexing
<b>NLM Catalog</b>	3	books, journals and more in the NLM Collections
<b>PubMed</b>	6,947	scientific & medical abstracts/citations
<b>PubMed Central</b>	11,510	full-text journal articles

#### Health

<b>ClinVar</b>	1,290	human variations of clinical significance
<b>dbGaP</b>	51	genotype/phenotype interaction studies
<b>GTR</b>	97	genetic testing registry
<b>MedGen</b>	1	medical genetics literature and links
<b>OMIM</b>	81	online mendelian inheritance in man
<b>PubMed Health</b>	15	clinical effectiveness, disease and drug reports

#### Genomes

<b>Assembly</b>	0	genome assembly information
<b>BioProject</b>	31	biological projects providing data to NCBI
<b>BioSample</b>	124	descriptions of biological source materials
<b>Clone</b>	1,362	genomic and cDNA clones
<b>dbVar</b>	2,605	genome structural variation studies
<b>Epigenomics</b>	5,161	epigenomic studies and display tools
<b>Genome</b>	2	genome sequencing projects by organism
<b>GSS</b>	3	genome survey sequences
<b>Nucleotide</b>	6,572	DNA and RNA sequences
<b>Probe</b>	1,368	sequence-based probes and primers
<b>SNP</b>	10,723	short genetic variations
<b>SRA</b>	105	high-throughput DNA and RNA sequence read archive
<b>Taxonomy</b>	0	taxonomic classification and nomenclature catalog

#### Genes

<b>EST</b>	3	expressed sequence tag sequences
<b>Gene</b>	2,233	collected information about gene loci
<b>GEO DataSets</b>	186	functional genomics studies
<b>GEO Profiles</b>	234,414	gene expression and molecular abundance profiles
<b>HomoloGene</b>	13	homologous gene sets for selected organisms
<b>PopSet</b>	5	sequence sets from phylogenetic and population studies
<b>UniGene</b>	9	clusters of expressed transcripts

#### Proteins

<b>Conserved Domains</b>	6	conserved protein domains
<b>Protein</b>	1,075	protein sequences
<b>Protein Clusters</b>	0	sequence similarity-based protein clusters
<b>Structure</b>	64	experimentally-determined biomolecular structures

#### Chemicals

<b>BioSystems</b>	865	molecular pathways with links to genes, proteins and chemicals
<b>PubChem BioAssay</b>	204	bioactivity screening studies
<b>PubChem Compound</b>	0	chemical information with structures, information and links
<b>PubChem Substance</b>	33	deposited substance and chemical information

# NCBI Gene



- ◆ Provides a single query interface to curated sequence and descriptive information about genetic loci
  - ◆ Nomenclature
  - ◆ Aliases
  - ◆ Sequence accessions
  - ◆ Phenotypes
  - ◆ EC numbers
  - ◆ OMIM numbers
  - ◆ UniGene clusters
  - ◆ Homology
  - ◆ Map locations
  - ◆ Web sites



Gene   [Advanced](#)

Display Settings:  Full Report Send to:

## CFTR cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7) [ *Homo sapiens* (human) ]

Gene ID: 1080, updated on 1-Feb-2015

**Summary**

**Official Symbol** CFTR provided by HGNC

**Official Full Name** cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7) provided by HGNC

**Primary source** [HGNC:HGNC:1884](#)

**See related** [Ensembl:ENSG000000001626](#); [HPRD:03883](#); [MIM:602421](#); [Vega:OTTHUMG000000023076](#)

**Gene type** protein coding

**RefSeq status** REVIEWED

**Organism** [Homo sapiens](#)

**Lineage** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorhini; Catarrhini; Hominidae; Homo

**Also known as** CF; MRP7; ABC35; ABCC7; CFTR/MRP; TNR-CFTR; dJ760C5.1

**Summary** This gene encodes a member of the ATP-binding cassette (ABC) transporter superfamily. ABC proteins transport various molecules across extra- and intra-cellular membranes. ABC genes are divided into seven distinct subfamilies (ABC1, MDR/TAP, MRP, ALD, OABP, GCN20, White). This protein is a member of the MRP subfamily that is involved in multi-drug resistance. The encoded protein functions as a chloride channel and controls the regulation of other transport pathways. Mutations in this gene are associated with the autosomal recessive disorders cystic fibrosis and congenital bilateral aplasia of the vas deferens. Alternatively spliced transcript variants have been described, many of which result from mutations in this gene. [provided by RefSeq, Jul 2008]

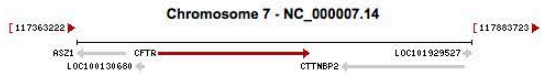
**Orthologs** [mouse](#) [all](#)

**Genomic context**

**Location:** 7q31.2 See CFTR in [MapViewer](#)

**Exon count:** 31

Annotation release	Status	Assembly	Chr	Location
106	current	GRCh38 ( <a href="#">GCF_000001405.26</a> )	7	NC_000007.14 (117470772..117668665)
105	previous assembly	GRCh37.p13 ( <a href="#">GCF_000001405.25</a> )	7	NC_000007.13 (117120017..117308719)



**Genomic regions, transcripts, and products**

**Genomic Sequence:**

Go to [reference sequence details](#) | [Graphics](#) | [FASTA](#) | [GenBank](#)

NC\_000007.14: 117M..118M (257Kbp) Find:  Tools

17,440 K | 117,460 K | 117,480 K | 117,500 K | 117,520 K | 117,540 K | 117,560 K | 117,580 K | 117,600 K | 117,620 K | 117,640 K | 117,660 K | 117,680 K

Genes, NCBI Homo sapiens Annotation Release 106

2 X1\_006715842.1 | N1\_000492.3 | XP\_006715905.1 | NP\_000483.3

1 LOC100130680

CFTR

CCDS Features, Release 17 (NCBI Annotation Release 106 compared to Ensembl Release 76)

CCDS5773.1

**Table of contents**

- Summary
- Genomic context
- Genomic regions, transcripts, and products
- Bibliography
- Phenotypes
- Variation
- Pathways from BioSystems
- Interactions
- General gene information
  - Markers, Related pseudogene(s), Homology, Gene Ontology
- General protein information
- NCBI Reference Sequences (RefSeq)
- Related sequences
- Additional links
  - Locus-specific Databases

- Related information**
- Order cDNA clone
  - 3D structures
  - BioAssay
  - BioAssay by Target (List)
  - BioAssay by Target (Summary)
  - BioAssay, by Gene target
  - BioAssays, RNAi Target, Active
  - BioAssays, RNAi Target, Tested
  - BioProjects
  - BioSystems
  - Books
  - CCDS
  - ClinVar
  - Conserved Domains
  - dbVar
  - Full text in PMC
  - Full text in PMC\_nucleotide
  - GAP
  - Gene neighbors
  - Genome
  - GEO Profiles
  - GTR
  - HomoloGene
  - Map Viewer
  - MedGen
  - Nucleotide
  - OMIM
  - Probe

# OMIM



- ◆ Online Mendelian Inheritance in Man
- ◆ Database of gene-linked genetic disorders
- ◆ Maintained at Johns Hopkins University
  - ◆ Dr. Victor A. McKusick
- ◆ Provides link to GeneTests
  - ◆ Laboratories that provide testing for specific genetic disorders

[Advanced Search](#)
**\*602421**

## CYSTIC FIBROSIS TRANSMEMBRANE CONDUCTANCE REGULATOR; CFTR

*Alternative titles; symbols*

ATP-BINDING CASSETTE, SUBFAMILY C, MEMBER 7; ABCC7

*HGNC Approved Gene Symbol: CFTR*
*Cytogenetic location: 7q31.2    Genomic coordinates (GRCh37): 7:117,120,016-117,308,718* (from NCBI)

### Gene-Phenotype Relationships

Location	Phenotype	Phenotype MIM number	Phenotype mapping key
7q31.2	Congenital bilateral absence of vas deferens	277180	3
	Cystic fibrosis	219700	3
	Sweat chloride elevation without CF		3
	{Bronchiectasis with or without elevated sweat chloride 1, modifier of}	211400	3
	{Hypertrypsinemia, neonatal}		3
	{Pancreatitis, idiopathic}	167800	3

### TEXT

#### Description

The CFTR gene encodes an ATP-binding cassette (ABC) transporter that functions as a low conductance Cl(-)-selective channel gated by cycles of ATP binding and hydrolysis at its nucleotide-binding domains (NBDs) and regulated tightly by an intrinsically disordered protein segment distinguished by multiple consensus phosphorylation sites termed the regulatory domain (summary by [Wang et al., 2014](#)).

#### Cloning and Expression

[Riordan et al. \(1989\)](#) isolated overlapping cDNA clones from epithelial cell libraries with a genomic DNA segment containing a portion of the putative gene causing cystic fibrosis (CF; 219700). Transcripts approximately 6,500 nucleotides in size were detectable in the tissues affected in patients with CF. The predicted protein consists of 2 similar motifs, each with a domain having properties consistent with membrane-association, and a domain believed to be involved in ATP binding. In CF patients, a deleted phenylalanine residue occurs at the center of the putative first nucleotide-binding fold (NBF). The predicted protein has 1,480 amino acids with a molecular mass of 168,138 Da. The characteristics are remarkably similar to those of the mammalian multidrug resistant P-glycoprotein (171050), which also maps to 7q, and to a number of other membrane-associated proteins. To avoid confusion with the previously named CF antigen (123885), [Riordan et al. \(1989\)](#) referred to the protein as cystic fibrosis transmembrane conductance regulator (CFTR).

Cystic fibrosis represents the first genetic disorder elucidated strictly by the process of reverse genetics (later called positional cloning), i.e., on the basis of map location but without the availability of chromosomal rearrangements or deletions such as those that have greatly facilitated previous success in the cloning of human disease genes in Duchenne muscular dystrophy (310200), retinoblastoma (180200), and chronic granulomatous disease (306400), for example. By use of a combination of chromosome walking and jumping, [Rommens et al. \(1989\)](#) succeeded in covering the CF region on 7q. The jumping technique was particularly useful in bypassing 'unclonable' regions, which are estimated to constitute 5% of the human genome. (Yeast artificial chromosome (YAC) vectors represent an alternative strategy.) The identification of undermethylated CpG islands was 1 tip-off; another was screening of a cDNA library constructed from cultured sweat gland cells of a non-CF individual. The CF gene proved to be about 250,000 bp long, a surprising finding since the absence of apparent genomic rearrangements in CF chromosomes and the evidence of a limited number of CF mutations predicted a small mutational target.

[Green and Olson \(1990\)](#) described a general strategy for cloning and mapping large regions of human DNA with yeast artificial chromosomes (YACs). By analyzing 30 YAC clones from the region of chromosome 7 containing the CFTR gene, a contig map spanning more than 1.5 Mbp was assembled. Individual YACs as large as 790 kb

#### Table of Contents for \*602421

- Title
- Gene-Phenotype Relationships
- Text
  - Description
  - Cloning and Expression
  - Gene Structure
  - Mapping
  - Gene Function
  - Biochemical Features
  - Molecular Genetics
  - Animal Model
  - History
- Allelic Variants
- Table View
- See Also
- References
- Contributors
- Creation Date
- Edit History

#### External Links for Entry:

- ▶ [Genome](#)
- ▶ [DNA](#)
- ▶ [Protein](#)
- ▶ [Gene Info](#)
- ▶ [Clinical Resources](#)
- ▶ [Variation](#)
- ▶ [Animal Models](#)
- ▶ [Cellular Pathways](#)



# Comparative Analysis



Similarities and Differences

# Comparative Analysis



- ◆ Identification of similarities
  - ◆ Primary sequence
  - ◆ Structure
  - ◆ Function
  
- ◆ Identification of differences
  - ◆ Gene complement
  - ◆ Genotypic differences resulting in phenotypic changes
  
- ◆ Phylogenetic inference
  - ◆ Predicting evolutionary history

# Finding Sequences by Similarity



# Homology vs. Similarity



- ◆ Similarity
  - ◆ Arises from:
    - ◆ Homology
    - ◆ Convergence
    - ◆ Gene Conversion
    - ◆ Chance
  
- ◆ Homology
  - ◆ Implies a common evolutionary origin
    - ◆ Homologs
    - ◆ Orthologs
    - ◆ Paralogs

# BLAST




- ◆ Basic Local Alignment Search Tool
- ◆ Search a sequence database for primary sequence similarities to some query sequence
- ◆ Provides a measure of the significance of the similarity
- ◆ Does not necessarily imply common evolutionary origin



# BLAST



- ◆ All search combinations possible
- ◆ nt vs. nt database
  - ◆ blastn
- ◆ protein vs. protein database
  - ◆ blastp
- ◆ translated nt vs. protein database
  - ◆ blastx
- ◆ protein vs. translated nt database
  - ◆ tblastn
- ◆ translated nt vs. translated nt database
  - ◆ tblastx



# BLAST Searching at NCBI





BLAST®

Basic Local Alignment Search Tool

[Home](#)[Recent Results](#)[Saved Strategies](#)[Help](#)[My NCBI](#)[\[Sign In\]](#) [\[Register\]](#)▶ [NCBI/ BLAST Home](#)BLAST finds regions of similarity between biological sequences. [more...](#)**New** Try [SmartBLAST](#) for an improved protein-protein search

## BLAST Assembled Genomes

Find Genomic BLAST pages:

**GO**
 [Human](#)  
 [Mouse](#)  
 [Rat](#)  
 [Cow](#)  
 [Pig](#)  
 [Dog](#)
 [Rabbit](#)  
 [Chimp](#)  
 [Guinea pig](#)  
 [Fruit fly](#)  
 [Honey bee](#)  
 [Chicken](#)
 [Zebrafish](#)  
 [Clawed frog](#)  
 [Arabidopsis](#)  
 [Rice](#)  
 [Yeast](#)  
 [Microbes](#)

## Basic BLAST

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a <b>nucleotide</b> database using a <b>nucleotide</b> query <i>Algorithms: blastn, megablast, discontinuous megablast</i>
<a href="#">protein blast</a>	Search <b>protein</b> database using a <b>protein</b> query <i>Algorithms: blastp, psi-blast, phi-blast, delta-blast</i>
<a href="#">blastx</a>	Search <b>protein</b> database using a <b>translated nucleotide</b> query
<a href="#">tblastn</a>	Search <b>translated nucleotide</b> database using a <b>protein</b> query
<a href="#">tblastx</a>	Search <b>translated nucleotide</b> database using a <b>translated nucleotide</b> query

## Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

- Get faster protein results with a graphical view using [SmartBLAST](#)
- Make specific primers with [Primer-BLAST](#)
- Cluster multiple sequences together with their database neighbors using [MOLE-BLAST](#)
- Find [conserved domains](#) in your sequence (ods)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IgBLAST)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align](#) two (or more) sequences using BLAST (bl2seq)
- Search [protein](#) or [nucleotide](#) targets in PubChem BioAssay
- Search [SRA by experiment](#)
- Constraint Based Protein [Multiple Alignment Tool](#)
- Needleman-Wunsch [Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)
- Search [trace archives](#)
- Search bacterial and fungal rRNA sequences with [Targeted Loci BLAST](#)

### Your Recent Results **New!**

[All Recent results...](#)

### News

#### [Searching Whole Genome Shotgun sequences](#)

It is now much easier to search WGS (Whole Genome Shotgun) with stand-alone BLAST on your own computer.

Wed, 20 Jan 2016 10:00:00 EST

[More BLAST news...](#)

### Tip of the Day

#### [Use Genomic BLAST to see the genomic context](#)

If you are interested in the evolution of a particular gene or gene family it is often interesting to examine the intro-exon structure even across species.

[More tips...](#)



BLAST®

Basic Local Alignment Search Tool

Home

Recent Results

Saved Strategies

Help

My NCBI

Welcome elliotl. [Sign Out]

NCBI/ BLAST/ blastp suite

## Standard Protein BLAST

blastn blastp blastx tblastn tblastx

## Enter Query Sequence

BLASTP programs search protein databases using a protein query. more...

[Reset page](#) [Bookmark](#)

Enter accession number(s), gi(s), or FASTA sequence(s)

[Clear](#)

Query subrange

From

To

Or, upload file

Choose File no file selected

Job Title

Enter a descriptive title for your BLAST search

 Align two or more sequences

## Choose Search Set

Database

UniProtKB/Swiss-Prot(swissprot)

Organism

Optional

Enter organism name or id—completions will be suggested

 Exclude

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude

Optional

 Models (XM/XP)  Uncultured/environmental sample sequences

Entrez Query

Optional

Enter an Entrez query to limit search

[Create custom database](#)

## Program Selection

Algorithm

- blastp (protein-protein BLAST)
- PSI-BLAST (Position-Specific Iterated BLAST)
- PHI-BLAST (Pattern Hit Initiated BLAST)
- DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

Choose a BLAST algorithm

**BLAST**

Search database UniProtKB/Swiss-Prot(swissprot) using Blastp (protein-protein BLAST)

 Show results in a new window

Algorithm parameters

Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

## General Parameters

Max target sequences

♦ 5000

Select the maximum number of aligned sequences to display

Short queries

 Automatically adjust parameters for short input sequences

Expect threshold

10

Word size

3

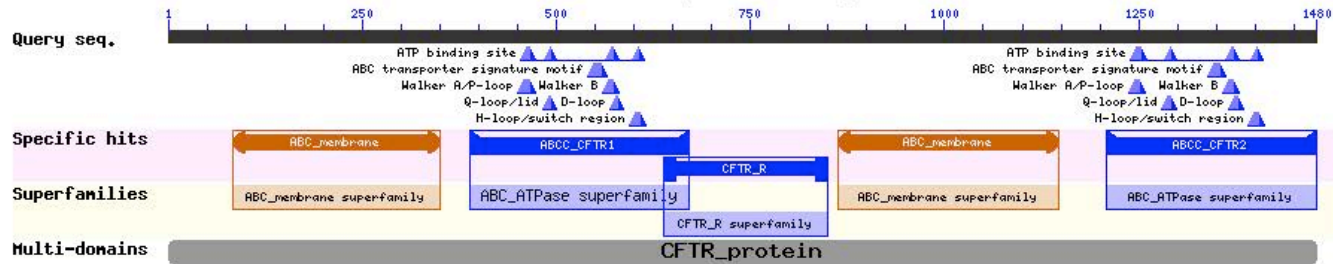
Max matches in a query range

0

NCBI/ BLAST/ blast suite/ Formatting Results - CZARABWB014 [\[Formatting options\]](#)

Job Title: gi|90421313|ref|NP\_000483.3| cystic fibrosis...

Putative conserved domains have been detected, click on the image below for detailed results.



Request ID	CZARABWB014
Status	Searching
Submitted at	Mon Feb 2 11:50:50 2015
Current time	Mon Feb 2 11:51:05 2015
Time since submission	00:00:15

This page will be automatically updated in 1 seconds

BLAST is a registered trademark of the National Library of Medicine.

NCBI/ BLAST/ blastp suite/ Formatting Results - CZARABWB014

Edit and Resubmit Save Search Strategies Formatting options Download

YouTube How to read this page Blast report description

g|90421313|ref|NP\_000483.3| cystic fibrosis...

**RID** CZARABWB014 (Expires on 02-03 23:50 pm)

**Query ID** |c|38978

**Description** g|90421313|ref|NP\_000483.3| cystic fibrosis transmembrane conductance regulator [Homo sapiens]

**Molecule type** amino acid

**Query Length** 1480

**Database Name** swissprot

**Description** Non-redundant UniProtKB/SwissProt sequences

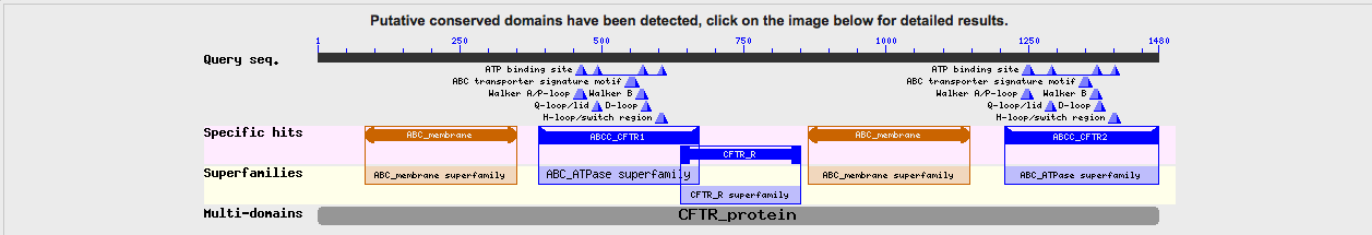
**Program** BLASTP 2.2.30+ Citation

Other reports: Search Summary Taxonomy reports Distance tree of results Multiple alignment

DELTA-BLAST, a more sensitive protein-protein search

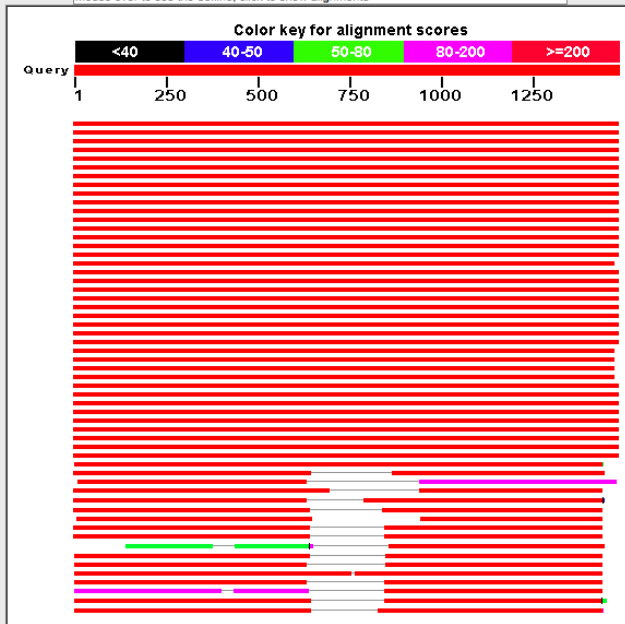
Graphic Summary

Show Conserved Domains



Distribution of 200 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



## Sequences producing significant alignments:

Select: All None Selected: 0

Alignments Download GenPept Graphics Distance tree of results Multiple alignment

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	3046	3046	100%	0.0	100%	<a href="#">P13569.3</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	3034	3034	100%	0.0	99%	<a href="#">Q2IBF6.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	3033	3033	100%	0.0	99%	<a href="#">Q2QLF5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	3017	3017	100%	0.0	99%	<a href="#">Q2IBE4.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	3016	3016	100%	0.0	99%	<a href="#">Q07DX5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2994	2994	100%	0.0	98%	<a href="#">Q7JII8.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2994	2994	100%	0.0	98%	<a href="#">Q9TSP5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2992	2992	100%	0.0	98%	<a href="#">Q9TUQ2.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2991	2991	100%	0.0	98%	<a href="#">Q2IBA1.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2972	2972	100%	0.0	98%	<a href="#">Q07DY5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2956	2956	100%	0.0	97%	<a href="#">Q2QLF9.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2954	2954	100%	0.0	97%	<a href="#">Q07DV2.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2953	2953	100%	0.0	97%	<a href="#">Q2QLB4.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2947	2947	100%	0.0	97%	<a href="#">Q09YH0.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2945	2945	100%	0.0	96%	<a href="#">Q09YK5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2796	2796	100%	0.0	94%	<a href="#">Q2QL83.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2783	2783	99%	0.0	94%	<a href="#">Q2QLA3.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2762	2762	100%	0.0	92%	<a href="#">Q6PQ22.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2760	2760	100%	0.0	92%	<a href="#">Q00554.4</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2734	2734	100%	0.0	93%	<a href="#">Q2QLH0.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2727	2727	100%	0.0	91%	<a href="#">Q07E42.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2703	2703	100%	0.0	92%	<a href="#">Q00PJ2.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2699	2699	100%	0.0	92%	<a href="#">Q07E16.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2694	2694	100%	0.0	90%	<a href="#">Q5U820.2</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2679	2679	100%	0.0	91%	<a href="#">Q10BU0.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2674	2674	100%	0.0	91%	<a href="#">Q00552.2</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2669	2669	99%	0.0	91%	<a href="#">Q07DW5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2667	2667	99%	0.0	91%	<a href="#">Q09YJ4.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2652	2652	99%	0.0	91%	<a href="#">Q00555.2</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2652	2652	99%	0.0	91%	<a href="#">P35071.2</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2649	2649	100%	0.0	90%	<a href="#">Q2IBB3.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2615	2615	100%	0.0	89%	<a href="#">Q2QLC5.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2580	2580	100%	0.0	87%	<a href="#">Q2QL74.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2551	2551	100%	0.0	86%	<a href="#">Q5DIZ7.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2494	2494	100%	0.0	83%	<a href="#">Q07DZ6.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2340	2340	100%	0.0	78%	<a href="#">P26361.2</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2291	2291	100%	0.0	77%	<a href="#">P26363.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2281	2281	100%	0.0	78%	<a href="#">P34158.3</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel</a>	2129	2129	100%	0.0	71%	<a href="#">P26362.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Multidrug resistance-associated protein 1; AltName: Full=ATP-binding cassette sub-family C member 1; AltName: Full=Leukotriene C(4) transporter; Short=MRP1; AltName: Full=MRP1</a>	488	565	96%	1e-144	26%	<a href="#">Q5F364.1</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Multidrug resistance-associated protein 4; AltName: Full=ATP-binding cassette sub-family C member 4; AltName: Full=MRP/cMOAT-related ABC transporter; Short=MRP4; AltName: Full=MRP4</a>	421	971	82%	3e-122	34%	<a href="#">O15439.3</a>
<input type="checkbox"/>	<a href="#">RecName: Full=Probable multidrug resistance-associated protein lethal(2)03659; AltName: Full=Wunen region A protein [Drosophila melanogaster]</a>	347	843	78%	8e-98	32%	<a href="#">P91660.3</a>
<input type="checkbox"/>	<a href="#">RecName: Full=ATP-binding cassette transporter abc2; Short=ABC transporter abc2; AltName: Full=ATP-energized glutathione S-conjugate pump abc2; AltName: Full=Glutathione S-conjugate pump abc2</a>	308	663	80%	2e-84	34%	<a href="#">Q10185.1</a>

RecName: Full=Cystic fibrosis transmembrane conductance regulator; Short=CFTR; AltName: Full=ATP-binding cassette sub-family C member 7; AltName: Full=Channel conductance-controlling ATPase; AltName: Full=Dogfish transmembrane conductance regulator; AltName: Full=cAMP-dependent chloride channel [Squalus acanthias]

Sequence ID: sp|P26362.1|CFTR\_SQUAC Length: 1492 Number of Matches: 1

Related Information

Range 1: 1 to 1492 GenPept Graphics Next Match Previous Match

Score Expect Method Identities Positives Gaps 2129 bits(5516) 0.0 Compositional matrix adjust. 1069/1497(71%) 1251/1497(83%) 22/1497(1%)

Query 1 MQRSPLEKASVSVKLFSSWTRPILRKKYRQRLELSDIYQIPSVDSADNLSEKLEREWDR 60
Sbjct 1 MQRSP+EKA+ SKLFF W RPIL+KGYRQ+LELSDIYQIPS DSAD LSE LEREWDR 60
MQRSPTEKANAFSKLFRWRPILRKKYRQKLELSDIYQIPSSDASELSEMLEREWDR

Query 61 LA-SKKNPKLINALRRCFFWRPFYGIPLYGEVTKAVQPLLGRIIASYPDNKEERSI 119
Sbjct 61 LA SKKNPKL+NALRRCFFWRPFYGI LY E TRAVOPL LGRIIAS Y+ N ER I 120
LATS KKNPKL+NALRRCFFWRPFLYGI LLYVFEVTKAVQPLCLGRIIAS YNAKNTYEREI

Query 120 AIYVLGIGLCLLFIVRTLLHFAIFGLHHIGMQMRIAIFSLIYKTKLSSRVLDKISIGO 179
Sbjct 121 AIYALGIGLCLLFIVRTLLHFAVFGLOHLGQMRIALFSLIYKTKLSSRVLDKIDTQ 180

Query 180 LVSLLSNLLKDFDEGLALAHFVWTAPOVALLMGLIWEELQASAFCGLFLIVLALFQAG 239
Sbjct 181 LVSLLSNLLKDFDEG+A-AHFVWTA+QV LLMGLI W L FCGLFLI+LALFQA 240
LVSLLSNLLKDFDEGVAHFWTA+VAVQVLLMGLIWEELTEFVFCGLGFLIMLALFQAW

Query 240 LGRMMKYRDRQAGKISERLVTSEMENIQSVKAYCWEAMEKMIENLQTEKLRKA 299
Sbjct 241 LG+ MM+YRD+RAGKI+ERL ITSE+I+NIQSVK YCWE+AMEK+I+++RQ ELKLRK 300
LGKMMQYRDRKAGKINERLAIITSEI+I+NIQSVKAYCWEAMEKI+IDDI+RQVELKLRKV

Query 300 AYVRYFNSSAFFSGFFVFLSVLPYALIKGIILRKKIPTTISFCIVLRMAVTRQFPAVQ 359
Sbjct 301 AYCRYFNSSAFFSGFFVFLSVVYAFIHTIKLRRIPTTISYNI+VLRMTVTRQFPAIQ 360

Query 360 TWYDSLGAINKIQDFLQKQEKTYLEYNLTTEVVMENVTAFWEEGFELFEKAKQNNNR 419
Sbjct 361 TWYDSLGAIKIQDFL R E+KT+EYNLT E V M NVTA W+EG GELFEK KQK++ R 420
TWYDSLGAIRKIQDFLHKEHRTVEYNLTTEVVMENVTASWDEGIGELFEKVKQNSER

Query 420 KTSNGDLSLFFSNFSLGTPVLKINFKIERGQLLAVAGSTGACKTSLMLVIMGELEPSE 479
Sbjct 421 K+NGDD LFFSNFSL TPVLK+I+FK+E+G+LLA+AGSTG+GK+SILM+IMGELEPS+ 480
KMANGDDGLFFSNFSLHVTPLKINISFKLEKGLLAIAGSTGSGKSSLLMIMGELEPESD

Query 480 GKIKHSGRISFCQFSWIMPGTIKENIIFGVSYDEYRYSVIKACQLEEDISKFAEKDNI 539
Sbjct 481 GKIKHSGRIS+ Q WIMPGTIK+NIIFG+SYDEYRY SV+ ACQLEEDI+ F KD 540
GKIKHSGRISYSFQVWIMPGTIKDNIFGLSYDEYRYSVNVACQLEEDITVFNKDKT

Query 540 VLGGGITLGGQQRARISLARAVYKADLYLSDSFPFGLDVLTEKIFESCVCVKLMAKNT 599
Sbjct 541 VLGGGITLGGQQRARISLARALYKADLYLSDSFPFGLDVTTEKIDIFESCCLKLMVNT 600

Query 600 RILVTSKMEHLKADKILLHEGSSYFYGTSELQNLQPDFSSKLMGCDSDPQFSAERN 659
Sbjct 601 RILVTSK+EHLKADKILL+HEG YFYGTSELO +PDFSS+L+G FD FSAERN 660
RILVTSKLEHLKADKILLHEGHCYFYGTSELQKQPDFSSQLLGSVHFDPSAERN

Query 660 SILTETLHRFSL---EGDAPVSWTETKQSFQ-TGEFGEKRKNS-ILNPINSIRKFSIV 714
Sbjct 661 SILTET R S+ +G S++ET+K SFQ EF EKRR+S I+NPI S +KFS+V 720
SILTETFRRCSSVSSGDGAGLGSYSETRKASFQPPPEFNEKRKSLVNPITSNKKFSLV

Query 715 QK---TFLQNGIIEEDSDEPLERRLSLVDPSEQEAAILPRISVISTGPTLQARRRQSVLNL 772
Sbjct 721 Q+ + Q NG+E+ + EF ER SL+P++E GE PR ++ + QA RROSV L 780
QTAMSYFQTNGMEDATSEFGERHFLIPENELGEPKPRSNIFKSELFPQARRRQSVLAL

Query 773 MTHSVNQGQNIHRKTTASTRKVSLAFOANL--TELDIYSRRLSQETGLEISEEINEEDLK 830
Sbjct 781 MTHS RK+S + Q N +E+DYSRRLS + EISEEINEEDLK 839
MTHSSTSPNKIHARRA-VRKMSMLSQTNFASSEIDYSRRLSEGDGSEISEEINEEDLK

Query 831 ECFDDMESIPAVTTWNTYLYRYTVHKSIFVLVWCLVIFLAEVAASLVVLLG----N 886
Sbjct 840 ECF D+ E TW+TYLRY+T +++L+VFLI CLVIFLAEVAASL LW++ N 899
ECFADEEIQNTTMSYLYRVTNRLVFLVLCVIFLAEVAASLGLWIIISGLAIN

Query 887 TPLQDKGNSTH-SRNNSYAVIITSTSSYYVYFIVYGVADTLAMGFFRGLPLVHTLIVTS 945
Sbjct 900 T Q ST S + ++ IT+ S YY+FYIYVG+AD+ LA+G RGLPLVHTL+TVS 959
TGSQNTDSTLDSLHVSFVKFITNGSHYIFIVYVGLADSFALGVIIRGLPLVHTLIVTS

Query 946 KILHMKHLHSLVQPMSTLNTLKAGGILNRFKSDIATLDDLLPLTFDFQILLIVIGAI 1005
Sbjct 960 KILHMKHLHSLVQPMSTLNTLKAGGILNRFKSDIATLDDLLPLTFDFVQILLIVIGAI 1019
KDLHKMLHSLVQPMSTLNTLKAGGILNRFKSDIATLDDLLPLTFDFVQILLIVIGAI

Query 1006 AVVAVLQPIYFVATVPVIVAFIMLRAYFLQTSQQLKQLESEGRSPIFTHLVTSKGLWTL 1065
Sbjct 1020 VV+VLOPY +A +FV V FIMLRAYFL+TSQQLKQLESE RSPIF+HL+SL+GLW+T 1079
CVVSVLQPIYLLAIFVAVIFIMLRAYFLRTSQQLKQLESEARSPIFTHLVTSKGLWTV

Query 1066 RAFGRQPYFTELFHKALNLTANWFLYLSLRLWFQMRIEMFVFFIAVTFISILITGEG 1125
Sbjct 1080 RAFGRQ YFTELFHKALNLTANWFLYLSLRLWFQMRI++FV+FFIAVTFI+I T G 1139
RAFGRQSYFTELFHKALNLTANWFLYLSLRLWFQMRI+IVVFFIAVTFIATATHDVG

Query 1126 EGRVGIILLAMNISTLQAVNNSIDVSLMRSVSRVFKIDPMDTEGKPTKSTKPKYKNG 1185
Sbjct 1140 EGVGIILLAMNISTLQAVNNSIDVGLMRSVSRVFKYIDIPPEGETKNRHRNANP 1199
EG+VGIILLAMNISTLQAVNNSIDV LMRVSRVFK+ID+P EG TK+ N

Query 1186 QLSKVMIIENSHVKDDIPSGGQMTVKDLTARYTEGGNAILNIFSI+SPQRVGLLGR 1245
Sbjct 1200 S Y+IEN H+ +WPSGQM +LTKRY G R++SFS+ QRQVGLLGR 1255
--SDVLVIENKHLTK--WPSGQMMVNNLTAKYTS+DGRAVLQDLFSVNVAGRVGLLGR

Query 1246 TGSKSTILLSAPRLLNTEGEIQDGVSWDSITLQQRKAFVGPQKVFIFSGTFRKRLND 1305
Sbjct 1256 TG+GKSTILLSA LRL+TEGEIQDGV+SW+S++LQ+WRKAFVGPQKVF+FGSTFRKRLND 1315
TGAGKSTILLSA LRL+TEGEIQDGVSWDSITLQQRKAFVGPQKVFIFSGTFRKRLND

Query 1306 PVEQWSDQEIWKVDEVLGRSVEIQFPGKLDVFLVDGGCVLSHGKQLMCLARSVLSKAK 1365
Sbjct 1316 PVEQWSD+EIWKV +EVLG+S+IEQFP KL+VFLVDGG +LS+GHKQLMCLARS+LSKAK 1375
PVEQWSDQEIWKVTEEVLGRSMIEQFPDKLNFVLDGGYILSNHGKQLMCLARSILSKAK

Query 1366 ILLDEPSAHLDPVTYQIIRTLKQAFADCTVILCEHRIEAMLECCQFVIEENKVRQYD 1425
Sbjct 1386 ILLDEP+AHLDVY+QIIR+TLK F++CTVIL EHR+EA+LECCOFLVIE V+O+D



► [NCBI/BLAST Home](#)

BLAST finds regions of similarity between biological sequences. [more...](#)

**New** [DELTA-BLAST](#), a more sensitive protein-protein search [Go](#)

## BLAST Assembled Genomes

Find Genomic BLAST pages:

Enter organism name or id--completions will be suggested [GO](#)

- |                       |                            |                             |
|-----------------------|----------------------------|-----------------------------|
| <a href="#">Human</a> | <a href="#">Rabbit</a>     | <a href="#">Zebrafish</a>   |
| <a href="#">Mouse</a> | <a href="#">Chimp</a>      | <a href="#">Clawed frog</a> |
| <a href="#">Rat</a>   | <a href="#">Guinea pig</a> | <a href="#">Arabidopsis</a> |
| <a href="#">Cow</a>   | <a href="#">Fruit fly</a>  | <a href="#">Rice</a>        |
| <a href="#">Pig</a>   | <a href="#">Honey bee</a>  | <a href="#">Yeast</a>       |
| <a href="#">Dog</a>   | <a href="#">Chicken</a>    | <a href="#">Microbes</a>    |

## Basic BLAST

Choose a BLAST program to run.

- |                                  |  |
|----------------------------------|--|
| <a href="#">nucleotide blast</a> | Search a <b>nucleotide</b> database using a <b>nucleotide</b> query<br><i>Algorithms: blastn, megablast, discontinuous megablast</i> |
| <a href="#">protein blast</a>    | Search <b>protein</b> database using a <b>protein</b> query<br><i>Algorithms: blastp, psi-blast, phi-blast, delta-blast</i>          |
| <a href="#">blastx</a>           | Search <b>protein</b> database using a <b>translated nucleotide</b> query  |
| <a href="#">tblastn</a>          | Search <b>translated nucleotide</b> database using a <b>protein</b> query  |
| <a href="#">tblastx</a>          | Search <b>translated nucleotide</b> database using a <b>translated nucleotide</b> query  |

## Specialized BLAST

Choose a type of specialized search (or database name in parentheses.)

- [Make specific primers with Primer-BLAST](#)
- [Cluster multiple sequences together with their database neighbors using MOLE-BLAST](#)
- [Find conserved domains](#) in your sequence (cds)
- Find sequences with similar [conserved domain architecture](#) (cdart)
- Search sequences that have [gene expression profiles](#) (GEO)
- Search [immunoglobulins and T cell receptor sequences](#) (IgBLAST)
- Screen sequence for [vector contamination](#) (vecscreen)
- [Align](#) two (or more) sequences using BLAST (bl2seq)
- Search [protein](#) or [nucleotide](#) targets in PubChem BioAssay
- Search [SRA by experiment](#)
- Constraint Based Protein [Multiple Alignment Tool](#)
- Needleman-Wunsch [Global Sequence Alignment Tool](#)
- Search [RefSeqGene](#)
- Search [trace archives](#)

### Your Recent Results **New!**

[All Recent results...](#)

### News

#### [BLAST 2.2.30+ released](#)

A new version of the stand-alone BLAST executables is now available.

Wed, 29 Oct 2014 15:00:00 EST

[More BLAST news...](#)

### Tip of the Day

#### [Use Genomic BLAST to see the genomic context](#)

If you are interested in the evolution of a particular gene or gene family it is often interesting to examine the intro-exon structure even across species.

[More tips...](#)

# Pairwise Sequence Comparison



# Detecting Similarity



- ◆ Is there a similarity between two sequences?
  - ◆ Identical symbols (nucleotides or amino acids)
  - ◆ Related symbols (amino acids)
  
- ◆ Do gaps/rearrangements allow for a higher degree of similarity?

# Symbol Comparison Tables (Scoring Matrices)



- ◆ How do we compare one sequence to another
  - ◆ What is a match?
  
- ◆ Define match values for all possible symbol comparisons
  - ◆ Nucleotides
  - ◆ Amino acids

# Nucleotide Tables



- ◆ Identical matches
  - ◆  $A=A, T=T, C=C, G=G, U=U, T=U$
  
- ◆ Matches with ambiguous bases
  - ◆ Based on ambiguity symbols
  - ◆  $Y=T, Y=C$
  - ◆  $R=A, R=G$
  
- ◆ Scores for matches or mismatches may differ depending on the table in use

# Amino Acid Tables



- ◆ Measure of similarity between amino acids
- ◆ Not a simple match/mismatch relationship
- ◆ Values vary depending on degree of relatedness
- ◆ Based on evolution, chemistry, or structure

# BLOSUM62 Table



- ◆ Default table for amino acid comparisons
- ◆ Many other similarity matrices are available

BLOSUM62 amino acid substitution matrix.

```
{
GAP_CREATE 12
GAP_EXTEND 4
}
```

	A	B	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	X	Y	Z
A	4	-2	0	-2	-1	-2	0	-2	-1	-1	-1	-1	-2	-1	-1	-1	1	0	0	-3	-1	-2	-1
B	-2	6	-3	6	2	-3	-1	-1	-3	-1	-4	-3	1	-1	0	-2	0	-1	-3	-4	-1	-3	2
C	0	-3	9	-3	-4	-2	-3	-3	-1	-3	-1	-1	-3	-3	-3	-3	-1	-1	-1	-2	-1	-2	-4
D	-2	6	-3	6	2	-3	-1	-1	-3	-1	-4	-3	1	-1	0	-2	0	-1	-3	-4	-1	-3	2
E	-1	2	-4	2	5	-3	-2	0	-3	1	-3	-2	0	-1	2	0	0	-1	-2	-3	-1	-2	5
F	-2	-3	-2	-3	-3	6	-3	-1	0	-3	0	0	-3	-4	-3	-3	-2	-2	-1	1	-1	3	-3
G	0	-1	-3	-1	-2	-3	6	-2	-4	-2	-4	-3	0	-2	-2	-2	0	-2	-3	-2	-1	-3	-2
H	-2	-1	-3	-1	0	-1	-2	8	-3	-1	-3	-2	1	-2	0	0	-1	-2	-3	-2	-1	2	0
I	-1	-3	-1	-3	-3	0	-4	-3	4	-3	2	1	-3	-3	-3	-3	-2	-1	3	-3	-1	-1	-3
K	-1	-1	-3	-1	1	-3	-2	-1	-3	5	-2	-1	0	-1	1	2	0	-1	-2	-3	-1	-2	1
L	-1	-4	-1	-4	-3	0	-4	-3	2	-2	4	2	-3	-3	-2	-2	-2	-1	1	-2	-1	-1	-3
M	-1	-3	-1	-3	-2	0	-3	-2	1	-1	2	5	-2	-2	0	-1	-1	-1	1	-1	-1	-1	-2
N	-2	1	-3	1	0	-3	0	1	-3	0	-3	-2	6	-2	0	0	1	0	-3	-4	-1	-2	0
P	-1	-1	-3	-1	-1	-4	-2	-2	-3	-1	-3	-2	-2	7	-1	-2	-1	-1	-2	-4	-1	-3	-1
Q	-1	0	-3	0	2	-3	-2	0	-3	1	-2	0	0	-1	5	1	0	-1	-2	-2	-1	-1	2
R	-1	-2	-3	-2	0	-3	-2	0	-3	2	-2	-1	0	-2	1	5	-1	-1	-3	-3	-1	-2	0
S	1	0	-1	0	0	-2	0	-1	-2	0	-2	-1	1	-1	0	-1	4	1	-2	-3	-1	-2	0
T	0	-1	-1	-1	-1	-2	-2	-2	-1	-1	-1	-1	0	-1	-1	-1	1	5	0	-2	-1	-2	-1
V	0	-3	-1	-3	-2	-1	-3	-3	3	-2	1	1	-3	-2	-2	-3	-2	0	4	-3	-1	-1	-2
W	-3	-4	-2	-4	-3	1	-2	-2	-3	-3	-2	-1	-4	-4	-2	-3	-3	-2	-3	11	-1	2	-3
X	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1
Y	-2	-3	-2	-3	-2	3	-3	2	-1	-2	-1	-1	-2	-3	-1	-2	-2	-2	-1	2	-1	7	-2
Z	-1	2	-4	2	5	-3	-2	0	-3	1	-3	-2	0	-1	2	0	0	-1	-2	-3	-1	-2	5



# Dot Plots

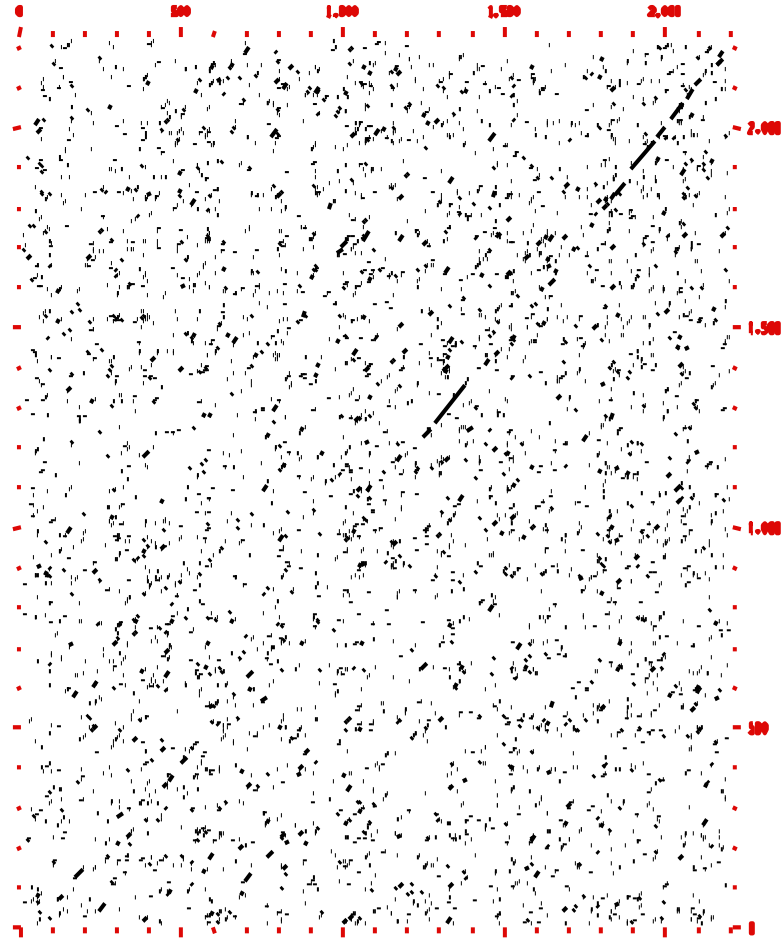


- ◆ Allow comparison of two sequences in all registers
- ◆ Produces a graph (Dotplot) of sequence similarities
- ◆ The human brain interprets the results

# Polio:Hepatitis A 20/10 (Protein)

DOTPLOT of: pol-tyo-20.pnt Density: 2532.95 February 9, 1992 17:48  
COMPARE Mirrors: 20 SIRrings: 10.9 Points: 12,002

Polgetpovl ck: 1,208, 1 to 2,227



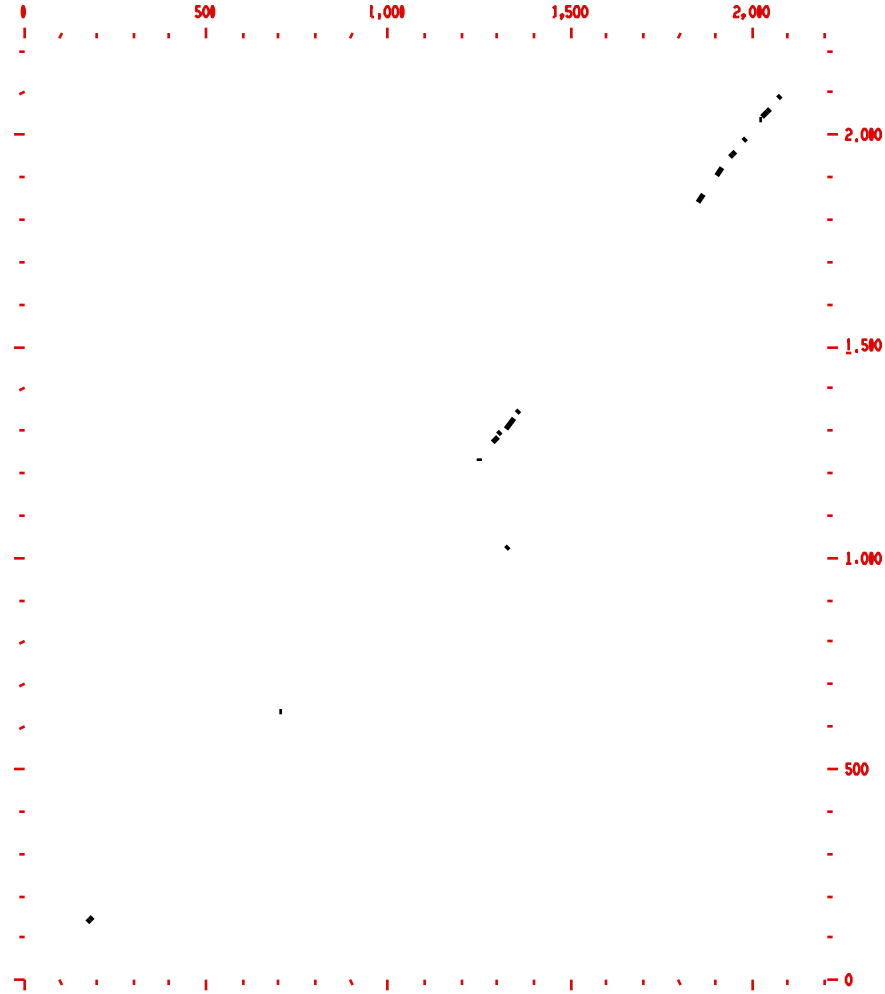
Polgetpol[m ck: 7,460, 1 to 2,207

# Polio:Hepatitis A 20/15 (Protein)

DOTPLOT of Pol-Hpa.Pnt:1 Density: 2532.95 February 9, 1992 17:57

COMPARE Window: 20 Stringency: 15.0 Points: 190

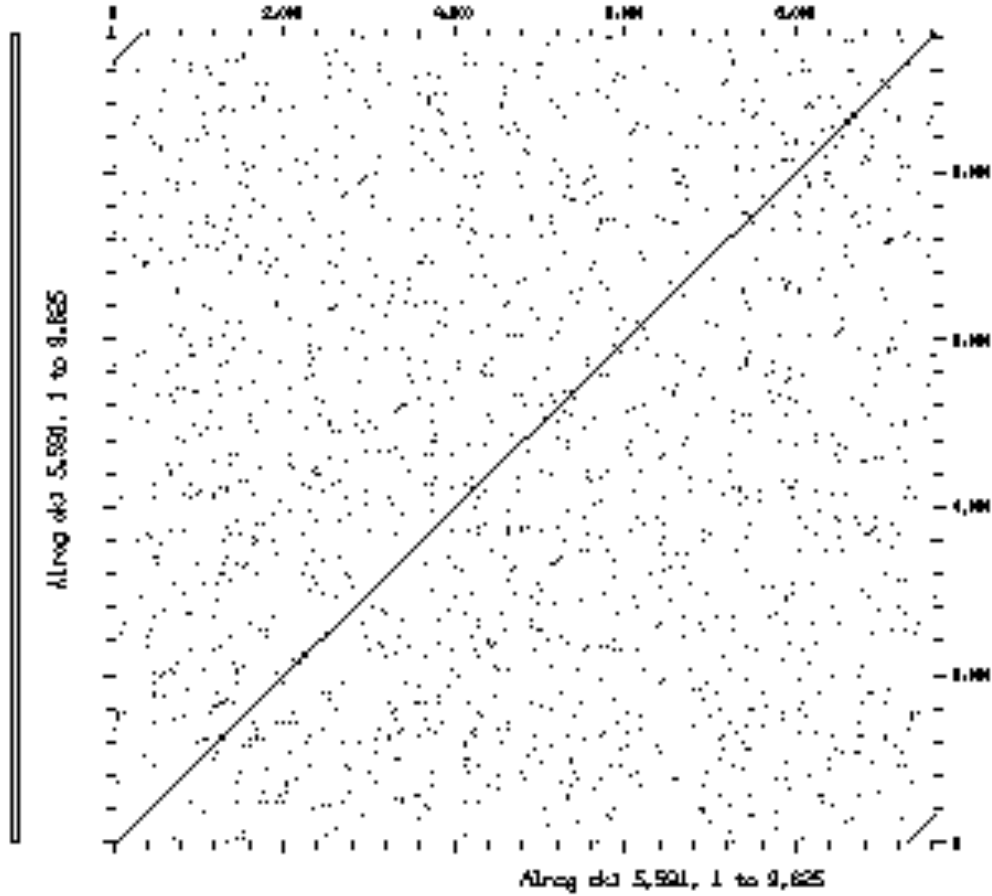
PolgsHpaV1 cl# 1.208, 1 to 2.227



PolgsPolm cl# 7,480, 1 to 2,207

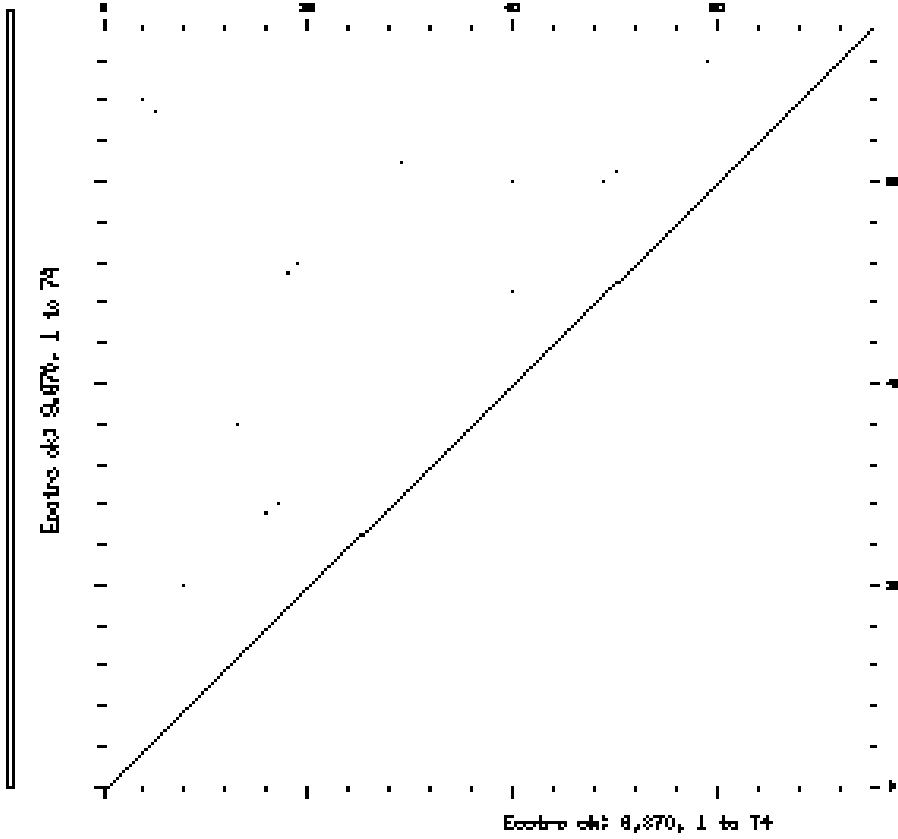
# Rous Sarcoma Virus vs. Itself

OUTPUT of Align\_Hand.Prt Date: 01/18/88 11:27 February 4, 1980 02:15  
COMPARE/Hand Hand-Bliss 3 Alphabet: 4 Ref: 12,384



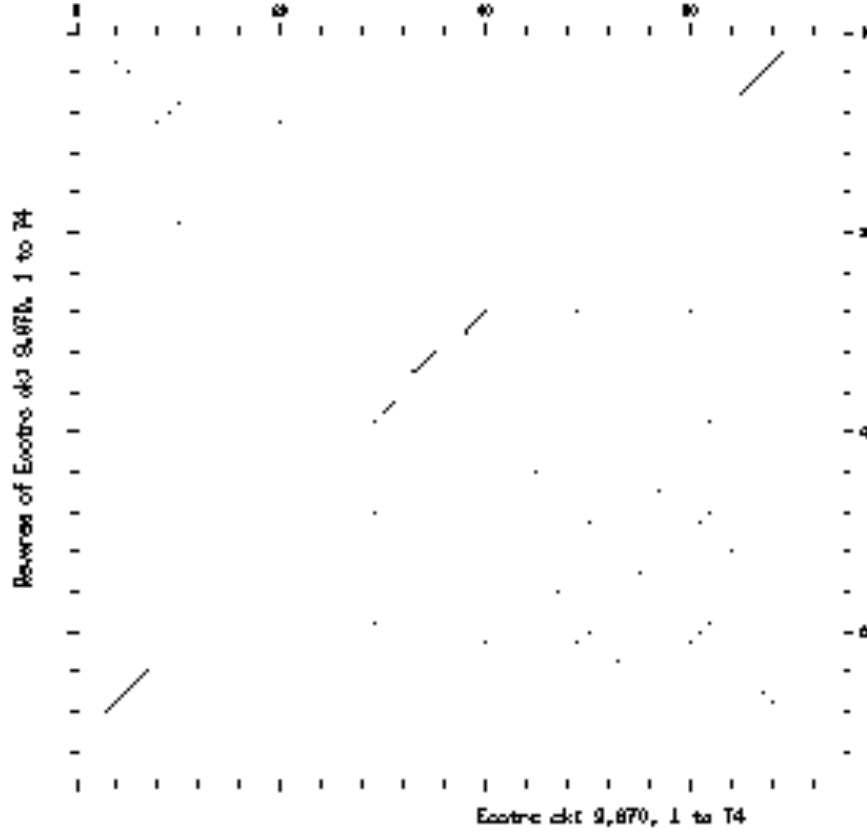
# E. coli tRNA-Cys vs. Itself

COMPUT off 04/05/91 08:30 02/02/91 14:36  
COMPARE/View 04/05/91 08:30 02/02/91 14:36



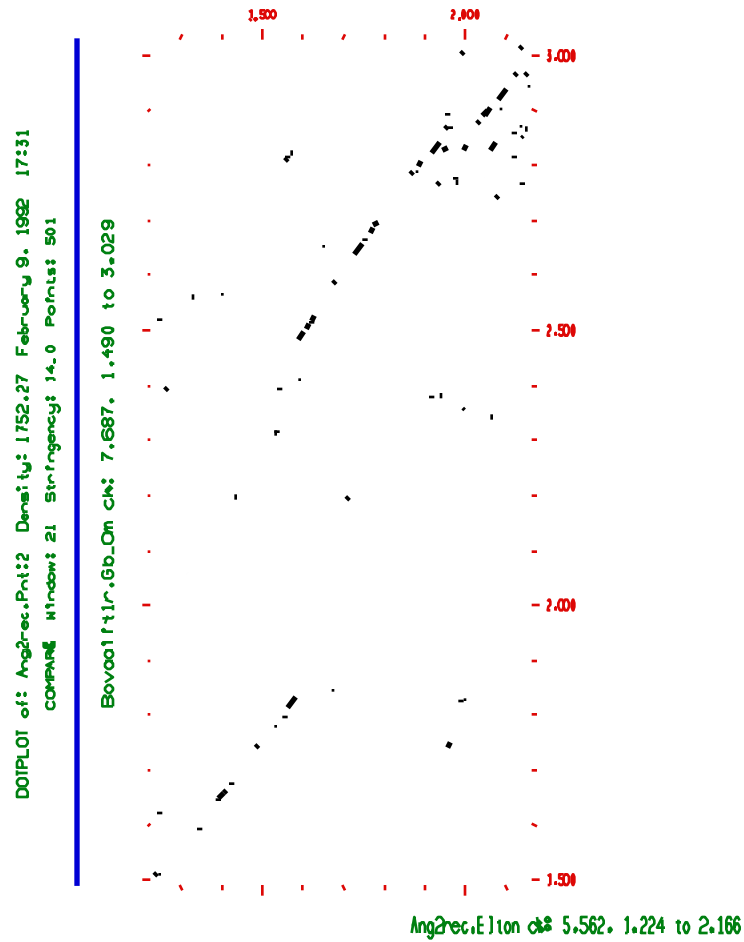
# E. coli tRNA-Cys vs. Its Reverse-Complement

OUTPUT of Ecolns.Pyt. Geneflyd ed.36 February 7, 1986 14948  
COMPARE/Word Hard-82aa 4 Alphabet: 4 Rev/In: 98



# Angiotensin II Receptor mRNA

## 3' non-coding end; Bovine:Rat



Π σ Pτ

# Computer Generated Alignments



- ◆ Align two (or more) sequences by the introduction of gaps
  - ◆ Maximize sequence identity/similarity
  - ◆ Minimize gaps



# Local vs. Global Alignments



- ◆ Local: What is the best region of similarity between two sequences?
  - ◆ Smith and Waterman
  - ◆ Not necessarily the whole sequence
  
- ◆ Global: What is the best possible alignment between two sequences?
  - ◆ Needleman–Wunsch
  - ◆ Always the whole sequence

# Quality for Nucleotide Alignments



- ◆ Matches Rewarded
- ◆ Mismatches Penalized
- ◆ Gaps Penalized
  - ◆ GapWeight
    - ◆ Penalty for the introduction of a new gap
  - ◆ GapLengthWeight
    - ◆ Penalty for extension of an existing gap

# Quality for Protein Alignments



- ◆ GapWeight and GapLengthWeight apply as for nucleotides
- ◆ A sum of comparison values between aligned amino acids is used in place of match and mismatch values



BLAST®

Basic Local Alignment Search Tool

My NCBI

[\[Sign In\]](#) [\[Register\]](#)[Home](#)[Recent Results](#)[Saved Strategies](#)[Help](#)[NCBI/BLAST/blastp suite](#)

## Align Sequences Protein BLAST

[blastn](#) [blastp](#) [blastx](#) [tblastn](#) [tblastx](#)BLASTP programs search protein subjects using a protein query. [more...](#)[Reset page](#) [Bookmark](#)

## Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#)

NP\_000483.3

Query subrange [Clear](#)

From

To

Or, upload file

Choose File no file selected

Job Title

Enter a descriptive title for your BLAST search [Clear](#) Align two or more sequences [Clear](#)

## Enter Subject Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#)

NP\_001038348.1

Subject subrange [Clear](#)

From

To

Or, upload file

Choose File no file selected

## Program Selection

Algorithm

 blastp (protein-protein BLAST)Choose a BLAST algorithm [Clear](#)**BLAST**Search [protein sequence](#) using [Blastp \(protein-protein BLAST\)](#) Show results in a new window[+ Algorithm parameters](#)

# Multiple sequence alignments



Get it right the first time!!

# Why MSAs?



- ◆ Starting point for many types of comparative analyses
- ◆ Identification of functional motifs
  - ◆ Regulatory
  - ◆ Protein domains
- ◆ Evolutionary inference
  - ◆ Evolutionary history
  - ◆ Selection pressures

# Conserved functional domains



- ◆ Sequences required for common function
- ◆ Conserved among different species

# Variable domains



- ◆ Sequences under selection
  - ◆ Antibody escape mutants



# Peptide motifs



- ◆ Active sites
- ◆ Binding motifs
- ◆ Protein modification motifs

# The Alignment



- ◆ Critical!!!
- ◆ For coding regions first align the protein sequences and then align the nucleotide sequences to the protein sequence alignment
- ◆ Construct trees for both the protein and nucleotide alignments

# Choose Your Program



- ◆ Simple nt or aa alignments
  - ◆ Clustal
- ◆ Alignment of many (hundreds) of sequences
- ◆ Large alignments (Whole genomes)
  - ◆ Sequence lengths from thousands to millions of residues

# MSA Programs



- ◆ Clustal (W and X)
- ◆ DIALIGN
- ◆ MAFFT
- ◆ MAUVE
- ◆ MAVID
- ◆ MUSCLE
- ◆ T-Coffee
- ◆ ...

# Choose your Parameters



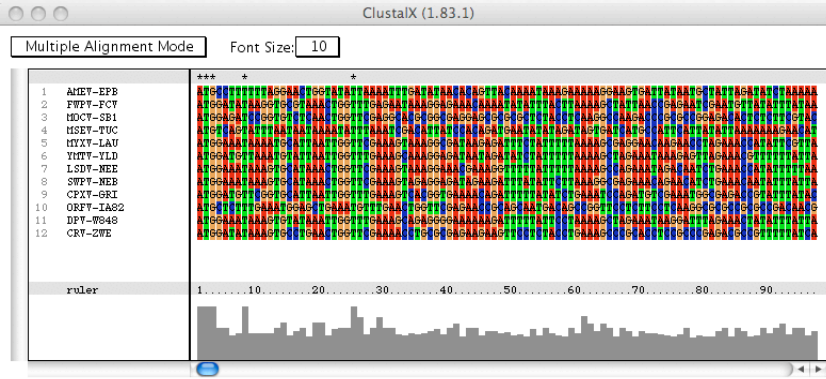
- ◆ Gap weight penalty
  - ◆ Penalty to introduce a new gap
- ◆ Gap length weight penalty
  - ◆ Penalty to extend a gap
- ◆ Substitution matrices
  - ◆ Evolutionary distance for various aa substitutions
- ◆ Start with the default settings
  - ◆ Inspect the alignment
  - ◆ Change parameters as needed

# Inspect Your Results

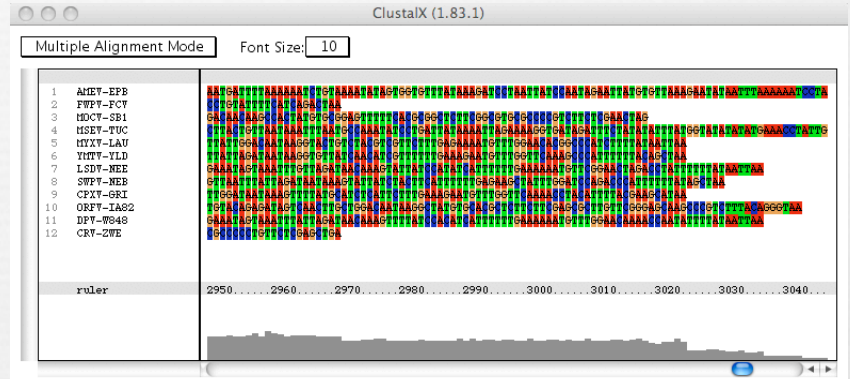


- ◆ Alignment Uncertainty and Genomic Analysis  
Science v319, p473, 1/25/2008
  - ◆ Learn to critically evaluate anything the computer tells you
  - ◆ Understand the limitations of any particular analysis tool
  - ◆ Understand the parameters that influence the behavior of any particular analysis tool

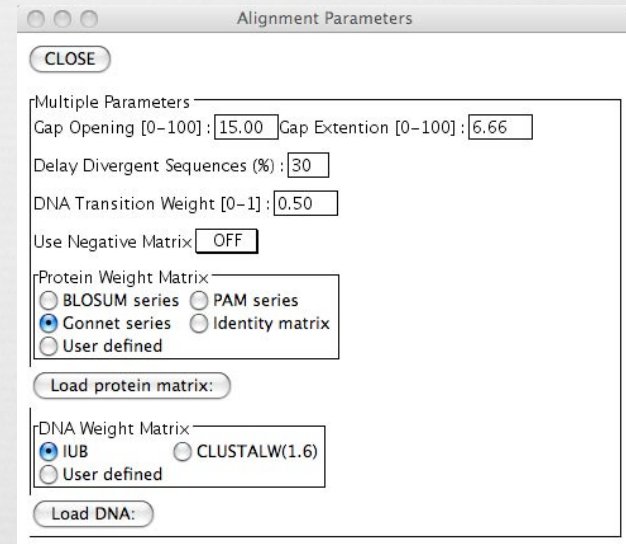
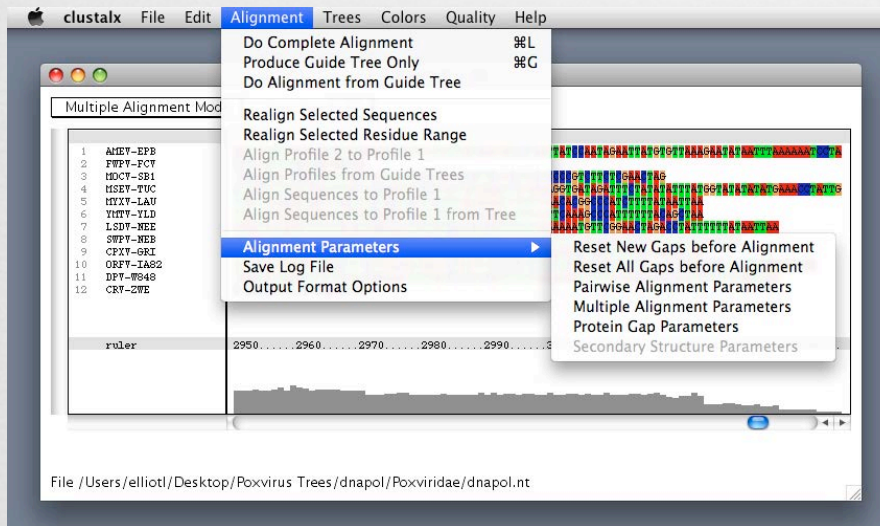
# Using ClustalX



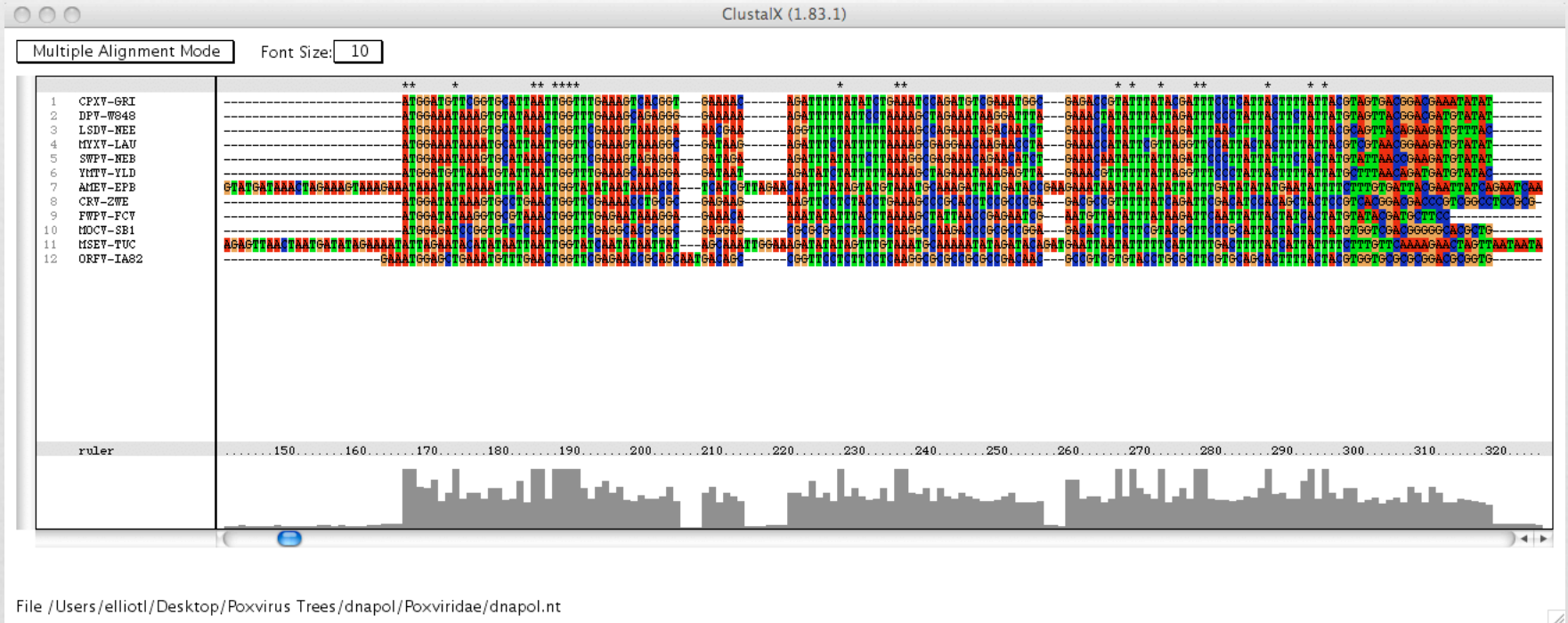
File /Users/elliott/Desktop/Poxvirus Trees/dnapol/Poxviridae/dnapol.nt



File /Users/elliott/Desktop/Poxvirus Trees/dnapol/Poxviridae/dnapol.nt

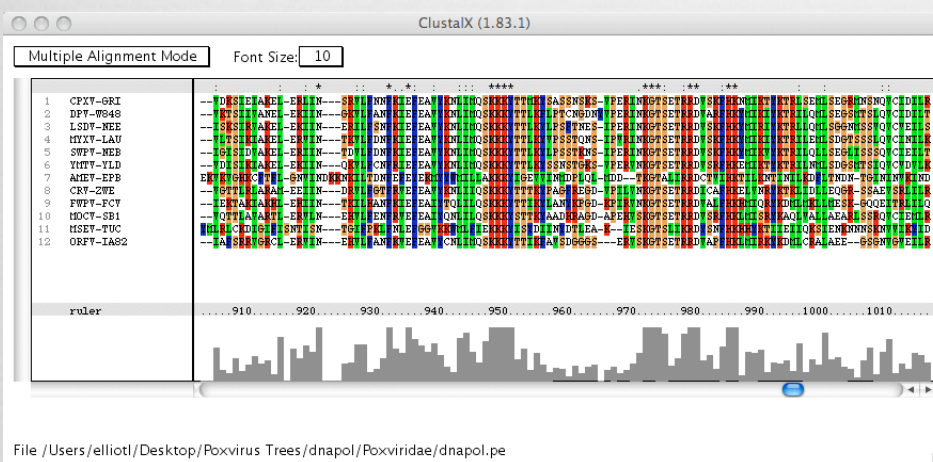
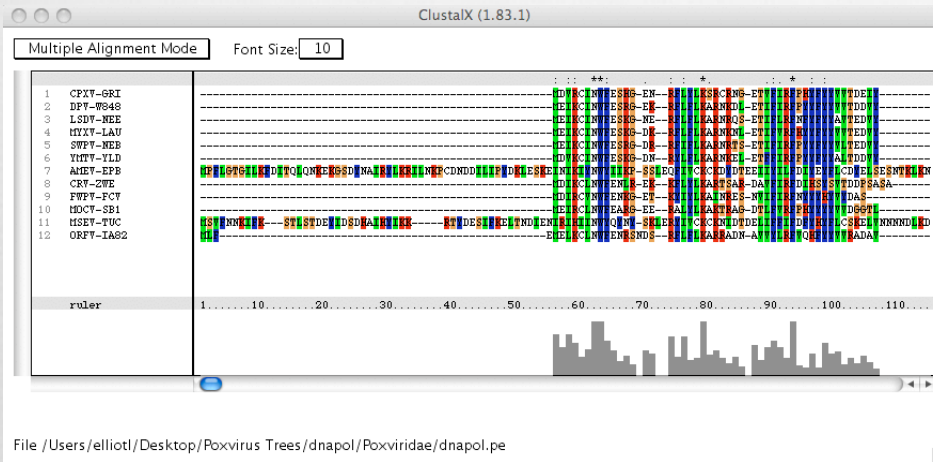


# Clustal NT Alignment





# Clustal AA Alignment



- ◆ Align Protein Seqs
- ◆ “Back translate” NT Seqs to AA alignment
- ◆ MEGA
- ◆ T-Coffee
- ◆ Provides codon-aligned NT MSA



# Genomics



# Genomic Resources



- ◆ NCBI Genome Resources
  - ◆ <http://www.ncbi.nih.gov/Genomes/>
  
- ◆ Ensembl
  - ◆ European Molecular Biology Laboratory (EMBL)
    - ◆ European Bioinformatics Institute
  - ◆ <http://www.ensembl.org/>
  
- ◆ UCSC Genome Browser
  - ◆ <http://genome.ucsc.edu>

Genomes Genome Browser Tools Mirrors Downloads My Data Help About Us

### Human (*Homo sapiens*) Genome Browser Gateway

The UCSC Genome Browser was created by the Genomics Bioinformatics Group of UC Santa Cruz. Software Copyright © The Regents of the University of California. All rights reserved.

group genome assembly position search term

Mar mammal Human Feb. 2009 (GRCh37/hg19) chr21:33,031,597-33,041,570 enter position, gene symbol or search terms submit

[Click here to reset](#) the browser user interface settings to their defaults.

[track search](#) [add custom tracks](#) [track hubs](#) [configure tracks and display](#)

#### Human Genome Browser – hg19 assembly (sequences)

The February 2009 human reference sequence (GRCh37) was produced by the [Genome Reference Consortium](#). For more information about this assembly, see [GRCh37](#) in the NCBI Assembly database.



Human Genome Browser  
(Graphic courtesy of GENIE)

#### Sample position queries

A genome position can be specified by the accession number of a sequenced genomic clone, an mRNA or EST or STS marker, a chromosomal coordinate range, or keywords from the GenBank description of an mRNA. The following list shows examples of valid position queries for the human genome. See the [User's Guide](#) for more information.

Request:	Genome Browser Response:
chr7	Displays all of chromosome 7
chrUn_g000212	Displays all of the unplaced contig g000212
20p13	Displays region for band p13 on chr 20
chr3:1-100000	Displays first million bases of chr 3, counting from p-arm telomere
chr3:100000+2000	Displays a region of chr3 that spans 2000 bases, starting with position 100000
RH18061;RH80175	Displays region between genome landmarks, such as the STS markers RH18061 and RH80175, or chromosome bands 15q11 to 15q13, or SNPs rs1042522 and rs1800370. This syntax may also be used for other range queries, such as between uniquely determined ESTs, mRNAs, refSeqs, etc.
rs1042522;rs1800370	
D16S3046	Displays region around STS marker D16S3046 from the Genethon/Marshfield maps. Includes 100,000 bases on each side as well.
AA205474	Displays region of EST with GenBank accession AA205474 in BRCA1 cancer gene on chr 17
AC008101	Displays region of clone with GenBank accession AC008101
AF083811	Displays region of mRNA with GenBank accession number AF083811
PRNP	Displays region of genome with HUGO Gene Nomenclature Committee identifier PRNP
NM_017414	Displays the region of genome with RefSeq identifier NM_017414
NP_059110	Displays the region of genome with protein accession number NP_059110
pseudogene mRNA	Lists transcribed pseudogenes, but not cDNAs
homeobox caudal	Lists mRNAs for caudal homeobox genes
zinc finger	Lists many zinc finger mRNAs
kruppel zinc finger	Lists only kruppel-like zinc fingers
huntington	Lists candidate genes associated with Huntington's disease
zahler	Lists mRNAs deposited by scientist named Zahler
Evans J.E.	Lists mRNAs deposited by co-author J.E. Evans

Use this last format for author queries. Although GenBank requires the search format *Evans J.E.*, internally it uses the format *Evans,J.E.*

#### Assembly Details

The GRCh37 build reference sequence is considered to be "finished", a technical term indicating that the sequence is highly accurate (with fewer than one error per 10,000 bases) and highly contiguous (with the only remaining gaps corresponding to regions whose sequence cannot be reliably resolved with current technology). Future work on the reference sequence will focus on improving accuracy and reducing gaps in the sequence. Statistics for the GRCh37 build assembly can be found on the NCBI [Build 37.1 Statistics](#) web page. For a glossary of assembly-related terms, please see the [GRC Assembly Terminology](#) page.

#### Note on chrM

Since the release of the UCSC hg19 assembly, the *Homo sapiens* mitochondrion sequence (represented as "chrM" in the Genome Browser) has been replaced in GenBank with the record [NC\\_012920](#). We have not replaced the original sequence, [NC\\_001807](#), in the hg19 Genome Browser. We plan to use the [Revised Cambridge Reference Sequence \(rCRS\)](#) in the next human assembly release.

#### Chromosome naming scheme

In addition to the "regular" chromosomes, the hg19 browser contains nine haplotype chromosomes and 59 unplaced contigs. If an unplaced contig is localized to a chromosome, the contig name is appended to the regular chromosome name, as in [chr11\\_q10q11.2\\_100000](#). If the chromosome is unknown, the contig is represented with the name "chrUn" followed by the contig identifier, as in [chrUn\\_g100211](#). Note that the chrUn contigs are no longer placed in a single, artificial chromosome as they have been in previous UCSC assemblies. See the [sequences](#) page for a complete list of hg19 chromosome names.

The nine haplotype chromosomes are:

name	accession	UCSC chr name
HSCHR6_MHC_APD_CTG1	GL000250.1	chr6_apd_hap1
HSCHR6_MHC_COX_CTG1	GL000251.1	chr6_cox_hap2
HSCHR6_MHC_DBB_CTG1	GL000252.1	chr6_dbb_hap3
HSCHR6_MHC_MANN_CTG1	GL000253.1	chr6_mann_hap4
HSCHR6_MHC_MCF_CTG1	GL000254.1	chr6_mcf_hap5
HSCHR6_MHC_QBL_CTG1	GL000255.1	chr6_qbl_hap6
HSCHR6_MHC_SSTO_CTG1	GL000256.1	chr6_ssto_hap7
HSCHR4_1_CTG9	GL000257.1	chr4_ctg9_hap1
HSCHR17_1_CTG5	GL000258.1	chr17_ctg5_hap1

See the [Wellcome Trust Sanger Institute MHC Haplotype Project](#) web site for additional information on the chr6 alternate haplotype assemblies.

The Y chromosome in this assembly contains two pseudoautosomal regions (PARs) that were taken from the corresponding regions in the X chromosome and are exact duplicates:

chrY:10001-2649520 and chrY:59034050-59363566  
chrX:60001-2699520 and chrX:154601044-155260560

For further information on GRCh37 build see the NCBI [GRCh37 release notes](#).

Bulk downloads of the sequence and annotation data are available via the Genome Browser [FTP server](#) or the [Downloads](#) page. We recommend that you use [rsync](#) for downloading large or multiple files.

The hg19 annotation tracks were generated by UCSC and collaborators worldwide. See the [Credits](#) page for a detailed list of the organizations and individuals who contributed to this release.

#### Genbank Pipeline Details

For the purposes of the Genbank alignment pipeline, this assembly is considered to be **finished**.

# One Final Word of Wisdom...



- ◆ “...although the computer is a wonderful helpmate for the sequence searcher and comparer, biochemists and molecular biologists must guard against the blind acceptance of any algorithmic output; given the choice, think like a biologist and not a statistician.”
  - ◆ - Russell F. Doolittle, 1990