

Argomenti del corso:

Introduzione alla bioinformatica

Strumenti della bioinformatica

Ricerca su banche dati

Allineamenti di sequenze

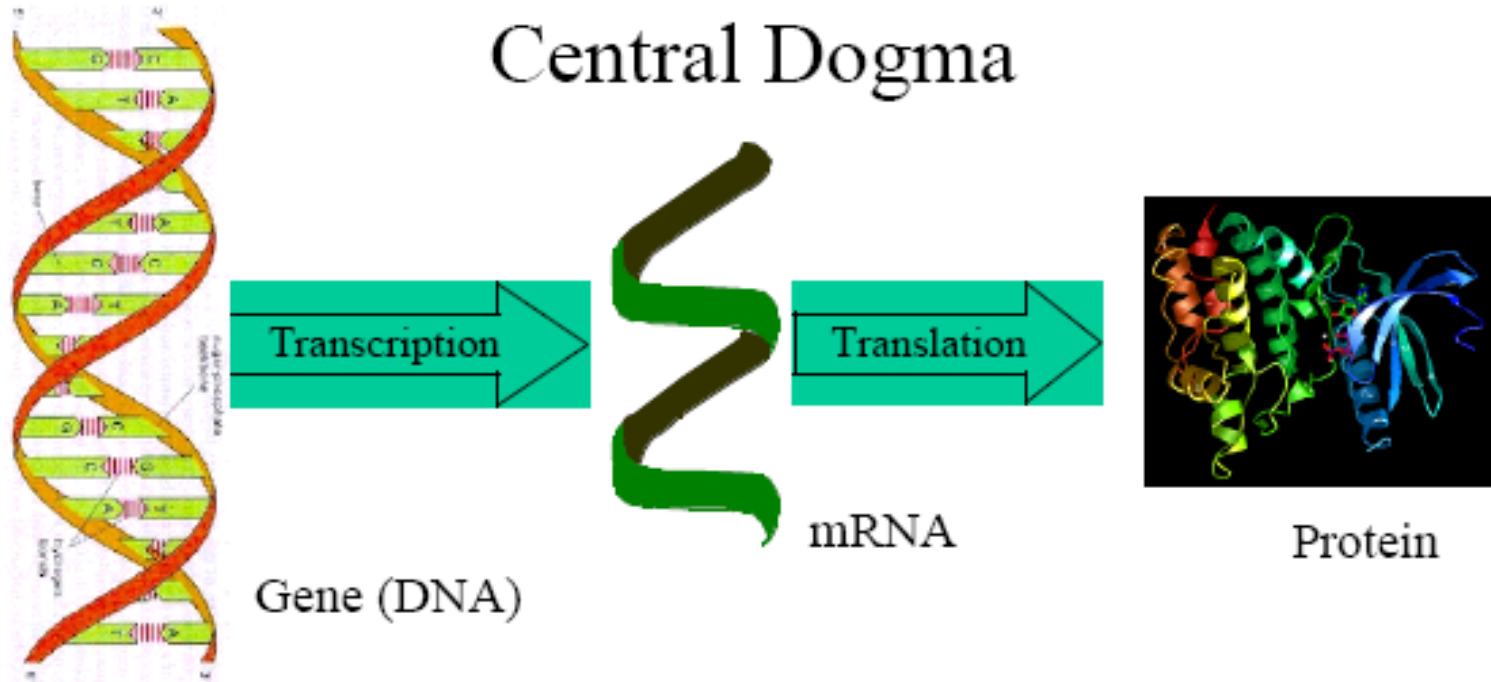
Analisi struttura proteine

Famiglie proteine

Alberi filogenetici

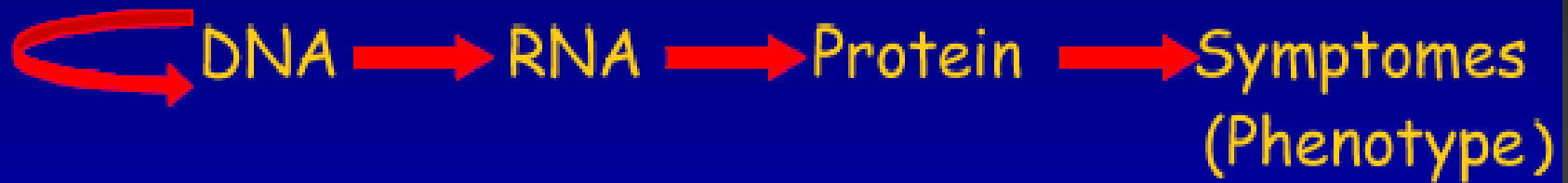
Metodi per la predizione della struttura

Central Dogma



Cells express different subset of the genes in different tissues and under different conditions

Central Paradigm of Molecular Biology





What is Bioinformatics ?

Bioinformatics is the use of computers for the acquisition, management, and analysis of biological information.

It incorporates elements of molecular biology, computational biology, database computing, and the Internet...

... bioinformatics is clearly a multi-disciplinary field including: computer systems management networking, database design, computer programming, molecular biology

From Using Computers for Molecular Biology, Stuart M. Brown, PhD, RCR, NYU Medical Center



Bioinformatics is a multifaceted discipline combining many scientific fields including computational biology, statistics, mathematics, molecular biology and genetics (Fenstermacher, 2005, p. 440).

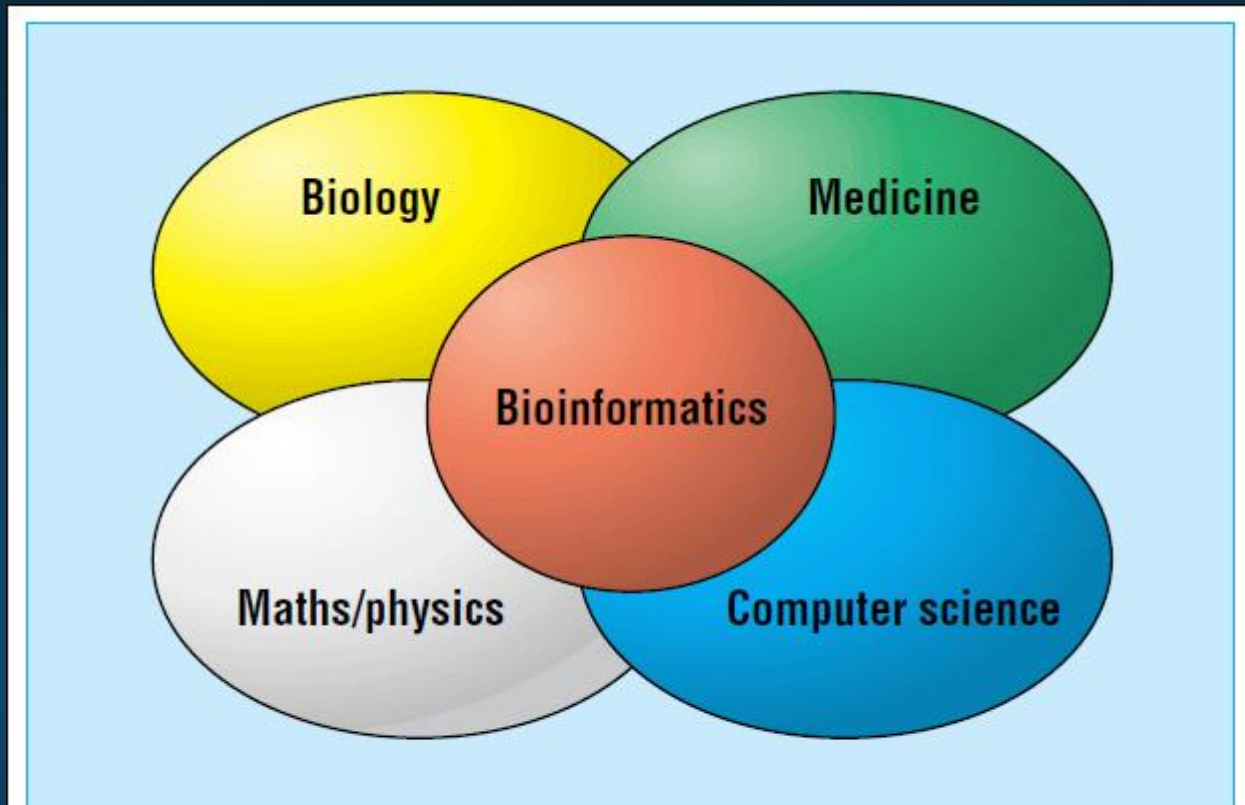


Fig 1 Interaction of disciplines that have contributed to the formation of bioinformatics

Bioinformatics: Origins & Definitions

Bioinformatics - a definition¹

(Molecular) bio – informatics: bioinformatics is conceptualising biology in terms of molecules (in the sense of physical chemistry) and applying "*informatics techniques*" (derived from disciplines such as applied maths, computer science and statistics) to *understand* and *organise* the *information* associated with these molecules, on a *large scale*. In short, bioinformatics is a management information system for molecular biology and has many *practical applications*.

¹ As submitted to the Oxford English Dictionary

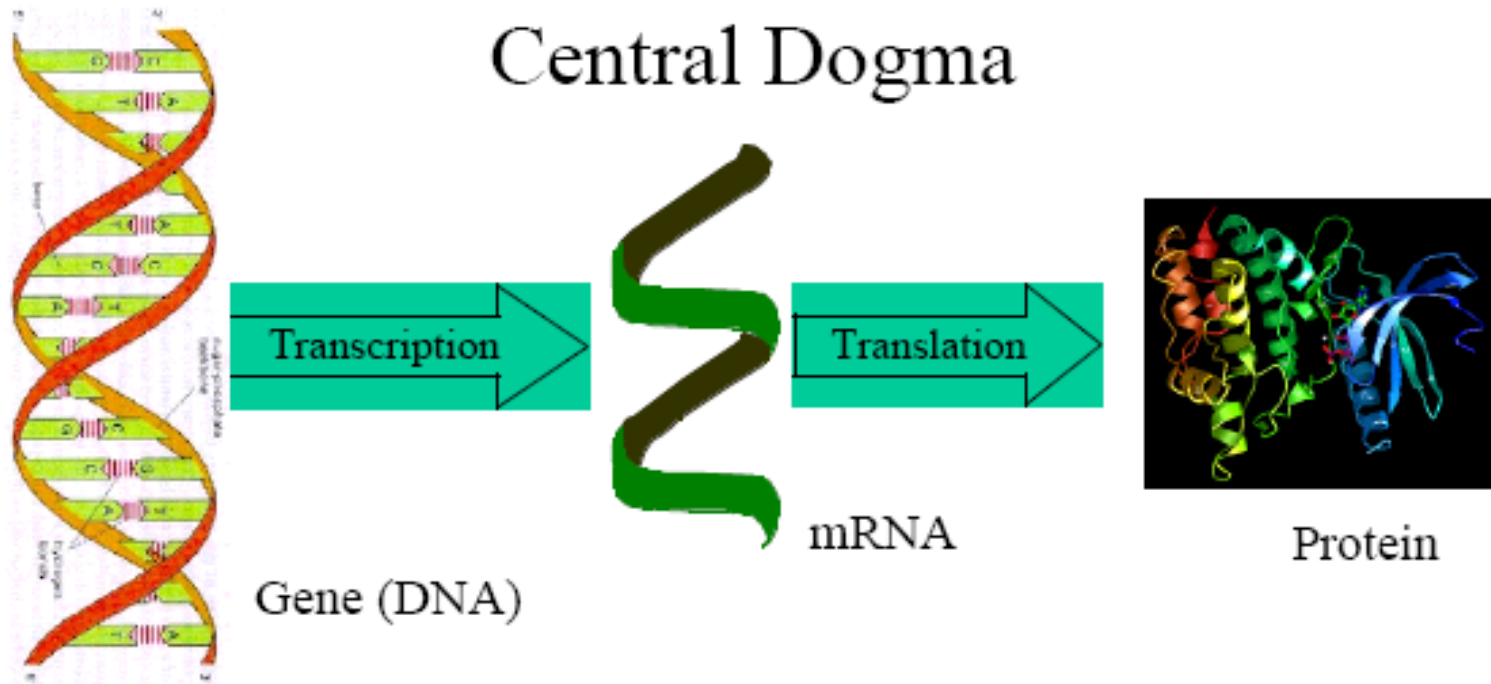
Bioinformatics has many definitions

... the study of how information is represented and analyzed in biological systems, starting at the molecular level ... concerned with understanding how basic biological systems conspire to create molecules, organelles, living cells, organs, and entire organisms (Altman & Mooney, 2006, p. 763)

... application of tools of computation and analysis to the capture and interpretation of biological data (Bayat, 2003, p. 1018)

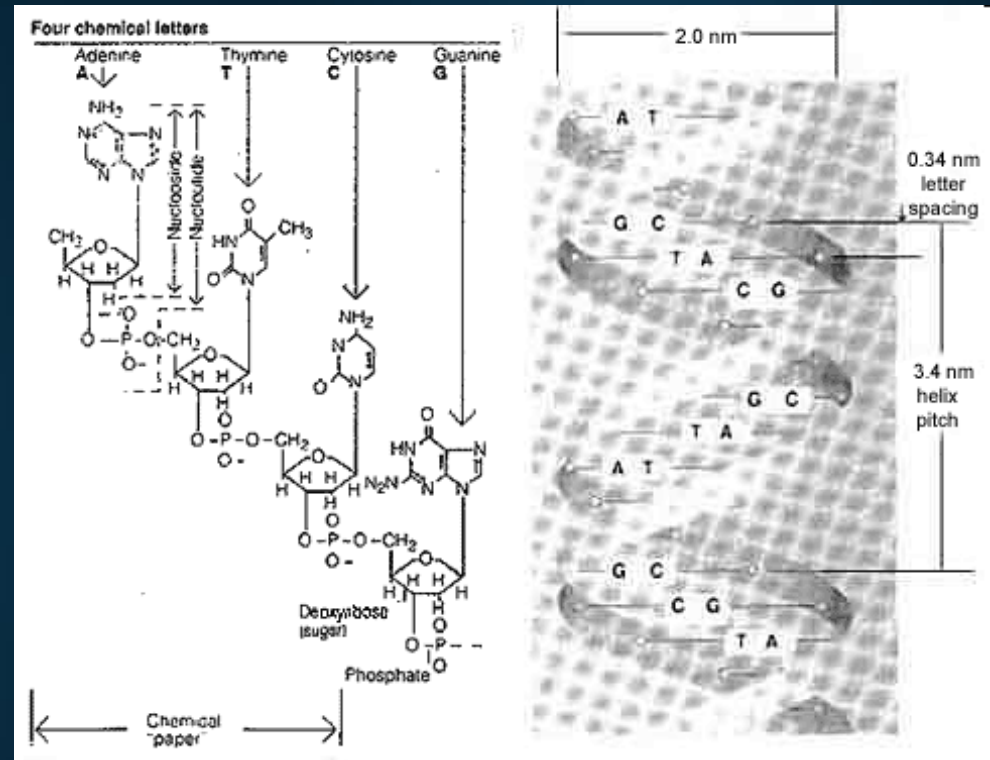
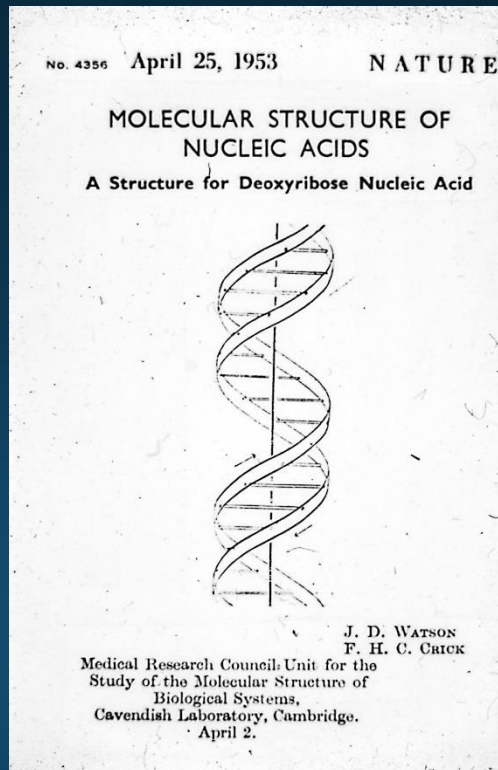


Central Dogma



Cells express different subset of the genes in different tissues and under different conditions

DNA is the nature's universal information storage medium



... increasingly, biological research relies on information science





The Human Genome Project

- Produced the human genome sequence
- Spawned a new field: genomics
- Spurred new technologies
- And now provides us an unparalleled opportunity to apply new knowledge, technologies, and approaches to health care

Guttmacher (2009)



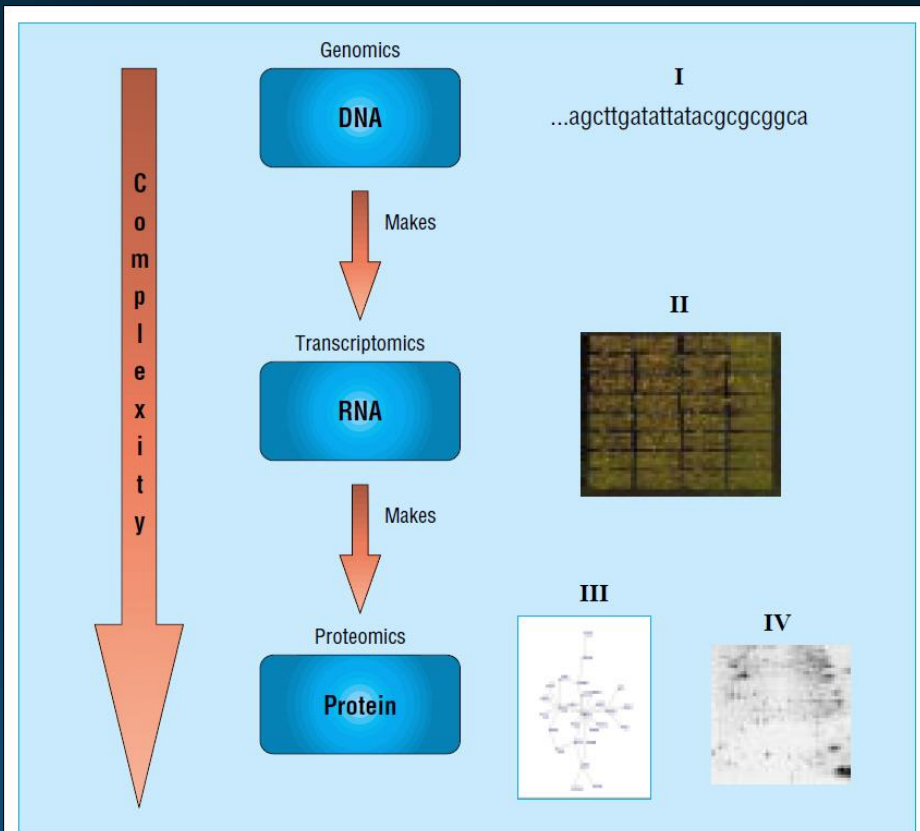


Fig 4 Schematic diagram representing complexity of genomic data processing. Analysis and interpretation of biological data considers information at every level from the genome (total genetic content) to the proteome (total protein content) and transcriptome (total messenger RNA content) of the cell. The images numbered I-IV to the right of the diagram represent relevant examples of DNA (image I is base pair nucleotides); RNA (image II is a microarray showing levels of gene expression); and protein (image III is a structure of a single protein; image IV is a two dimensional gel electrophoresis showing separation of all proteins of a cell—each spot corresponds to a different protein chain)

Bioinformatics
supports
“-omics”
research

...from Bayat (2002), p 1020.



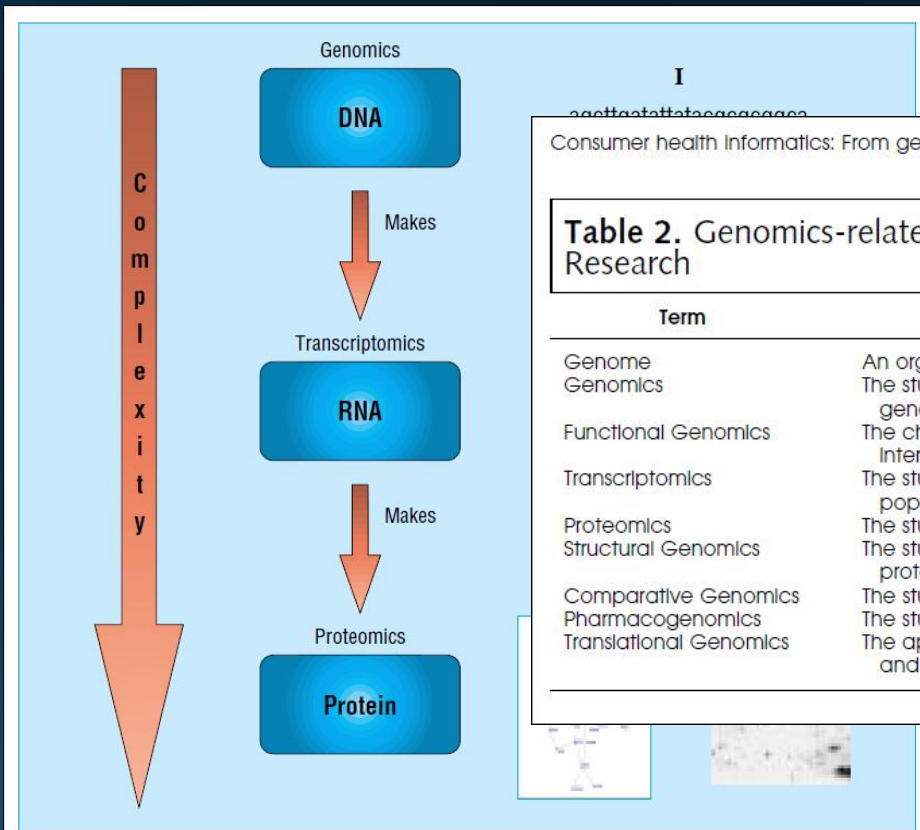


Fig 4 Schematic diagram representing complexity of genomic data processing. Analysis and interpretation of biological data considers information at every level from the genome (total genetic content) to the proteome (total protein content) and transcriptome (total messenger RNA content) of the cell. The images numbered I-IV to the right of the diagram represent relevant examples of DNA (image I is base pair nucleotides); RNA (image II is a microarray showing levels of gene expression); and protein (image III is a structure of a single protein; image IV is a two dimensional gel electrophoresis showing separation of all proteins of a cell—each spot corresponds to a different protein chain)

Table 2. Genomics-related Terminology and Description of Related Research

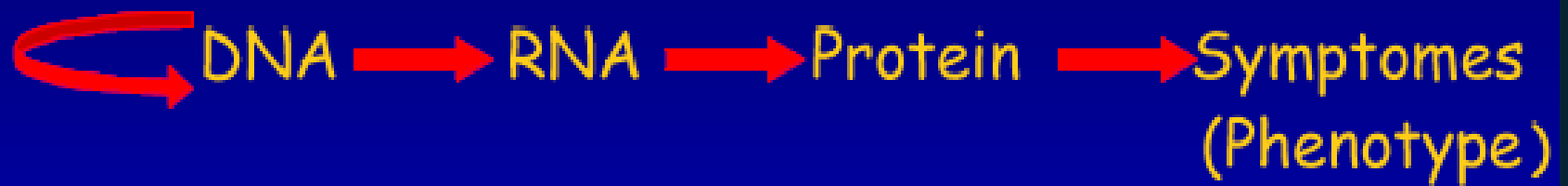
Term	Description of Related Research
Genome	An organism's total set of genes.
Genomics	The study of an organism's entire genome, including the interactions between genes and the interactions between genes and environment.
Functional Genomics	The characterization of genes to determine their structure, function, and interactions.
Transcriptomics	The study of the expression level of mRNA (transcript) in a given cell population under given conditions.
Proteomics	The study of proteins to establish their structure and function.
Structural Genomics	The study and subsequent development of 3D structures of one or more proteins from each protein family.
Comparative Genomics	The study of the relationship or comparison of genomes across species.
Pharmacogenomics	The study of the influence of genetic variation on drug response.
Translational Genomics	The application or translation of genomics and genomics-related discovery and technology into clinically useful information and tools.

...from McDaniel, Schutte, & Keller (2008), p. 220

...from Bayat (2002), p 1020.



Central Paradigm of Molecular Biology



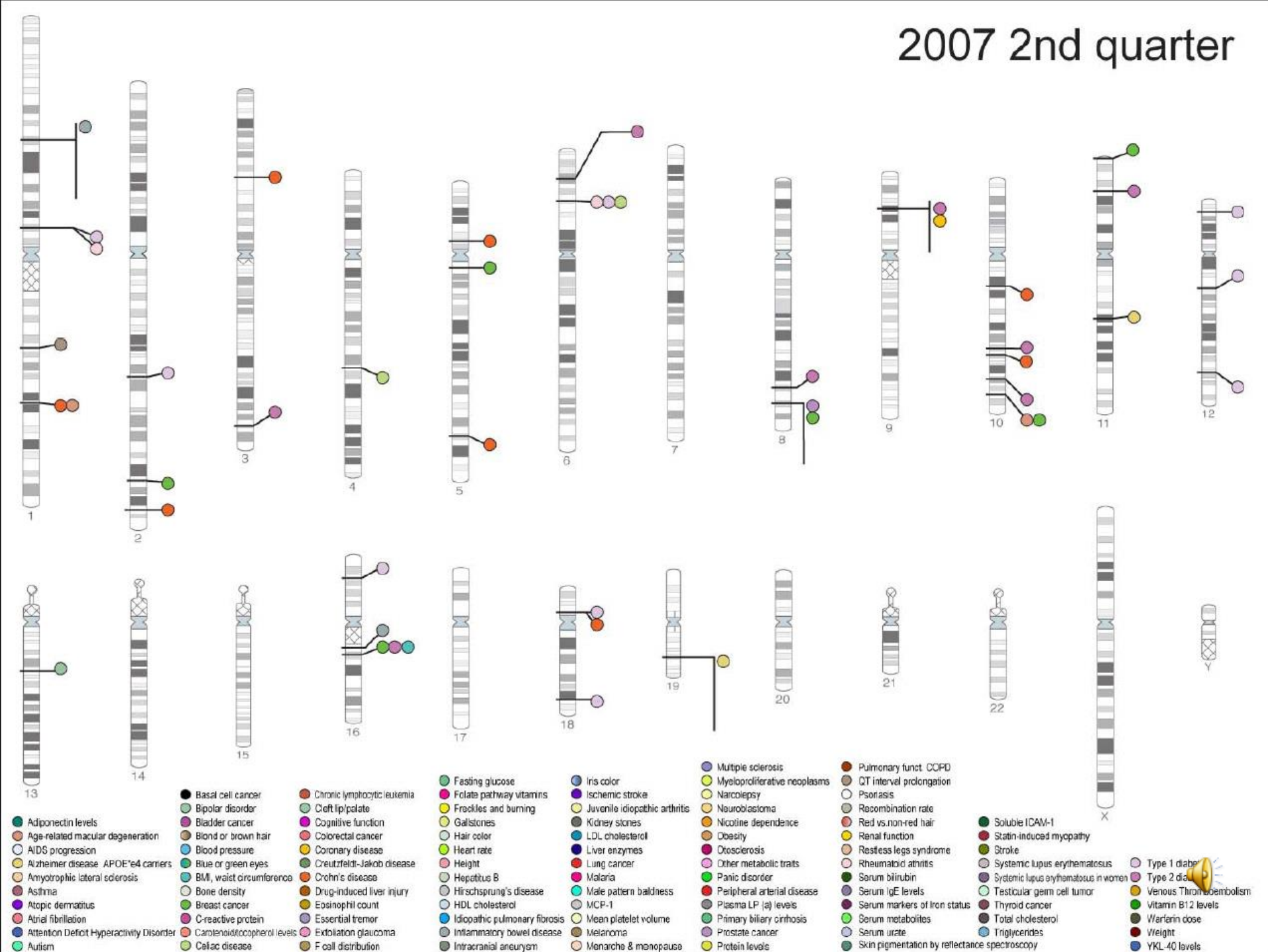


Bioinformatics Data

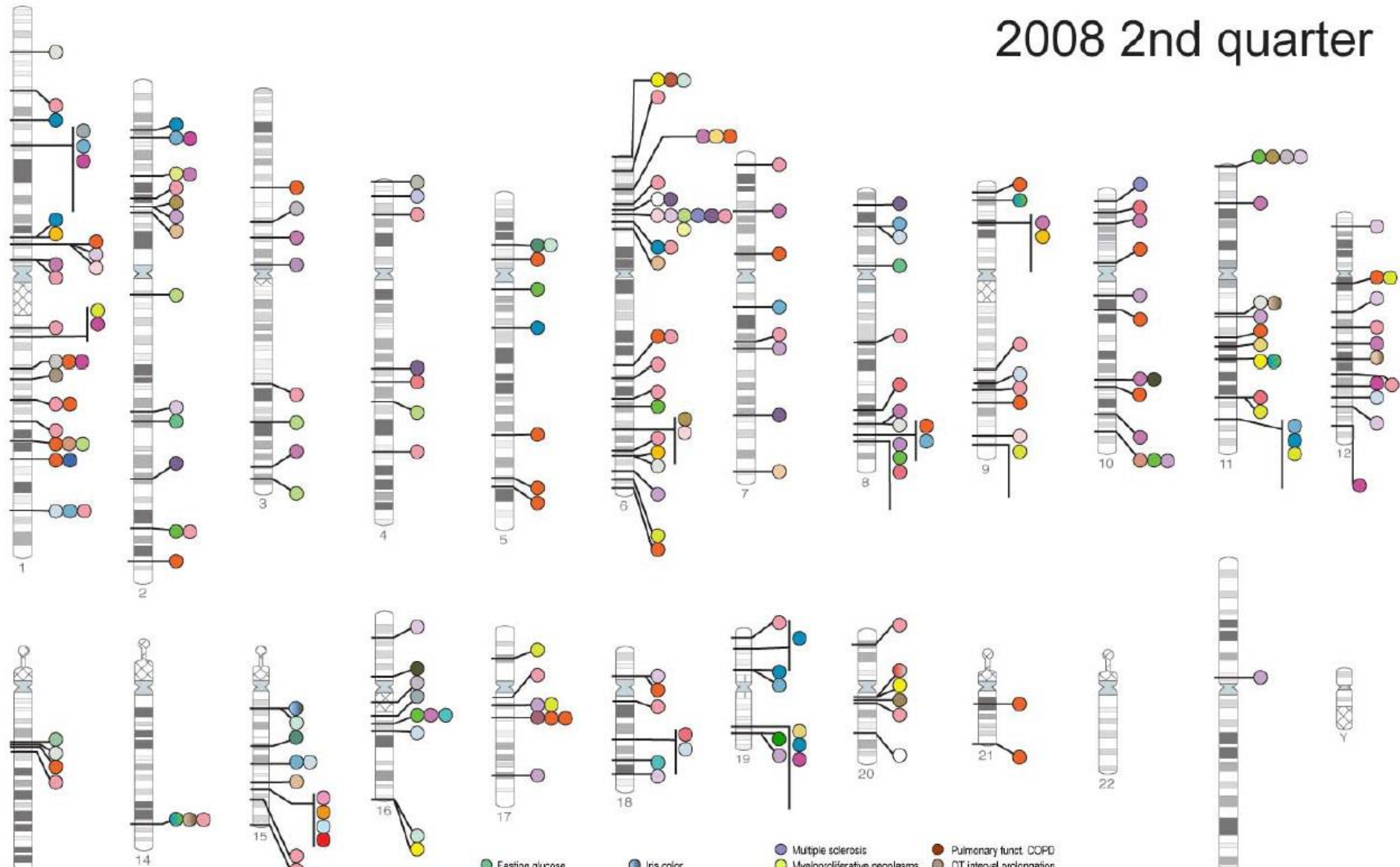
- Bioinformatics deals with any type of data that is of interest to biologists
 - DNA and protein sequences
 - Gene expression (microarray)
 - Raw data collected from field or laboratory experiment
 - Images, virtual models, Software
 - Articles from literature and databases of citations
- Each type of data can exist in many incompatible computer formats
- The analysis of DNA sequence data has come to dominate the field of bioinformatics, but the term can be applied to any type of biological data that can be recorded as numbers or images and handled by computers



2007 2nd quarter

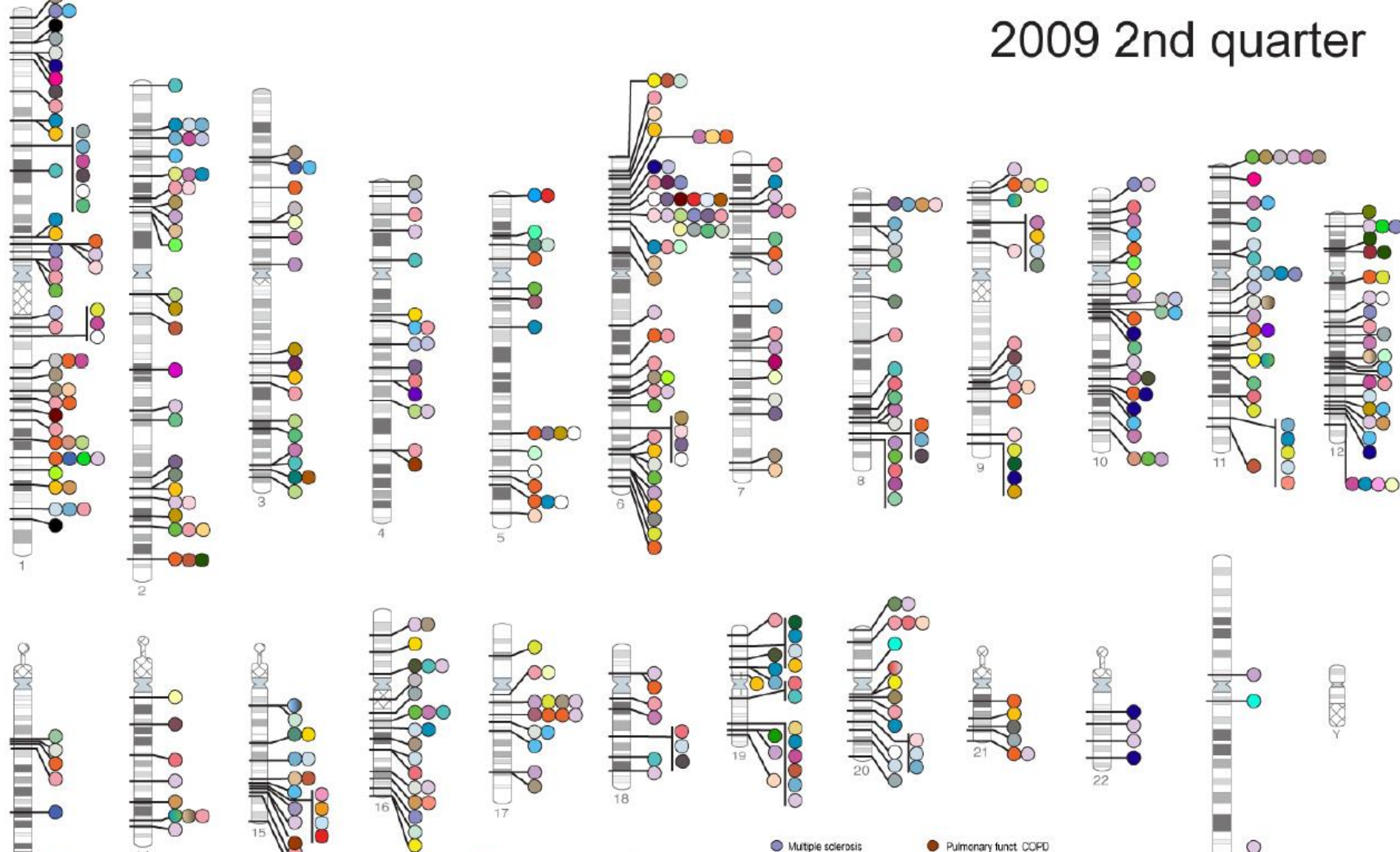


2008 2nd quarter

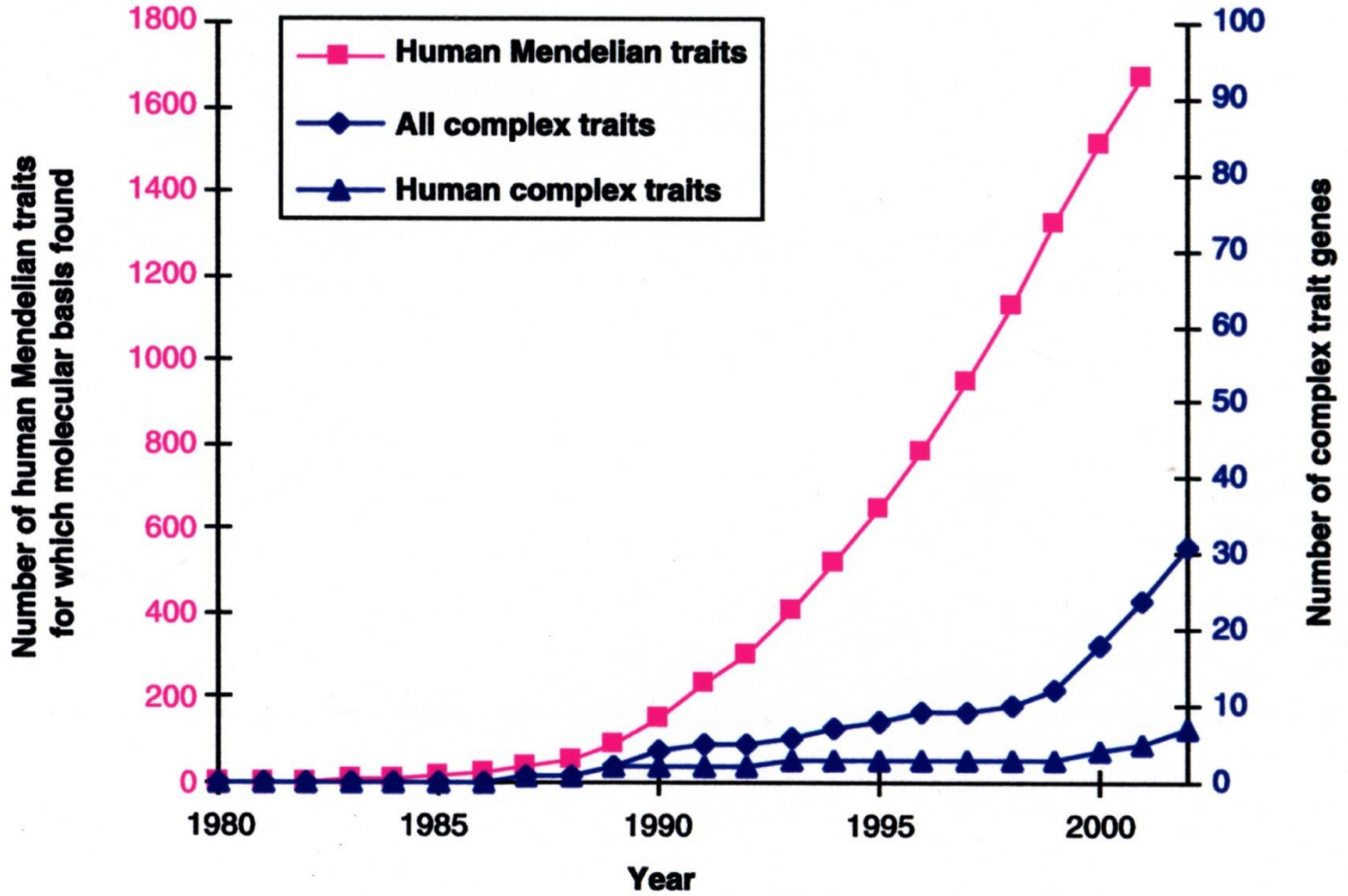


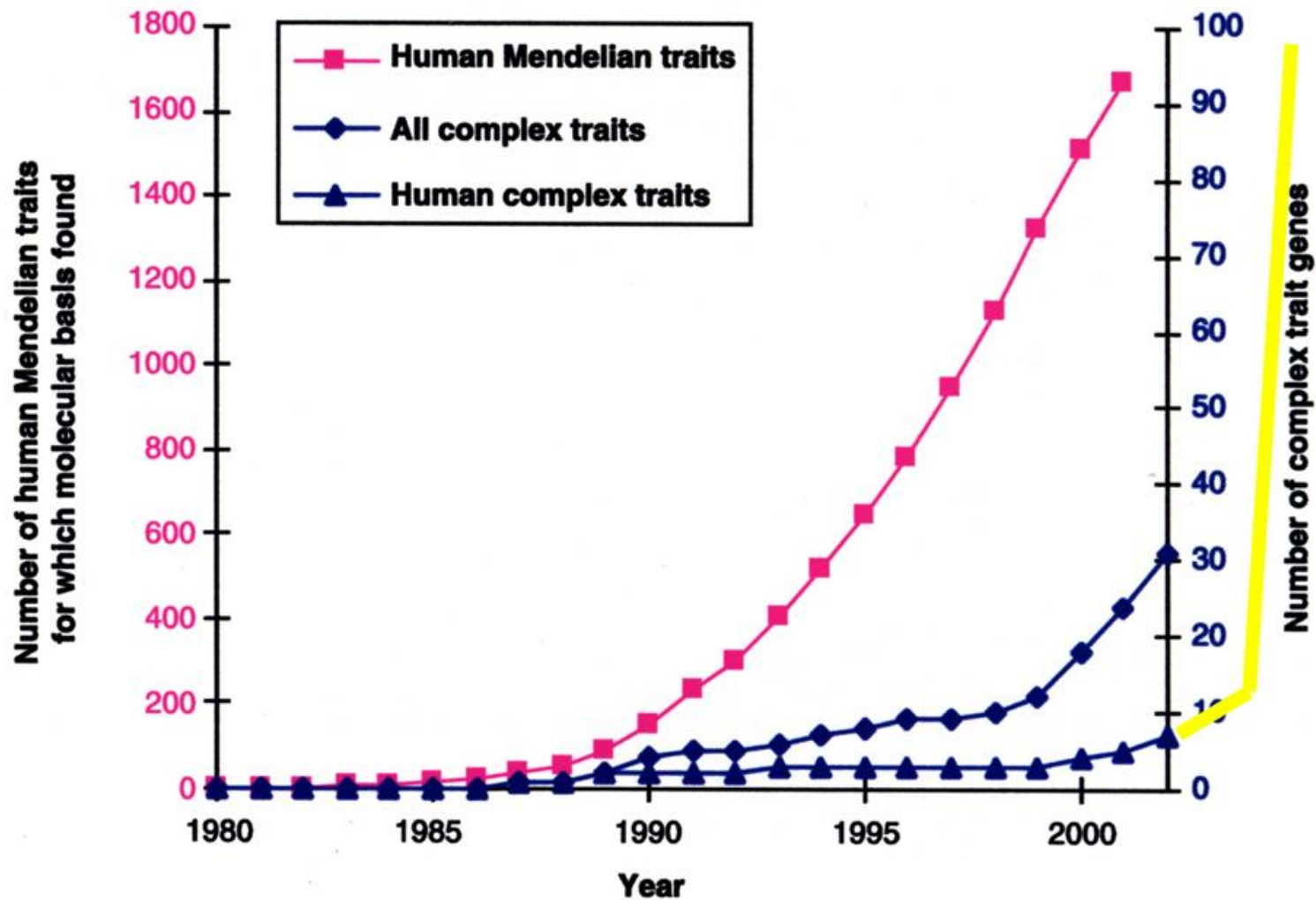
- | | | | | | | |
|---|------------------------------|--------------------------------|---------------------------------|---------------------------------|--------------------------------|---|
| ● Adiponectin levels | ● Basal cell cancer | ● Chronic lymphocytic leukemia | ● Fasting glucose | ● Iris color | ● Multiple sclerosis | ● Pulmonary funct. COPD |
| ● Age-related macular degeneration | ● Bipolar disorder | ● Cleft lip/palate | ● Folate pathway vitamins | ● Ischemic stroke | ● Myeloproliferative neoplasms | ● QT interval prolongation |
| ● AIDS progression | ● Bladder cancer | ● Cognitive function | ● Freckles and burning | ● Juvenile idiopathic arthritis | ● Narcolepsy | ● Psoriasis |
| ● Alzheimer disease. APOE ϵ 4 carriers | ● Blond or brown hair | ● Colorectal cancer | ● Hair color | ● Kidney stones | ● Neuroblastoma | ● Recombination rate |
| ● Amyotrophic lateral sclerosis | ● Blood pressure | ● Coronary disease | ● Heart rate | ● LDL cholesterol | ● Nicotine dependence | ● Red vs. non-red hair |
| ● Asthma | ● Blue or green eyes | ● Creutzfeldt-Jakob disease | ● Height | ● Liver enzymes | ● Obesity | ● Renal function |
| ● Atopic dermatitis | ● BMI, waist circumference | ● Crohn's disease | ● Hepatitis B | ● Lung cancer | ● Osteoarthritis | ● Restless legs syndrome |
| ● Atrial fibrillation | ● Bone density | ● Drug-induced liver injury | ● Hirschsprung's disease | ● Malaria | ● Other metabolic traits | ● Rheumatoid arthritis |
| ● Attention Deficit Hyperactivity Disorder | ● Breast cancer | ● Eosinophil count | ● HDL cholesterol | ● Male pattern baldness | ● Panic disorder | ● Serum bilirubin |
| ● Autism | ● C-reactive protein | ● Essential tremor | ● Idiopathic pulmonary fibrosis | ● Mean platelet volume | ● Peripheral arterial disease | ● Serum IgE levels |
| | ● Carotenoid/retinoid levels | ● Exfoliation glaucoma | ● Inflammatory bowel disease | ● Melanoma | ● Plasma LP (a) levels | ● Serum markers of iron status |
| | ● Celiac disease | ● F cell distribution | ● Intracranial aneurysm | ● Menarche & menopause | ● Primary biliary cirrhosis | ● Thyroid cancer |
| | | | | | ● Prostate cancer | ● Total cholesterol |
| | | | | | ● Protein levels | ● Triglycerides |
| | | | | | | ● Skin pigmentation by reflectance spectroscopy |
| | | | | | | ● Soluble ICAM-1 |
| | | | | | | ● Statin-induced myopathy |
| | | | | | | ● Stroke |
| | | | | | | ● Systemic lupus erythematosus |
| | | | | | | ● Systemic lupus erythematosus in women |
| | | | | | | ● Testicular germ cell tumor |
| | | | | | | ● Type 1 diabetes |
| | | | | | | ● Type 2 diabetes |
| | | | | | | ● Venous Thrombembolism |
| | | | | | | ● Vitamin B12 levels |
| | | | | | | ● Warfarin dose |
| | | | | | | ● Weight |
| | | | | | | ● YKL-40 levels |

2009 2nd quarter



- | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|----------------------|---------------------|--------------------|------------------|-----------------------|------------------|----------------------|----------------------------|----------------|-----------------|----------------------|----------------------------------|------------------|--------------------------------|--------------------|----------------------|---------------------|--------------------|-----------------------------|-------------------|-----------------------------|--------------------|--------------------|------------------------|-----------------------|-------------------|---------------------------|------------------------|--------------|--------------|--------------|----------|---------------|--------------------------|-------------------|---------------------------------|------------------------------|-------------------------|--------------|-------------------|---------------------------------|-----------------|-------------------|-----------------|-----------|-------------------------|---------|------------------------|------------|------------------------|----------------------|--------------------------------|--------------|-----------------|-----------------------|-----------|---------------|--------------------------|------------------|-------------------------------|------------------------|-----------------------------|-------------------|------------------|-------------------------|----------------------------|-------------|----------------------|------------------------|------------------|--------------------------|------------------------|-------------------|--------------------|--------------------------------|---------------|---|------------------|---------------------------|----------|--------------------------------|---|------------------------------|------------------|---------------------|-----------------|-------------------|-------------------|--------------------------|----------------------|-----------------|----------|-----------------|
| ● Adiponectin levels | ● Basal cell cancer | ● Bipolar disorder | ● Bladder cancer | ● Blond or brown hair | ● Blood pressure | ● Blue or green eyes | ● BMI, waist circumference | ● Bone density | ● Breast cancer | ● C-reactive protein | ● Carotenoid/retinophenol levels | ● Celiac disease | ● Chronic lymphocytic leukemia | ● Cleft lip/palate | ● Cognitive function | ● Colorectal cancer | ● Coronary disease | ● Creutzfeldt-Jakob disease | ● Crohn's disease | ● Drug-induced liver injury | ● Eosinophil count | ● Essential tremor | ● Exfoliation glaucoma | ● F cell distribution | ● Fasting glucose | ● Folate pathway vitamins | ● Freckles and burning | ● Gallstones | ● Hair color | ● Heart rate | ● Height | ● Hepatitis B | ● Hirschsprung's disease | ● HDL cholesterol | ● Idiopathic pulmonary fibrosis | ● Inflammatory bowel disease | ● Intracranial aneurysm | ● Iris color | ● Ischemic stroke | ● Juvenile idiopathic arthritis | ● Kidney stones | ● LDL cholesterol | ● Liver enzymes | ● Malaria | ● Male pattern baldness | ● MCP-1 | ● Mean platelet volume | ● Melanoma | ● Menarche & menopause | ● Multiple sclerosis | ● Myeloproliferative neoplasms | ● Narcolepsy | ● Neuroblastoma | ● Nicotine dependence | ● Obesity | ● Otsclerosis | ● Other metabolic traits | ● Panic disorder | ● Peripheral arterial disease | ● Plasma LP (a) levels | ● Primary biliary cirrhosis | ● Prostate cancer | ● Protein levels | ● Pulmonary funct. COPD | ● QT interval prolongation | ● Psoriasis | ● Recombination rate | ● Red vs. non-red hair | ● Renal function | ● Restless legs syndrome | ● Rheumatoid arthritis | ● Serum bilirubin | ● Serum IgE levels | ● Serum markers of iron status | ● Serum urate | ● Skin pigmentation by reflectance spectroscopy | ● Soluble ICAM-1 | ● Statin-induced myopathy | ● Stroke | ● Systemic lupus erythematosus | ● Systemic lupus erythematosus in women | ● Testicular germ cell tumor | ● Thyroid cancer | ● Total cholesterol | ● Triglycerides | ● Type 1 diabetes | ● Type 2 diabetes | ● Venous Thromboembolism | ● Vitamin B12 levels | ● Warfarin dose | ● Weight | ● YKL-40 levels |
|----------------------|---------------------|--------------------|------------------|-----------------------|------------------|----------------------|----------------------------|----------------|-----------------|----------------------|----------------------------------|------------------|--------------------------------|--------------------|----------------------|---------------------|--------------------|-----------------------------|-------------------|-----------------------------|--------------------|--------------------|------------------------|-----------------------|-------------------|---------------------------|------------------------|--------------|--------------|--------------|----------|---------------|--------------------------|-------------------|---------------------------------|------------------------------|-------------------------|--------------|-------------------|---------------------------------|-----------------|-------------------|-----------------|-----------|-------------------------|---------|------------------------|------------|------------------------|----------------------|--------------------------------|--------------|-----------------|-----------------------|-----------|---------------|--------------------------|------------------|-------------------------------|------------------------|-----------------------------|-------------------|------------------|-------------------------|----------------------------|-------------|----------------------|------------------------|------------------|--------------------------|------------------------|-------------------|--------------------|--------------------------------|---------------|---|------------------|---------------------------|----------|--------------------------------|---|------------------------------|------------------|---------------------|-----------------|-------------------|-------------------|--------------------------|----------------------|-----------------|----------|-----------------|



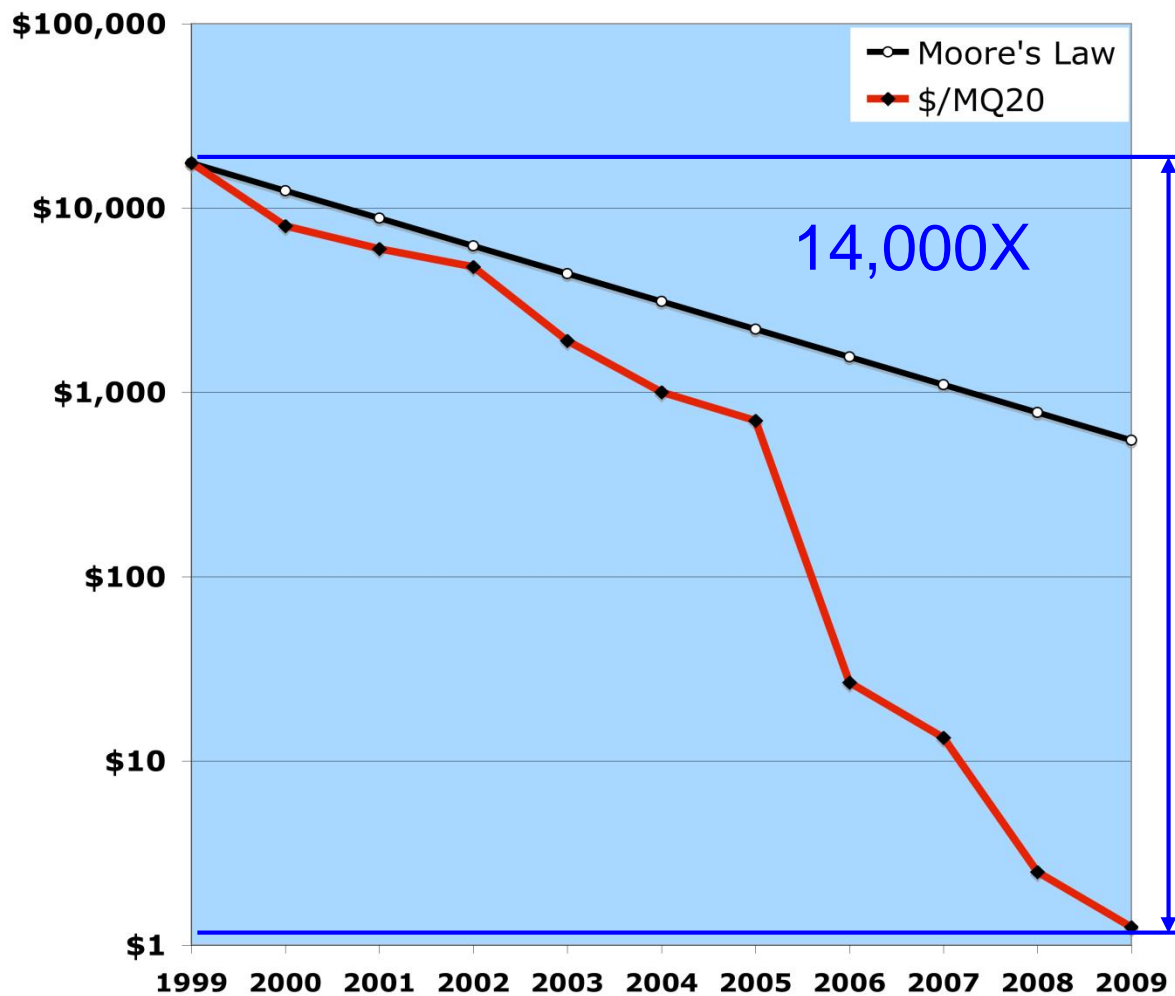


Glazier et al., Science 298:2345-9, 2002





Cost per Q20 Megabase (\$)





An information explosion...

- Lots of data in genome
- More data in when we attempt to
 - discern structure of data
 - relate to transcriptomics, proteomics
 - relate to structure, physiology
 - relate to disease
 - relate to variation
- Automated discovery, experiments
- Biomedical knowledge (coming)
- Clinical knowledge (coming)



[Some] Research Projects

- **The Human Genome Project** -- old news, 6 years ago
- **International HapMap Project** -- www.hapmap.org
- **The 1000 Genomes Project** – www.1000.genomes.org
- Encyclopedia of DNA Elements (**ENCODE**) Project
- **The Cancer Genome Atlas** (TCGA)
- **Human Microbiome Project** (HMP) – www.hmpdacc.org
- The **eMERGE** (Electronic Medical Records and Genomics) Network





Common Features of Projects

- High throughput
- Use of technology, in particular
 - Automation (Robotics, AI)
 - Databases
 - Visualization, simulation/computational models
 - Groupware: Coordination and communication
- Public domain tools
- Open sharing of data





Some Challenges

- Volume of data is staggering
 - How to store and collect sequence information?
 - RDBMSs don't handle sequence data well
 - Better handled by Object Oriented DBM
- How to analyze and display the data
 - Automated algorithms
 - Contextual visualization methods
 - Clusters, profiles, etc
- Sequence data is meaningless without context
 - Not well suited to printed medical record





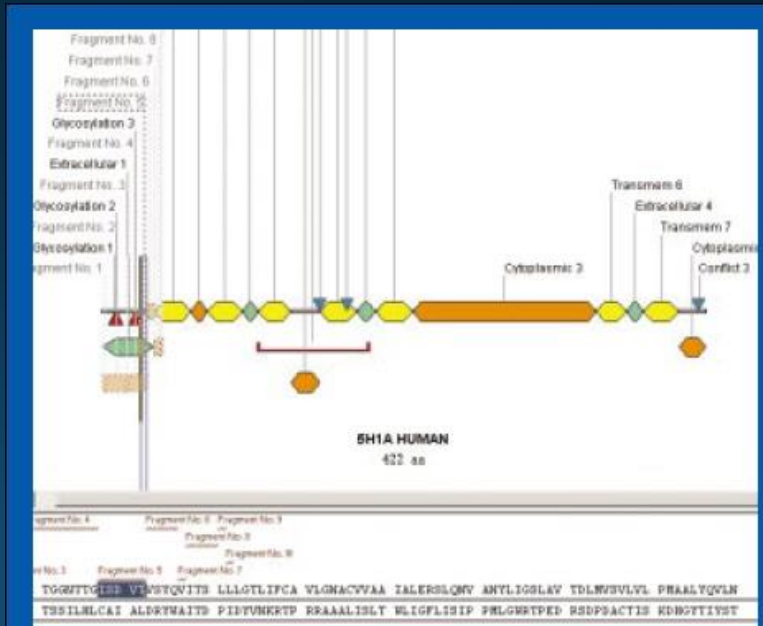
General Informatics Techniques/Tools in Bioinformatics

- Storage
 - Databases
 - Building, Querying
 - Complex data
 - Annotations
 - Citations
- Standards
- Interoperability
- Knowledge Management
 - Classification
 - Vocabularies
 - Ontologies
- Communications
- Process Workflow
- Discovery and Analyses
 - Text String Comparison
 - Text search
 - Statistical analysis
 - Finding Patterns
 - AI / Machine Learning
 - Clustering
 - Data mining
 - Geometric
 - Robotics
 - Graphics (Surfaces, Volumes)
 - Comparison and 3D Matching (Vision, Recognition)
 - Physical Simulation
 - Newtonian Mechanics
 - Electrostatics
 - Numerical Algorithms
 - Simulation



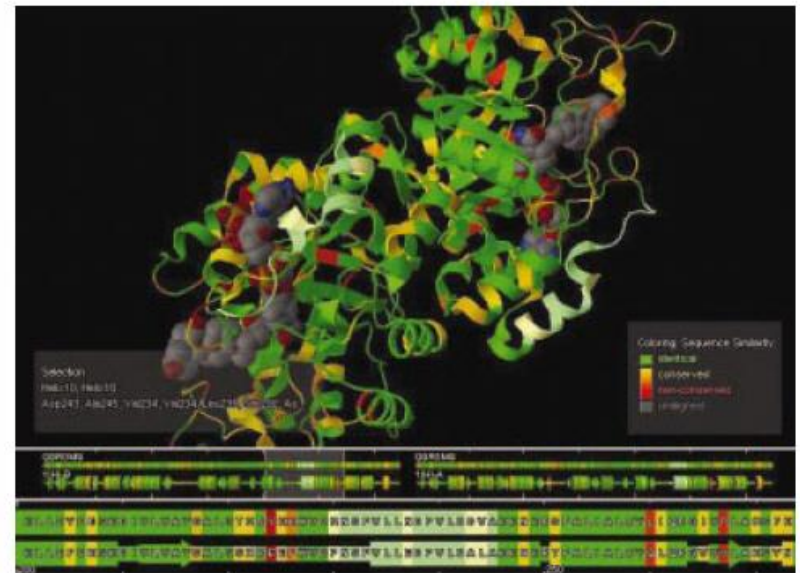
Bioinformatics: Tools

- Annotation



InforMax's BioAnnotator uses locally stored databases to find protein motifs.

- Visualization

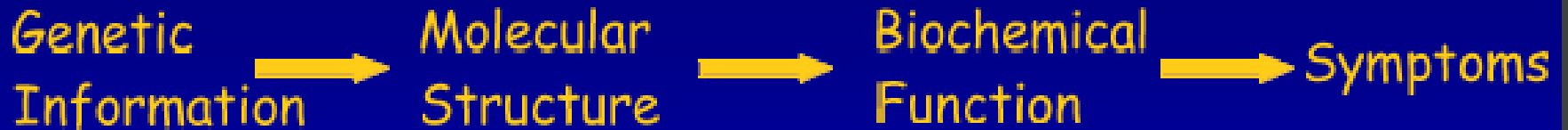


Structure prediction: modelling a sequence homolog in LION's SRS 3D.

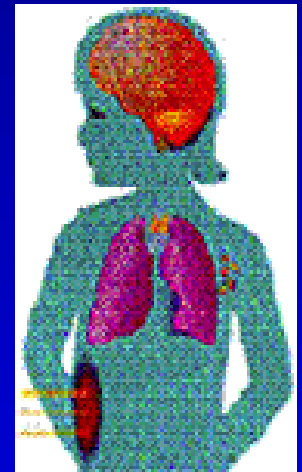
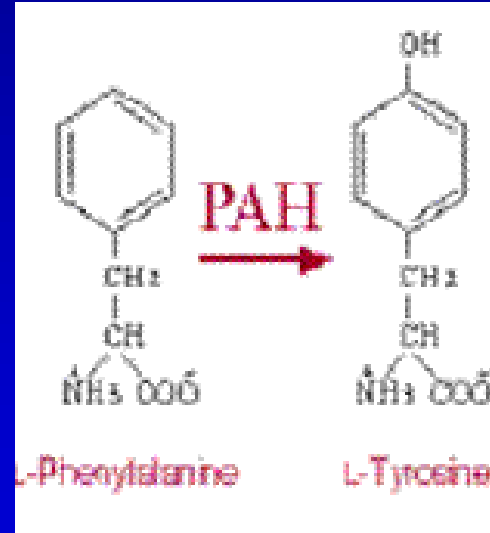
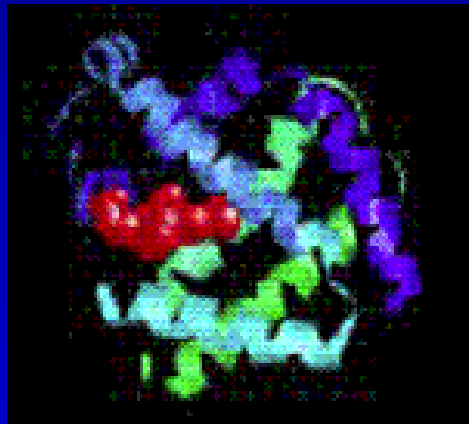
- user friendly, in the public domain, and increasingly integrated
- commercial tools streamline tasks, access proprietary databases



Central Paradigm of Bioinformatics



SRRAINXRIYA
VSYQTVSRVUN
VSTATVSEALA
GVTTTVSHVIN
SGVSAVSAIIN
GVSENRRLIN
TAYATIHVVE
GSQPTVSEELA
MSIATITRGSN
ISRETVGRILK
FDISRLSHLFR
LRPSRLAHLFR
MTVETISRLLG
TLEFHLHPLFK



Central Paradigm of Bioinformatics

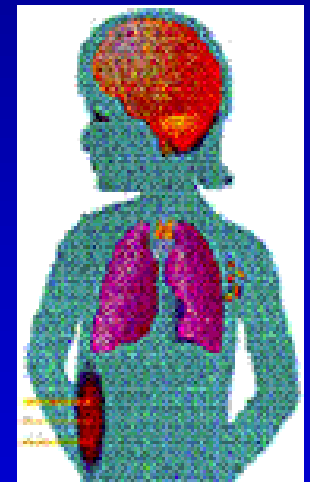
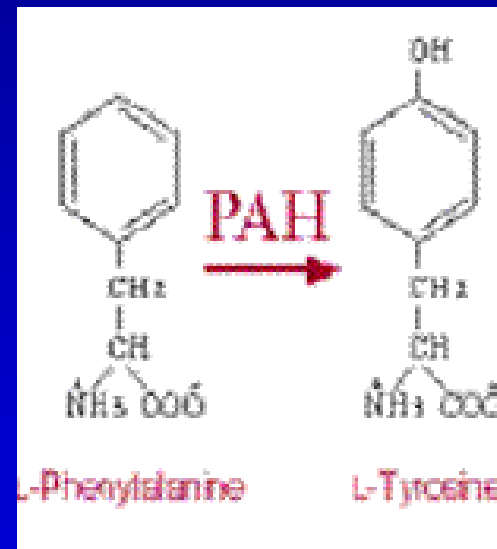
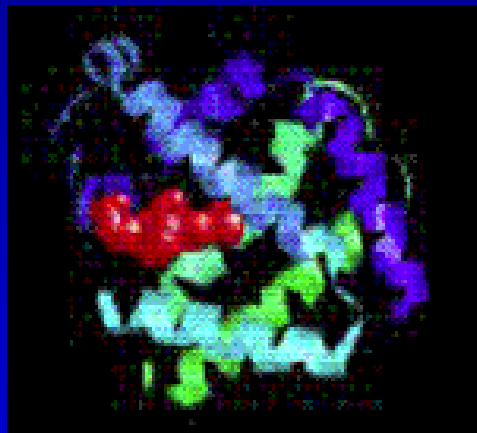
Genetic
Information

Molecular
Structure

Biochemical
Function

Symptoms

SRRAINKHIYA
VSYQTVSHVYN
VSTATVSRALA
GYPTTVSHVIN
SGVSAVSAILN
QVSEMTARDLN
TAYATEHVREVE
OSQPTVSRRLA
MSIATETROSN
ISRETVGRILK
FDISRLSHLFR
LRPSRLAHLFR
MTVETESALLG
TLEFELHPLFX



Cosa è la Bioinformatica?

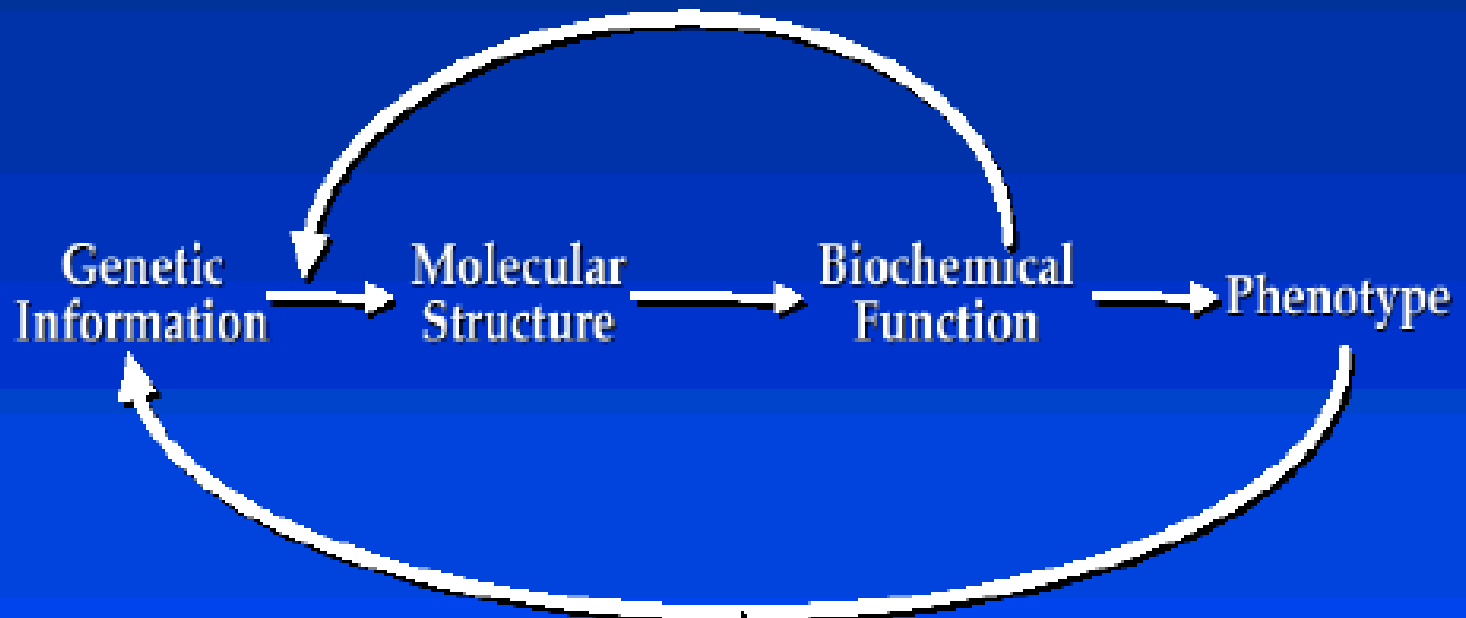
Un campo delle scienze in cui **Biologia, Informatica e Biotecnologie** confluiscono in un'unica disciplina che permetta:

- **Sviluppo di nuovi algoritmi e funzioni statistiche per individuare relazioni tra informazioni provenienti da diversi gruppi di dati**
- **Analisi e interpretazione di vari tipi di dati**
- **Sviluppo e implementazione di strumenti per accedere e gestire con efficienza diversi tipi di informazioni**

Perché usare la Bioinformatica

Crescita esponenziale di informazioni biologiche (sequenziamento automatico, DNA chips, identificazione proteine, spettrometria di massa ecc.) ha richiesto lo sviluppo di algoritmi che permettessero l'analisi di dati e la loro archiviazione (tools e database)

Central Paradigm of Bioinformatics



Perché è così difficile riuscire ad ottenere informazioni solo dalla sequenza:

L'informazione genetica è ridondante

Codice genetico

Sostituzione amminoacidica

Variazioni introni-esoni

Variazione del modulo di lettura

L'informazione strutturale è ridondante

Cambiamenti conformazionali

Differenti strutture possono mostrare similare funzione

Differenti sequenze danno simili strutture

Geni singoli possono avere diverse funzioni

Agire come enzimi metabolici e regolatori

I geni hanno 1D ma la funzione dipende dalla struttura 3D

Search

GenBank



for

Go

NCBI Homepage

SITE MAP

About NCBI

general and contact information

GenBank

sequence submission support and software

Molecular databases

sequences, structures and taxonomy

Literature databases

PubMed, OMM and PubMed Central

NEW

Genomic biology

the human genome, whole genomes and related resources

Tools

for data mining

Research at NCBI

What does NCBI do?

Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease.

Draft Human Genome

Explore [human genome resources](#) or browse the human genome sequence using the [Map Viewer](#).

DART: A new tool

Archives

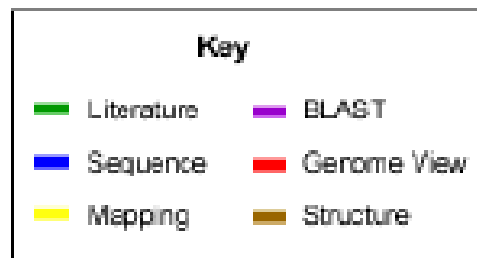
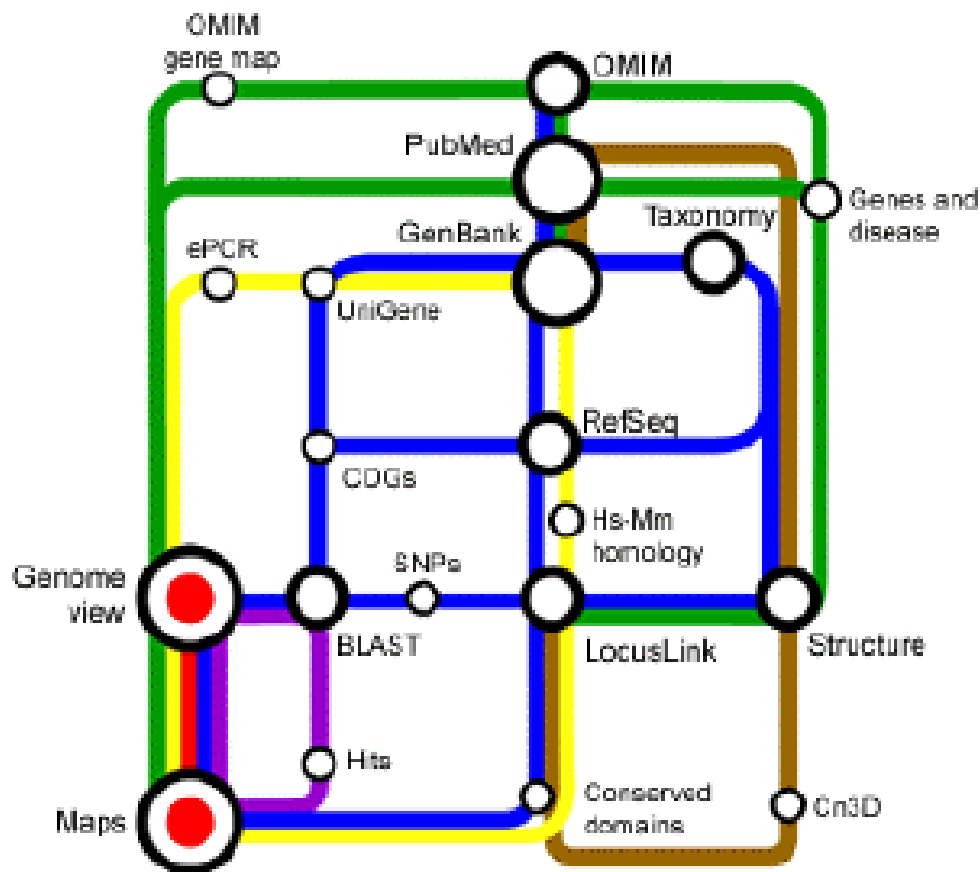


Want to locate protein neighbors by domain architecture? Learn about NCBI's new Domain Architecture Retrieval Tool...

<http://www.ncbi.nlm.nih.gov/>

Hot Spots

- ▶ Cancer genome anatomy project
- ▶ Clusters of orthologous groups
- ▶ Coffee Break
- ▶ Electronic PCR
- ▶ Gene expression omnibus
- ▶ Genes and disease
- ▶ Human genome resources
- ▶ Human map viewer
- ▶ Human/mouse homology maps

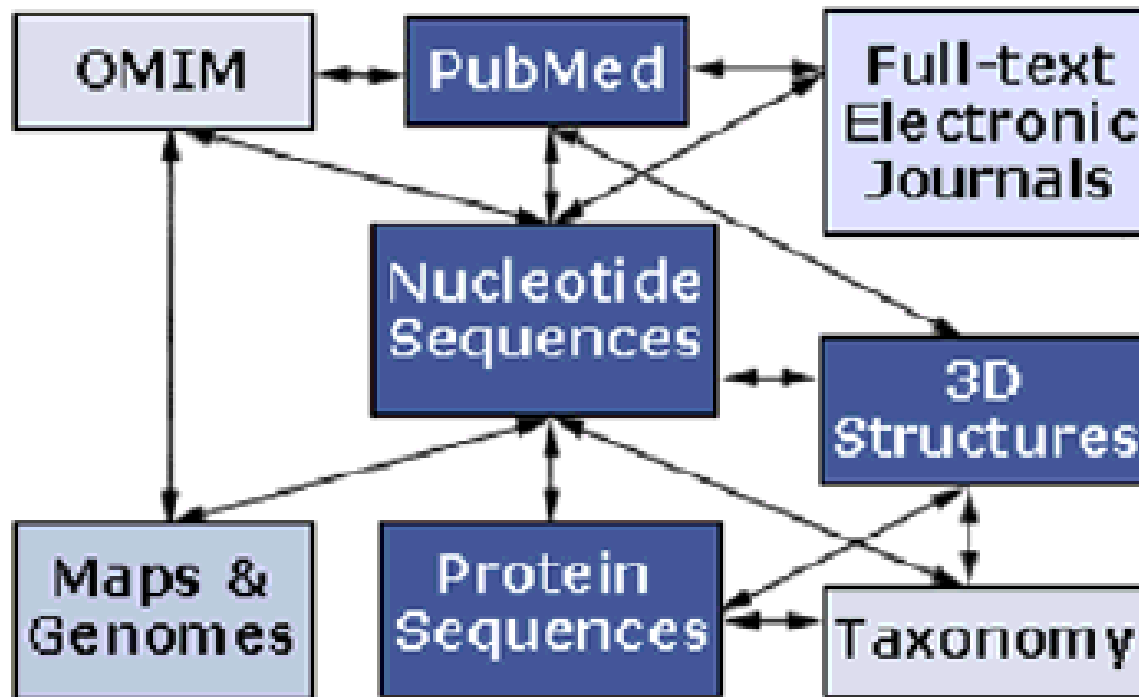


[Back to showcase](#)

<http://www.ncbi.nlm.nih.gov/Tour/tour.html>

ENTREZ

A search and retrieval system for information integration.



Entrez is a retrieval system for searching several linked databases.

It provides access to:

MedLine, abstracts
and links to journals

PubMed: The biomedical literature (PubMed)

GenBank, EMBL, DDBJ,
RefSeq

Nucleotide: Sequence databases (GenBank)

GenBank, DDBJ, EMBL, PDB,
PIR, Swiss-prot, RefSeq

Protein: Sequence databases

Structure: Three-dimensional macromolecular structures

Genome: Complete genome assemblies

NCBI's MMDB-
derived from PDB

PopSet: Population study data sets

Graphical view,
mapping data

Taxonomy: Organisms in GenBank

OMIM: Online Mendelian Inheritance in Man

Population and
phylogenetic studies

MIM in Entrez

NCBI's Taxonomy,
tree structures

<http://www.ncbi.nlm.nih.gov/Entrez/>



National Center for Biotechnology Information

National Library of Medicine

National Institutes of Health

PubMed

Entrez

BLAST

OMIM

Taxonomy

Structure

Search

GenBank

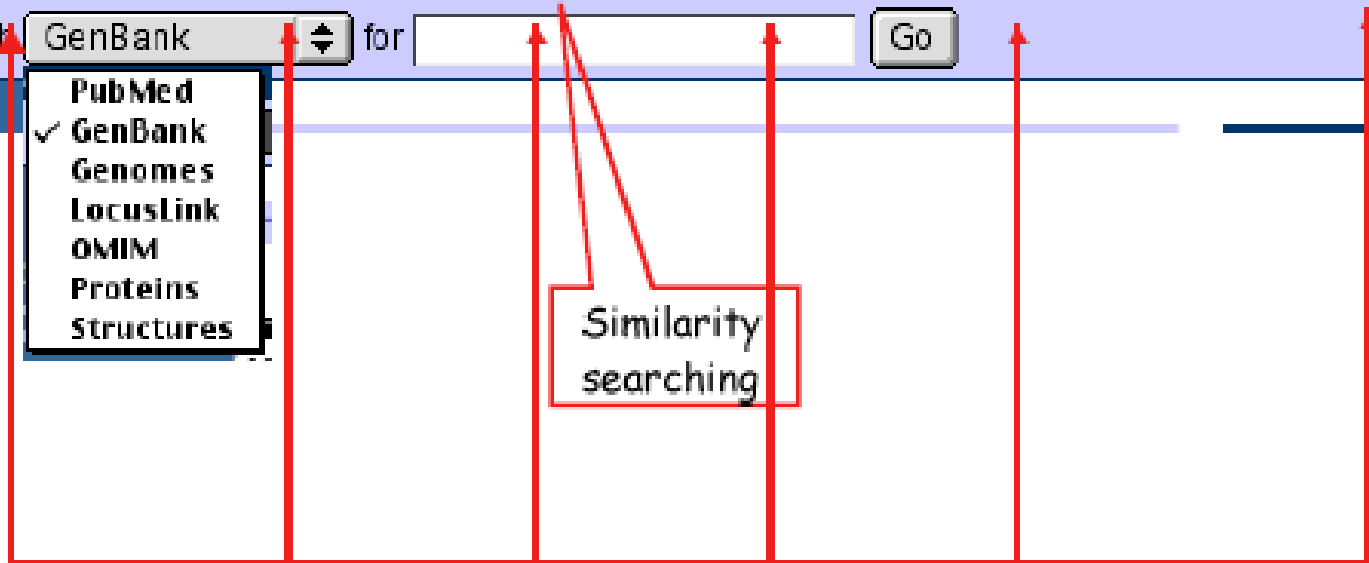
for

Go

- PubMed
- ✓ GenBank
- Genomes
- LocusLink
- OMIM
- Proteins
- Structures

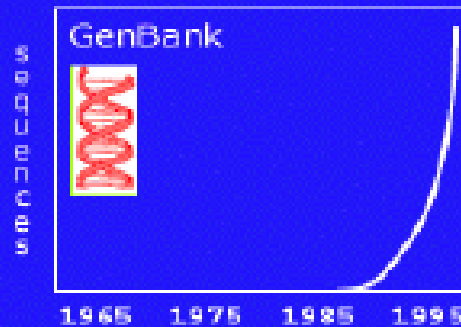
Similarity
searching

NCBI

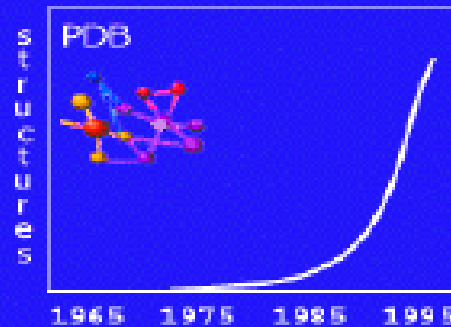




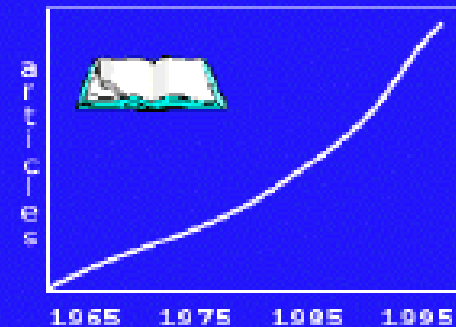
Nucleotide sequences



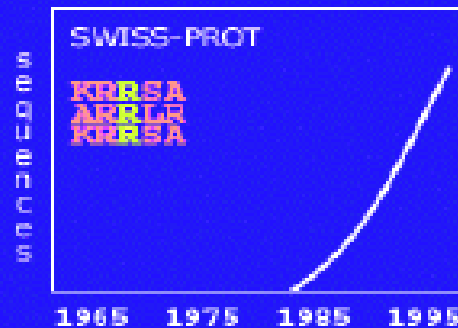
Protein structures



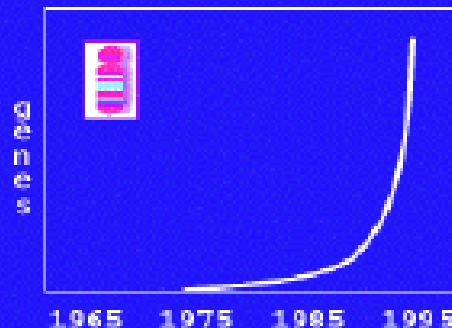
Bibliographic



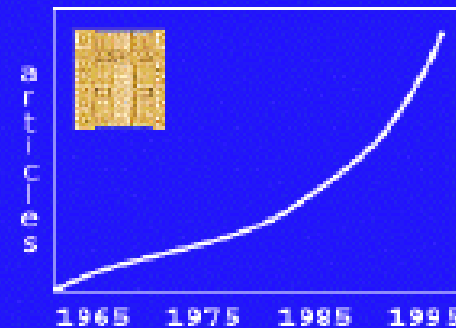
Protein sequences



Mapped human gene



Genetic bibliography



- Exponential growth of biological information: growth of **sequences**, **structures**, and **literature**.
- Efficient **storage** and **management tools** were most important.

NCBI bioinformatics tools - 1-

PubMed Entrez **BLAST** OMIM Taxonomy Structure

Search for

BLAST

The **Basic Local Alignment Search Tool (BLAST)**, for comparing gene and protein sequences against others in public databases, now comes in several

flavors including [PSI-BLAST](#), [PHI-BLAST](#), and [BLAST 2 sequences](#). Specialized BLASTs are also available for [human](#), [microbial](#), and [malaria](#) genomes, as well as for [vector contamination](#), [immunoglobulins](#), and [tentative human consensus sequences](#).

[BLAST](#)

[COGs](#)



Clusters of Orthologous Groups (COGs) currently covers 21 complete genomes from 17 major phylogenetic lineages. A COG is a cluster of very similar proteins found in at least three species.

The presence or absence of a protein in different genomes can tell us about the evolution of the organisms, as well as point to new drug targets.

[Map Viewer](#)

[LocusLink](#)



Map Viewer shows integrated views of chromosome maps

currently for human, mouse, and *Drosophila*. Used to view the NCBI assembly of the complete human genome, Map Viewer is a valuable tool for the identification and localization of genes that contribute to human disease.



LocusLink

combines descriptive and sequence information on human genes through a single query interface. LocusLink covers information on official nomenclature, aliases, sequence accession numbers, phenotypes, EC numbers, OMIM numbers, UniGene clusters, map information, and relevant web sites.

[UniGene](#)

[ORF finder](#)

[Electronic PCR](#)

NCBI bioinformatics tools -2-

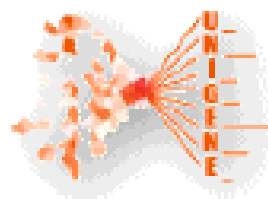
[VAST search](#)

[CCAP](#)

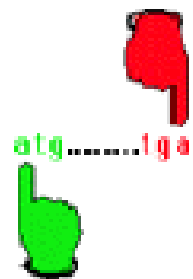
[Human-Mouse](#)

[VecScreen](#)

[CGAP](#)



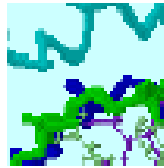
A [UniGene](#) cluster is a non-redundant set of sequences that represents a unique human, mouse, or rat gene. Well-characterized genes, as well as thousands of expressed sequence tag (EST) sequences have been included. Each cluster record also contains information such as the tissue types in which the gene has been expressed and map location. UniGene can assist in gene discovery, gene mapping projects, and large-scale expression analysis.



[ORF finder](#) identifies all possible ORFs in a DNA sequence by locating the standard and alternative stop and start codons. The deduced amino acid sequences can then be used to BLAST against GenBank. ORF finder is also packaged in the sequence submission software [Sequin](#).

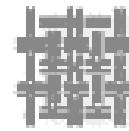
electronic
PCR
001101011AGCGT

[Electronic PCR](#) allows you to search your DNA sequence for sequence tagged sites (STSs), which have been used as landmarks in various types of genomic maps. It compares the query sequence against data in NCBI's [UnSTS](#), a unified, non-redundant view of STSs from a wide range of sources.



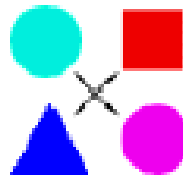
VAST search is a structure-structure similarity search service. It compares 3D coordinates of a newly determined

protein structure to those in the MMDB/PDB database. VAST Search computes a list of similar structures that can be browsed interactively, using molecular graphics to view superimpositions and alignments.



The Cancer Chromosome Aberration Project (CCAP) compiles information on the distinct chromosome

aberrations that are associated with different cancers. The identification of chromosomal abnormalities by clinicians can enable the diagnosis of, classification of, and treatment selection for a given cancer.



The **Human-Mouse Homology Maps** compare genes in homologous segments of DNA from human and mouse sources, sorted by position in each genome. A total of 1793 loci are presented, most of which are genes. This map should be interpreted as a reflection of probable, not confirmed, homology relationships due to the lack of further information available for about half the loci.



VecScreen is a tool for identifying segments of a nucleic acid sequence that may be of vector, linker or adapter origin

prior to sequence analysis or submission. VecScreen was developed to combat the problem of vector contamination in public sequence databases.



The Cancer Genome Anatomy Project

(CGAP) aims to decipher the molecular anatomy of cancer cells. CGAP develops profiles of cancer cells by comparing gene expression in normal, precancerous, and malignant cells from a wide variety of tissues.

Introduzione

- Esistono diverse banche dati in cui sono depositate sequenze biologiche
 - di nucleotidi (DNA)
 - ✓ GenBank (National Institute of Health, Bethesda, MD)
 - ✓ EMBL (European Bioinformatics Institute, Hinxton, UK)
 - ✓ DNA Database of Japan (Mishima, Japan)
 - di amminoacidi (proteine)
 - ✓ PIR (National Biomedical Research Foundation, Washington, DC)
 - ✓ SwissProt/TrEMBL (Università di Ginevra, CH)
 - ✓ PRFDB (Protein Research Foundation, Osaka, Japan)
- Esiste la necessità di evidenziare, tramite confronto, la similarità tra una sequenza studiata e quelle contenute in una banca dati

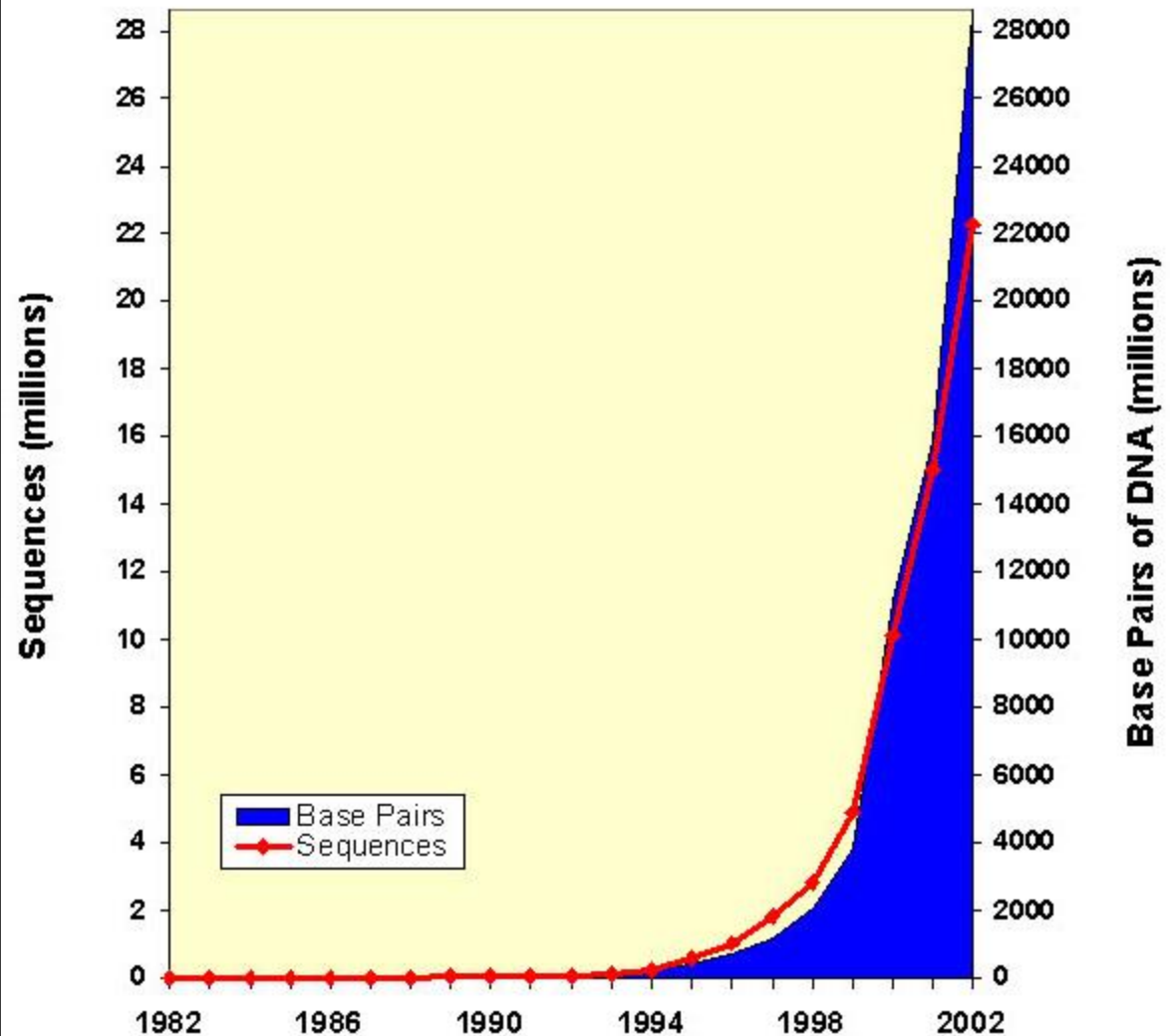
BANCHE DATI DI SEQUENZE GENOMICHE

*GenBank deriva dalla
collaborazione di diversi database
di sequenze tra cui EMBL e DDBJ*

Growth of GenBank

What is GenBank?

GenBank® is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (Nucleic Acids Research 2002 Jan 1;30(1):17-20). There are approximately 22,617,000,000 bases in 18,197,000 sequence records as of August 2002 (see GenBank growth statistics). As an example, you may view the record for a *Saccharomyces cerevisiae* gene. The complete release notes for the current version of GenBank are available. A new release is made every two months. GenBank is part of the International Nucleotide Sequence Database Collaboration, which comprises the DNA DataBank of Japan (DDBJ), the European Molecular Biology Laboratory (EMBL), and GenBank at NCBI. These three organizations exchange data on a daily basis.





NCBI

Submit to GenBank

PubMed

Entrez

BLAST

OMIM

Books

TaxBrowser

Structure

Search

Entrez

▼ for

Go

NCBI

[SITE MAP](#)

Guide to NCBI resources

[Accession numbers](#)

For manuscript citation

[BankIt](#)

[Sequin](#)

[SequinMacroSend](#)

Upload .sqn files directly

[TBL2ASN](#)

Command line program

[Special submissions](#)

▶ **Submitting Sequence Data to GenBank**

The most important source of new data for GenBank® is direct submissions from scientists. GenBank depends on its contributors to help keep the database as comprehensive, current, and accurate as possible. NCBI provides timely and accurate processing and biological review of new entries and updates to existing entries, and is ready to assist authors who have new data to submit.

▶ **Receiving an Accession Number for your Manuscript**

Most journals now expect that DNA and amino acid sequences that appear in articles will be submitted to a sequence database before publication. Soon after submission, you will receive an accession number from the

▶ **Submit now!!**

[Sequin](#)

Stand-alone sequence submission tool

[BankIt](#)

For quick and simple submissions

[VecScreen](#)

Vector contamination screening tool

[dbEST](#)

[dbGSS](#)

[dbSTS](#)

Submit to GenBank divisions

▶ **GenBank**

[GenBank](#)

Overview of the

BANCHE DATI DI SEQUENZE PROTEICHE

SWISS-PROT

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

TrEMBL

The TrEMBL database contains the translations of all coding sequences (CDS) present in the EMBL Nucleotide Sequence Database, which are not yet integrated into Swiss-Prot.

SP-TrEMBL (Swiss-Prot TrEMBL) Contains the entries which should eventually be incorporated into Swiss-Prot and can be considered as a preliminary section of Swiss-Prot as all SP-TrEMBL entries have been assigned Swiss-Prot accession numbers.

REM-TrEMBL (REMaining TrEMBL) Contains the entries that we do not want to include in Swiss-Prot. REM-TrEMBL entries have no accession numbers.

[Site Map](#)[Search ExPASy](#)[Contact us](#)Search for 

ExPASy Proteomics Server

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

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

Databases

- [Swiss-Prot and TrEMBL](#) - Protein knowledgebase
- [PROSITE](#) - Protein families and domains
- [SWISS-2DPAGE](#) - Two-dimensional polyacrylamide gel electrophoresis
- [ENZYME](#) - Enzyme nomenclature
- [SWISS-3DIMAGE](#) - 3D images of proteins and other biological macromolecules
- [SWISS-MODEL Repository](#) - Automatically generated protein models

- [GermOnLine](#) - Knowledgebase on germ cell differentiation
- [Ashbya Genome Database](#)
- [Links to many other molecular biology databases](#)

Tools and software packages

- [Proteomics and sequence analysis tools](#)
 - ◊ [Proteomics](#) [[PeptIdent](#), [PeptideMass](#), ...]
 - ◊ [DNA -> Protein](#) [[Translate](#)]
 - ◊ [Similarity searches](#) [[BLAST](#)]
 - ◊ [Pattern and profile searches](#) [[ScanProsite](#)]
 - ◊ [Post-translational modification and topology prediction](#)
 - ◊ [Primary structure analysis](#) [[ProtParam](#), [pI/MW](#), [ProtScale](#)]
 - ◊ [Secondary and tertiary structure prediction](#) [[SWISS-MODEL](#), [Swiss-PdbViewer](#)]
 - ◊ [Alignment](#) [[T-COFFEE](#), [SIM](#)]
 - ◊ [Biological text analysis](#)
- [ImageMaster / Melanie](#) - Software for 2-D PAGE analysis
- [Roche Applied Science's Biochemical Pathways](#)

Education and services

- [The ExPASy FTP server](#)

Documentation

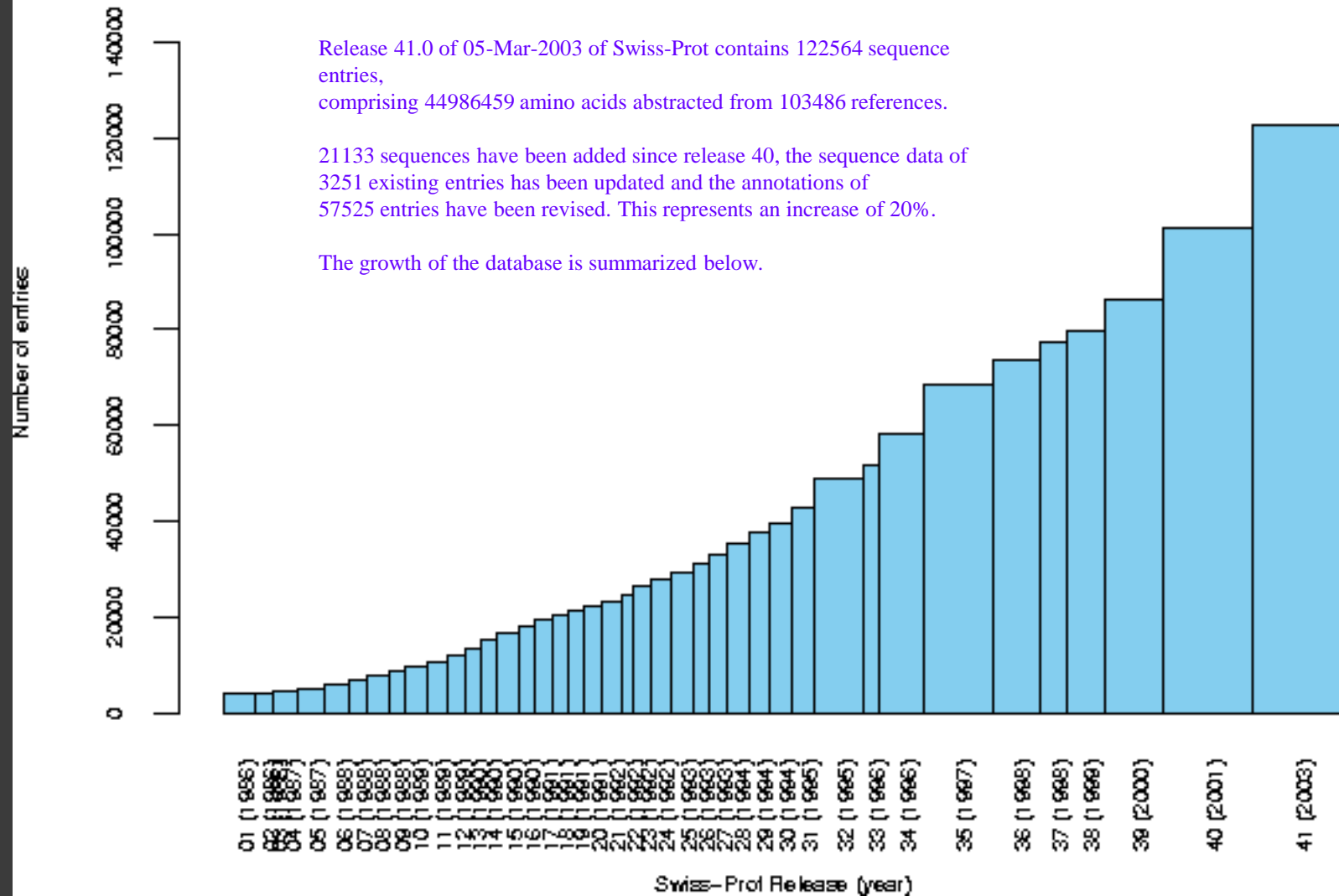
- [What's New on ExPASy](#)

Size of the Swiss-Prot database

Release 41.0 of 05-Mar-2003 of Swiss-Prot contains 122564 sequence entries, comprising 44986459 amino acids abstracted from 103486 references.

21133 sequences have been added since release 40, the sequence data of 3251 existing entries has been updated and the annotations of 57525 entries have been revised. This represents an increase of 20%.

The growth of the database is summarized below.

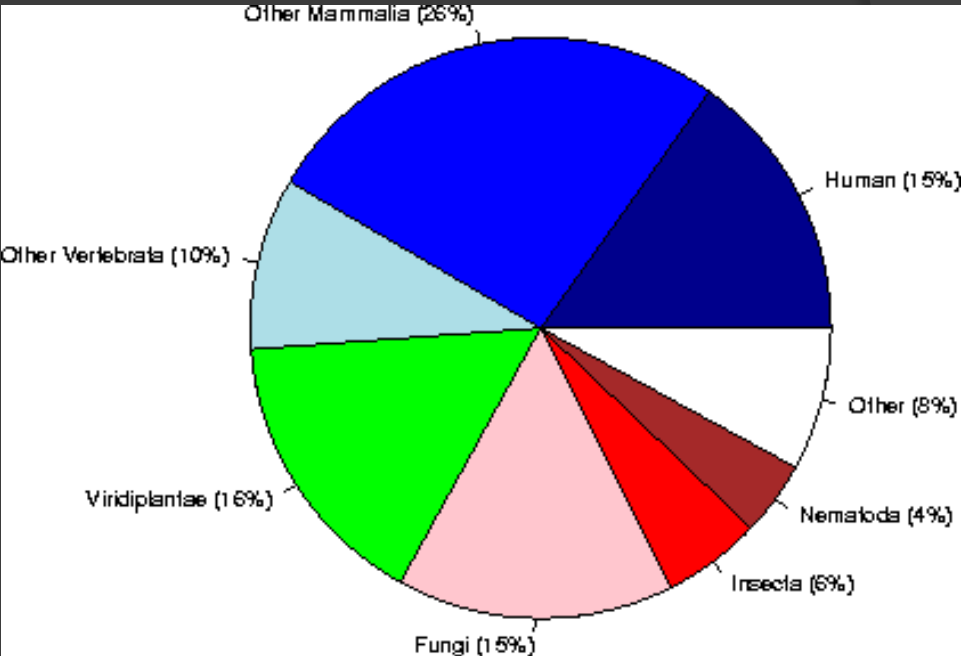
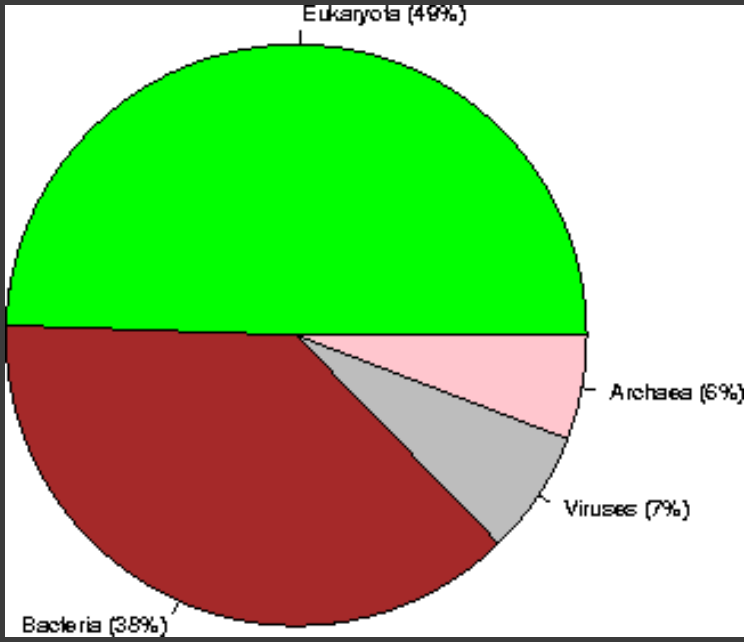


Taxonomic distribution of the sequences

Kingdom	sequences (% of the database)
Archaea	7119 (6%)
Bacteria	46344 (38%)
Eukaryota	60623 (49%)
Viruses	8478 (7%)

Within Eukaryota:

Category	sequences (% of Eukaryota)(% of the complete database)
Human	9172 (15%) (7%)
Other Mammalia	16041 (26%) (13%)
Other Vertebrata	5806 (10%) (5%)
Viridiplantae	9581 (16%) (8%)
Fungi	9337 (15%) (8%)
Insecta	3352 (6%) (3%)
Nematoda	2504 (4%) (2%)
Other	4830 (8%) (4%)



Public Databases:

- NCBI database
- PIR International NBRF USA, MIPS Germania, JIPID Giappone
- SwissProt: Embl (European Bioinformatics Institute) and Swiss Institute
- PDB Databases

Alignment/Similarity Search Tools:

- BLAST
- FASTA
- ClustalW and ClustalX

Pattern/Motif/Properties/Protein Family Finding Tools:

- MotifScan
- Prosite
- Psort
- ProtScale
- SOSUI
- Protein structure classification CATH and SCOP

Molecule structure prediction & visualization Tools:

- Cn3D
- MolMol and DeepView
- 2D Structure Prediction Tools PHD, PsiPred, GORIV
- SwissModel
- PHYRE

Boutique banche dati specializzate:

- Proteine chinasi
- Proteasi dell'HIV
- Virus icosaedrici
- Immunologia

Componente della
United States National
Library of Medicine

The image shows the NCBI website interface. At the top, the NCBI logo is on the left, and the text "National Center for Biotechnology Information" is centered, with "National Library of Medicine" and "National Institutes of Health" below it. A navigation bar contains links for PubMed, All Databases, BLAST, OMIM, Books, TaxBrowser, and Structure. Below this is a search bar with a dropdown menu set to "All Databases" and a "Go" button. The main content area is divided into several sections:

- SITE MAP**: Includes links for Alphabetical List and Resource Guide.
- About NCBI**: An introduction to NCBI.
- GenBank**: Sequence submission support and software.
- Literature databases**: PubMed, OMIM, Books, and PubMed Central.
- Molecular databases**: Sequences, structures, and taxonomy.
- Genomic biology**: The human genome, whole genomes, and related.

On the right side, there are two main sections:

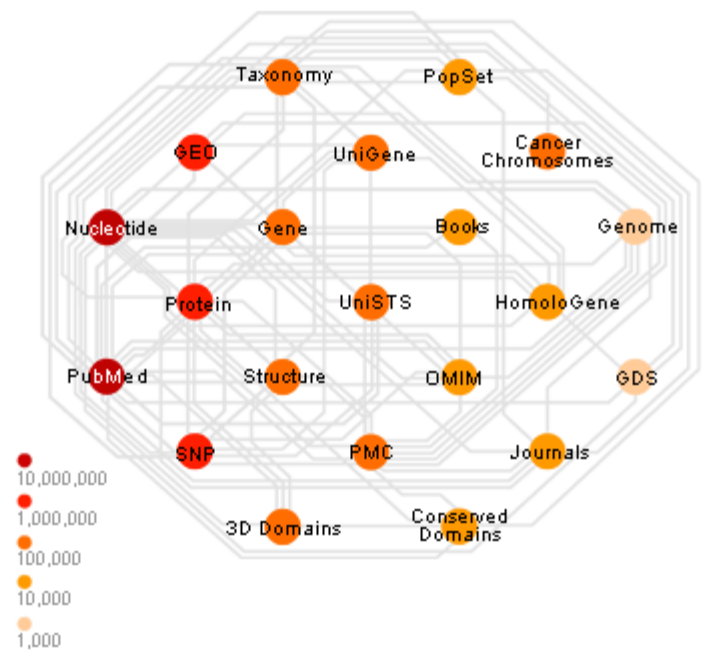
- What does NCBI do?**: A text block explaining the center's mission, established in 1988, and listing its activities: creating public databases, conducting research in computational biology, developing software tools, and disseminating biomedical information. It includes a "More..." link.
- Hot Spots**: A list of featured resources:
 - Assembly Archive
 - Clusters of orthologous groups
 - Coffee Break, Genes & Disease, NCBI Handbook
 - Electronic PCR
 - Entrez Home
 - Entrez Tools
 - Gene expression omnibus (GEO)
 - Human genome resources
 - Malaria genetics & genomics
 - Map Viewer
 - dbMHC

Two callout boxes are present:

- 100 Gigabases**: A blue callout box with a yellow background containing text: "GenBank and its collaborating databases, the European Molecular Biology Laboratory and the DNA Data Bank of Japan, have reached a milestone of 100 billion bases from over 165,000 organisms. See the [press release](#) or find more information on [GenBank](#)."
- Influenza Virus Resource**: A green callout box with a dark green background containing text: "The Influenza Virus Resource enables comparison of influenza virus strains and provides a reference for viral sequences. The resource contains data from the NIAID Influenza Genome Sequencing Project and GenBank, as well as pre-computed alignments of flu sequences."

- NCBI
 - Site Map
Guide to NCBI resources
 - Entrez Help
Help documentation for the Entrez system
 - Entrez Tutorial
 - Entrez Global Query
Search a subset of Entrez databases
 - Entrez Tools
Links to advanced Entrez tools such as Batch Entrez and E-Utilities
 - NCBI Handbook
In-depth guide to NCBI resources

Entrez is the integrated, text-based search and retrieval system used at NCBI for the major databases, including PubMed, Nucleotide and Protein Sequences, Protein Structures, Complete Genomes, Taxonomy, and others. Click on the graphic below for a more detailed view of Entrez integration.



Entrez offre accesso attraverso le seguenti divisioni di banche dati:

- Protein
- Peptide
- Nucleotide
- Structure
- Genome
- Popset (info su popolazioni)
- OMIM (eredità mendeliana nell'uomo)

Protein Database

3Ddomains

Protein

PROW

RefSeq

Conserved Domain

Structure (MMDB)

NCBI Structure

PubMed Entrez BLAST OMIM Books TaxBrowser Entrez Structure

Search Entrez Structure for Go

MMDB - Entrez's Structure Database

NCBI's Entrez includes a database of experimentally determined three-dimensional biomolecular structures. Most 3D-structure data are obtained from X-ray crystallography and NMR-spectroscopy, they provide a wealth of information on the biological function, on mechanisms linked to the function, and on the evolutionary history of and relationships between macromolecules. Our goals in adding structure data to Entrez are to make this information easily accessible to biologists, and to facilitate comparative analysis involving 3-D structure.

NCBI's structure database is called MMDB (Molecular Modeling DataBase), and it is a subset of three-dimensional structures obtained from the Protein Data Bank (PDB), excluding

Searching MMDB:

The structure database may be queried directly, using specific fields such as author names, or text occurring anywhere in the structure description. Entry points for queries the Search Bar at top of all Structure Group WWW page or the WWW-Ent interface to the 3-I structure databases

Alternatively you can use a PDB 4-character code a numerical MMDB-Id to retrieve structure summary pages directly: PDB/MMDB Code(s) Go

NCBI RefSeq

PubMed All Databases BLAST OMIM Books Taxonomy Structure

Search All Databases for Go

NCBI Reference Sequences

The Reference Sequence (RefSeq) collection aims to provide a comprehensive, integrated, non-redundant set of sequences, including genomic DNA, transcript (RNA), and protein products, for major research organisms.

RefSeq standards serve as the basis for medical, functional, and diversity studies; they provide a stable reference for gene identification and characterization, mutation analysis, expression studies, polymorphism discovery, and comparative analyses. RefSeqs are used as a reagent for the functional annotation of some genome sequencing projects, including those of human and mouse.

Site contents

- Information
 - NCBI Handbook
 - Overview | FAQ
 - Accessions | Status
 - Entrez Queries
- FTP
 - RefSeq Release
 - Catalog | Notes
 - Genomes
 - BLAST databases
- Statistics
 - Release Statistics
 - Feedback

PROW Protein Reviews On The Web

Index of information available from PROW

Current guides: expanded format including Summary Sentence and Abstract
Past guides: older guides with excellent information, some data may be dated

CD molecule	Alternate Names	Current Guides	Past Guides	Entrez Gene	Assigning Workshop
CD1a	R4; HTA1		CD1a	909	
CD1b	R1		CD1b	910	
CD1c	M241; R7		CD1c	911	
CD1d	R3		CD1d	912	
CD1e	R2		CD1e	913	
CD2	CD2R; E-rosette receptor; T11; LFA-2		CD2	914	
CD3delta	CD3d			915	
CD3epsilon	CD3e			916	
CD3gamma	CD3g			917	
CD4	L3T4; W3/25		CD4	920	

Entrez Protein

PubMed Nucleotide Protein Genome Str

for Go Clear

Limits Preview/Index History Clipboard Details

The protein entries in the Entrez search and retrieval system have been compiled from a variety of sources, including SwissProt, PIR, PRF, PDB, and translations from annotated coding regions in GenBank and RefSeq.

Human Genome

Explore [human genome resources](#) or browse the human genome sequence using the [Map Viewer](#).

NCBI Conserved Domains

HOME SEARCH SITE MAP PubMed Entrez CDD Structure Protein Taxonomy BLAST Help?

Search across Entrez databases GO CLEAR Help

A Conserved Domain Database and Search Service, v2.06

Proteins often contain several modules or domains, each with a distinct evolutionary origin and function. NCBI's Conserved Domain Database is a collection of multiple sequence alignments for ancient domains and full-length proteins. The CD-Search service may be used to identify the conserved domains present in a protein query sequence:

Submit Query Search Database CDD v2.06 - 11530 PSSMs

Enter a Protein query as Accession, GI, or Sequence in FASTA format:

Find CDDs in Entrez:

Read about the [FASTA](#) format description. Click [here](#) for advanced options.

Additional protein information

In addition to Protein sequences, other protein-related information is available via Entrez. Search the [Structure](#) database by choosing, "Structure" from the Entrez pull down menu, [Conserved Domains Database](#) (CDD) by choosing, "Domains", and [3D Domains](#) by choosing, the "3D Domains" option.

Retrieve taxonomy information

The Entrez protein database is cross-linked to the [Entrez taxonomy database](#). This allows you to find taxonomy information for the species from which a protein sequence was derived. First, look up a protein in Entrez. A "Taxonomy" link appears to the right of each entry that is linked to the Entrez taxonomy database. To view all non-redundant taxonomy links for a search result, select "Taxonomy Links" from the drop-down menu above the search results and click on the "Display" button to the left of that menu.

Search Protein for hemagglutinin Go Clear

Limits Preview/Index History Clipboard Details

Display Summary Show 20 Sort by Relevance Send

All: 15512 bacteria: 716 RefSeq: 297

Items 1 - 20 of 15512

- 1:** [CAA40728](#) Reports
 hemagglutinin [Influenza A virus]
 gi|516371|emb|CAA40728.1|[516371]
- 2:** [CAG28944](#) Reports
 hemagglutinin [Influenza A virus (A/swan/Germany/62/81(H7N7))]
 gi|55056922|emb|CAG28944.1|[55056922]
- 3:** [CAG28943](#) Reports
 hemagglutinin [Influenza A virus (A/duck/Germany/15/80(H7N7))]
 gi|55056920|emb|CAG28943.1|[55056920]
- 4:** [AAA43099](#) Reports
 hemagglutinin
 gi|323995|gb|AAA43099.1|[323995]
- 5:** [CAA91080](#) Reports
 hemagglutinin [Influenza A virus (A/Mongolia/231/85(H1N1))]
 gi|1177512|emb|CAA91080.1|[1177512]

About Entrez

Entrez Protein
 Help | FAQ

Entrez Tools

Check sequence
 revision history

LinkOut

My NCBI

Related resources
 BLAST

Reference sequence
 project

Search for Genes

Clusters of
 orthologous groups

Protein reviews on the
 web

Codice EMB
European Bioinformatics Institute

Lunghezza amminoacidica

Data di deposizione

Origine

Sequenza amminoacidica

NCBI Protein

Search Protein for [] Go Clear

Limits Preview/Index History Clipboard

Display GenPept Show 5 Send to

Range from begin to end Features: SNP CDD MGC HPRD STS tRNA Refresh

1: CAG28943 Reports hemagglutinin [In...[gi:55056920]

LOCUS CAG28943 177 aa linear VRI 31-OCT-2004

DEFINITION hemagglutinin [Influenza A virus (A/duck/Germany/15/80 (H7N7))].

ACCESSION CAG28943

VERSION CAG28943.1 GI:55056920

DBSOURCE emb1 accession AJ704797.1

KEYWORDS

SOURCE Influenza A virus (A/duck/Germany/15/80 (H7N7))

ORGANISM Influenza A virus (A/duck/Germany/15/80 (H7N7))
Viruses; ssRNA negative-strand viruses; Orthomyxoviridae;
Influenzavirus A.

REFERENCE 1

AUTHORS Starick,E. and Werner,O.

TITLE Experiences in the determination of avian influenza virus (AIV) hemagglutinin subtypes H1-H13 by reverse transcription (RT)-PCR

JOURNAL Unpublished

REFERENCE 2 (residues 1 to 177)

AUTHORS Starick,E.

TITLE Direct Submission

JOURNAL Submitted (10-MAY-2004) Starick E., Fr.-Löffler-Institutes, Fed Res Centre Virus Dis Animals, Boddenblick 5a, 17493 Greifswald-Insel riems, GERMANY

FEATURES Location/Qualifiers

source 1..177
/organism="Influenza A virus (A/duck/Germany/15/80 (H7N7))"
/virion
/isolate="A/duck/Germany/15/80"
/db_xref="taxon:278126"
/segment="4"
/country="Germany"

Protein 1..177
/product="hemagglutinin"

mat peptide <1..>177
/product="hemagglutinin"

CDS 1..177
/gene="HA"
/coded_by="AJ704797.1:<1..>532"
/experiment="experimental evidence, no additional details"
/product="hemagglutinin"

mat peptide <1..>177
/product="hemagglutinin"

CDS 1..177
/gene="HA"
/coded_by="AJ704797.1:<1..>532"
/experiment="experimental evidence, no additional details recorded"

ORIGIN

1 eqtklygsgs klitvgsnsy qgsfvpspga rpvngqsgs idfhwlmnp ndvtvtsing
61 afiapdrasf lrgksmgigs dvqvdanceg dcyhsaggtil snlpfqins ravgkcprrv
121 kqeslllatg mknvpeipkg rgllfgaiagf iengweglvd gwygfrhqlna qgegtnh



PIR Protein Information Resource

About PIR

Databases

Search and Retrieval

Download

Support

AN INTEGRATED PUBLIC RESOURCE OF PROTEIN INFORMATICS TO SUPPORT GENOMIC AND PROTEOMIC RESEARCH AND SCIENTIFIC DISCOVERY

Since 1984, PIR has produced the **Protein Sequence Database (PSD)** of functionally annotated protein sequences, which grew out of the Atlas of Protein Sequence and Structure (1965-1978) edited by Margaret Dayhoff. Now a part of the [UniProt](#) effort, sequences and annotations in PIR-PSD have been integrated into UniProt Knowledgebase. Release 80.00 (31-Dec-2004) is the final release for PSD.

iProClass, a central point for exploration of protein information, provides summary descriptions of protein family, function and structure for PIR-PSD, Swiss-Prot, and TrEMBL sequences, with links to over 50 biological databases. [Release 2.82](#), 12-Dec-2005, contains 2568,372 entries.

PIR-NREF, a comprehensive database for sequence searching and protein identification, contains non-redundant protein sequences from PIR-PSD, Swiss-Prot, TrEMBL, RefSeq, GenPept, and PDB. [Release 1.82](#), 12-Dec-2005, contains 3,254,342 entries.



← Click here to try the new PIR beta site



Text Search Protein Databases:

Find an Exact Peptide Match:

Type in a string of single letter amino acid code (at least 3 letters)



PIR has recently joined forces with [EBI](#) (European Bioinformatics Institute) and [SIB](#) (Swiss Institute of Bioinformatics) to establish the [UniProt](#) (United Protein Databases), the central resource of protein sequence and function.

Combinazione efficace
Di una banca dati accurata
Software per la ricerca di informazioni
Banco di lavoro per lo studio di sequenze



PIR Search Results

[Site Map](#)[Site Search](#)Text Search Protein Databases: [About PIR](#)[Databases](#)[Search & Retrieval](#)[Download](#)[Support](#)

Thu Jan 26 10:45:30 EST 2006

[← examples](#)

Search for

6763 protein sequences in total. [page 1 \(50/page\)](#) [Next](#)[HELP](#)For sequence analyses, pick a method (**radiobutton**) below, select a sequence(s) (**checkbox**) in *Protein ID* column, and GO.
 BLAST
 FASTA
 HMM Search
 Pattern Match
 Multiple Alignment
 Domain Display

<input type="checkbox"/> Protein ID <small>check all</small>	Matched	Protein Name	Length	Organism Name /Taxon Group	PIRSF ID /Family ID	Pfam ID	PC Motif ID	PDB ID
<input type="checkbox"/> NREF: NF02311106 iProClass: Q5G7N6 9CAUD UniProt: Q5G7N6 9CAUD	Protein Name =>hemagglutinin Protein Name =>hemagglutinin	Putative hemagglutinin protein	801	Listonella pelagia phage phiHSIC Virus				
<input type="checkbox"/> NREF: NF00058625 iProClass: O55889 9PARA UniProt: O55889 9PARA	Paper Title =>hemagglutinin-neuraminidase	RNA dependent RNA polymerase (Fragment)	916	Recombinant PIV3/PIV1 virus Virus	SF000830	PF00946		
<input type="checkbox"/> NREF: NF00058624 iProClass: Q77222 9PARA UniProt: Q77222 9PARA	Protein Name =>Hemagglutinin Keyword =>hemagglutinin PIRSF Name =>hemagglutinin-neuraminidase Pfam Name =>Hemagglutinin-neuraminidase Paper Title =>hemagglutinin-neuraminidase	Hemagglutinin	575	Recombinant PIV3/PIV1 virus Virus	SF001072	PF00423		
<input type="checkbox"/> NREF: NF00058623 iProClass: O55888 9PARA UniProt: O55888 9PARA	Paper Title =>hemagglutinin-neuraminidase	Fusion glycoprotein	555	Recombinant PIV3/PIV1 virus Virus		PF00523		
<input type="checkbox"/> NREF: NF00244898 iProClass: Q83350 9PARA UniProt: Q83350 9PARA	Keyword =>hemagglutinin PIRSF Name =>hemagglutinin Pfam Name =>Hemagglutinin-neuraminidase	Haemagglutinin protein (H)	607	Morbillivirus Virus	SF003926	PF00423		
<input type="checkbox"/> NREF: NF00244895 iProClass: Q83352 9PARA UniProt: Q83352 9PARA	Keyword =>hemagglutinin PIRSF Name =>hemagglutinin Pfam Name =>Hemagglutinin-neuraminidase	Haemagglutinin protein (H)	607	Morbillivirus Virus	SF003926	PF00423		



PIR Search Results

Site Map Site Search

Text Search Protein Databases:

About PIR

Databases

Search & Retrieval

Download

Support

Thu Jan 26 10:46:56 EST 2006

[examples](#)

Search for AND AND AND

12 protein sequences in total.

[HELP](#)

For sequence analyses, pick a method (radiobutton) below, select a sequence(s) (checkbox) in Protein ID column, and GO.

BLAST FASTA HMM Search Pattern Match Multiple Alignment Domain Display

<input type="checkbox"/> Protein ID <small>check all</small>	Matched	Protein Name	Length	Organism Name /Taxon Group	PIRSF ID /Family ID	Pfam ID	PC Motif ID	PDB ID
<input type="checkbox"/> NREF:NF00272830 iProClass:Q84097_9INFB UniProt:Q84097_9INFB	Paper Title =>hemagglutinin Paper Title =>influenza	Influenza B/Hong Kong/8/73 hemagglutinin (HA) (seg 4) RNA, complete cds	582	Influenza B virus Virus		PF00509		
<input type="checkbox"/> NREF:NF00272800 iProClass:Q84102_9INFB PIR-PSD:S11738 UniProt:Q84102_9INFB	Paper Title =>hemagglutinin Paper Title =>influenza	HA haemagglutinin precursor, genomic RNA, strain B/NIB/25/88 (Fragment)	378	Influenza B virus Virus	SF003927 FAM0004468	PF00509		
<input type="checkbox"/> NREF:NF00263104 iProClass:HMIVT3 PIR-PSD:HMIVT3	Paper Title =>hemagglutinin Paper Title =>influenza	hemagglutinin precursor	565	Influenza A virus Virus	SF003927 FAM0004463			
<input type="checkbox"/> NREF:NF00263701 iProClass:HMIVT7 PIR-PSD:HMIVT7	Paper Title =>hemagglutinin Paper Title =>influenza	hemagglutinin precursor	560	Influenza A virus Virus	SF003927 FAM0004462			
<input type="checkbox"/> NREF:NF00263284 iProClass:HMIV84 PIR-PSD:HMIV84	Paper Title =>hemagglutinin Paper Title =>influenza	hemagglutinin precursor	561	Influenza A virus Virus	SF003927 FAM0004462			
<input type="checkbox"/> NREF:NF00263110 iProClass:P13102 PIR-PSD:HMIVT2	Paper Title =>hemagglutinin Paper Title =>influenza	Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain]	566	Influenza A virus Virus	SF003927 FAM0004463	PF00509		



PIR BLAST/SSearch Results

Site Map Site Search

Text Search Protein Databases:



About PIR

Databases

Search & Retrieval

Download

Support

Sort by E-value

Re-Load

[View Superfamily Summary](#)

Search

Protein Name

Query sequence:

>Your input sequence:

```

1 QDLPGNDNSTATLCLGHHAVPNGTLVKTITDDQIEVTNATELVQSSSTGKICNNPRLIDGIDCTLIDALLGDPHCDVFPQ
81 NETWDLFVRSKAFSNCYPYDVPDYASLRSLVASSGTFLEITEGFTWTGVTQNGRSNACKRPGSGGFFSRLNULTKSGST
161 YPVLNVTMPNNDNFDKLYIWIHHPSTNQEQTSLYVQASGRVTVSTRSQQTIIIPNIGSRPVRGLSSRSIYWTIVKPG
241 DVLVINSNGNLIAPRGYFKMRTGKSSIMRSDAPIDTCISECITPNGSIIPNDKPFQNVNKIITYGACPKYVKQNTLKLATGM
321 RNVPEKQTGLFGAIAAGFIENGWEHIDGWYGRHQNSEGTGQAADLKSTQAAIDQINGKLNRIEKTNEKFHQIEKEFSE
401 VEGRIQDLEKYVEDTKIDLWSYNAELLVALENQHTIDLTDSEMNKLFKTRRQLRENAEEMGNCGCFKIYHKCDNACIESI
481 IRNGTYDHDVYRDEALNNRFQIKG

```

200 match(es) shown in the following table:

[HELP](#)

For sequence analyses, pick a method (radio button) below, select a sequence(s) (checkbox) in Protein ID column, and GO.

BLAST FASTA HMM Search Pattern Match Multiple Alignment Domain Display

<input type="checkbox"/> Protein ID	Protein Name	Organism	Taxon Group	PIRSE ID	e-value	Length	Ov.lap	%idn	Query Sequence
<input type="checkbox"/> NREF: NF00263909 iProClass: Q3AVI6_9SYNE PIR-PSD: HMIVHA UniProt: HEMA_IAAIC	Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain]	Influenza A virus	Virus	SF003927	0.0	566	504	99	
<input type="checkbox"/> NREF: NF00262744 iProClass: Q5K598_SALSA UniProt: Q67132_9INFA	Hemagglutinin	Influenza A virus	Virus		0.0	566	504	99	>NF00263909 Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain] [Influenza A virus] Length = 566 Score = 1034 bits (2673), Expect = 0.0 Identities = 502/504 (99%), Positives = 502/504 (99%), Gaps = 1/504 (0%)
<input type="checkbox"/> NREF: NF01662375 iProClass: Q47Z20_COLP3 UniProt: Q91MA7_IAHO1	Hemagglutinin	Influenza A virus	Virus		0.0	566	504	99	Query: 1 QDLPGNDNSTATLCLGHHAVPNGTLVKTITDDQIEVTNATELVQSSSTGKICNNPRLID 60 Sbjct: 17 QDLPGNDNSTATLCLGHHAVPNGTLVKTITDDQIEVTNATELVQSSSTGKICNNPRLID 76
<input type="checkbox"/> NREF: NF01661952 iProClass: Q487E6_COLP3 UniProt: Q910M5_IAHO1	Hemagglutinin	Influenza A virus	Virus		0.0	566	504	98	Query: 61 GIDCTLIDALLGDPHCDVFPQNETWDLFVRSKAFSNCYPYDVPDYASLRSLVASSGTFLE 120 Sbjct: 77 GIDCTLIDALLGDPHCDVFPQNETWDLFVRSKAFSNCYPYDVPDYASLRSLVASSGTFLE 136
<input type="checkbox"/> NREF: NF00263705 iProClass: Q4ZE13_9CAUD	Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain;	Influenza A virus	Virus	SF003927	0.0	566	504	98	Query: 121 ITEGFTWTGVTQNGRSNACKRPGSGGFFSRLNULTKSGSTYPVLNVTMPNNDNFDKLYIW 180 Sbjct: 137 ITEGFTWTGVTQNGRSNACKRPGSGGFFSRLNULTKSGSTYPVLNVTMPNNDNFDKLYIW 196
									Query: 181 GIHHPSTNQEQTSLYVQASGRVTVSTRSQQTIIIPNIGSRPVRGLSSRSIYWTIVKPG 240 Sbjct: 197 GIHHPSTNQEQTSLYVQASGRVTVSTRSQQTIIIPNIGSRPVRGLSSRSIYWTIVKPG 256
									Query: 241 DVLVINSNGNLIAPRGYFKMRTGKSSIMRSDAPIDTCISECITPNGSIIPNDKPFQNVNKI 300 Sbjct: 257 DVLVINSNGNLIAPRGYFKMRTGKSSIMRSDAPIDTCISECITPNGSIIPNDKPFQNVNKI 316
									Query: 301 TYGACPKYVKQNTLKLATGHRNVPEKQTGLFGAIAAGFIENGWEHIDGWYGRHQNSE 359 Sbjct: 317 TYGACPKYVKQNTLKLATGHRNVPEKQTRGLFGAIAAGFIENGWEHIDGWYGRHQNSE 376
									Query: 360 TQQAADLKSTQAAIDQINGKLNRIEKTNEKFHQIEKEFSEVEGRIQDLEKYVEDTKIDL 419 Sbjct: 377 TQQAADLKSTQAAIDQINGKLNRIEKTNEKFHQIEKEFSEVEGRIQDLEKYVEDTKIDL 436
									Query: 420 WSYNAELLVALENQHTIDLTDSEMNKLFKTRRQLRENAEEMGNCGCFKIYHKCDNACIES 479 Sbjct: 437 WSYNAELLVALENQHTIDLTDSEMNKLFKTRRQLRENAEEMGNCGCFKIYHKCDNACIES 496
									Query: 480 IRNGTYDHDVYRDEALNNRFQIKG 503 Sbjct: 497 IRNGTYDHDVYRDEALNNRFQIKG 520

Search for 

ExPASy Proteomics Server

The ExPASy (**Ex**pert **P**rotein **A**nalysis **S**ystem) proteomics server of the [Swiss Institute of Bioinformatics](#) (SIB) is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE ([Disclaimer](#) / [References](#)).

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

Databases

- [Swiss-Prot and TrEMBL](#) - Protein knowledgebase
- [PROSITE](#) - Protein families and domains
- [SWISS-2DPAGE](#) - Two-dimensional polyacrylamide gel electrophoresis
- [ENZYME](#) - Enzyme nomenclature
- [SWISS-MODEL Repository](#) - Automatically generated protein models
- [GermOnLine](#) - Knowledgebase on germ cell differentiation
- [Ashbya Genome Database](#)
- [Links to many other molecular biology databases](#)

Tools and software packages

- [Proteomics and sequence analysis tools](#)
 - ◊ [Proteomics](#) [[Aldente](#) (PMF), [Phenyx](#) (MS/MS), [FindMod](#), [PeptideMass](#), ...]
 - ◊ [DNA -> Protein](#) [[Translate](#)]
 - ◊ [Similarity searches](#) [[BLAST](#)]
 - ◊ [Pattern and profile searches](#) [[ScanProsite](#)]
 - ◊ [Post-translational modification and topology prediction](#)
 - ◊ [Primary structure analysis](#) [[ProtParam](#), [pI/MW](#), [ProtScale](#)]
 - ◊ [Secondary and tertiary structure prediction](#) [[SWISS-MODEL](#), [Swiss-PdbViewer](#)]
 - ◊ [Alignment](#) [[T-COFFEE](#), [SIM](#)]
 - ◊ [Phylogenetic analysis](#)
 - ◊ [Biological text analysis](#)
- [ImageMaster / Melanie](#) - Software for 2-D PAGE analysis
- [MSight](#) - Mass Spectrometry Imager
- [Roche Applied Science's Biochemical Pathways](#)

Sistema di ricerca e analisi delle informazioni dello Swiss Institute of Bioinformatics in collaborazione con EBI (European Bioinformatics Institute) gestendo anche le banche dati SwissProt e TrEMBL traduzione delle sequenze nucleotidiche di EMBL

Search in UniProt Knowledgebase (Swiss-Prot and TrEMBL) for: hemagglutinin

UniProtKB/Swiss-Prot Release 48.9 of 24-Jan-2006

UniProtKB/TrEMBL Release 31.9 of 24-Jan-2006

-
- Number of sequences found in [UniProt Knowledgebase \(Swiss-Prot^{\(326\)} and TrEMBL^{\(6938\)}\)](#): **7259**
 - Note that the selected sequences can be saved to a file to be later retrieved; to do so, go to the [bottom](#) of this page.
 - For more directed searches, you can use the Sequence Retrieval System [SRS](#).
-

Search in UniProtKB/Swiss-Prot: There are matches to 326 out of 206586 entries

[AFAF_ECOLI \(Q47037\)](#)

Dr hemagglutinin AFA-III operon regulatory protein afaF. {GENE: Name=afaF} - Escherichia coli

[FHAB_BORPE \(P12255\)](#)

Filamentous hemagglutinin. {GENE: Name=fhaB; OrderedLocusNames=BP1879} - Bordetella pertussis

[FMA1_ECOLI \(P08180\)](#)

Afimbrial adhesin AFA-I precursor (Dr hemagglutinin AFA-I). {GENE: Name=afaE1; Synonyms=afaE, afaE-1} - Escherichia coli

[FMDR_ECOLI \(P24093\)](#)

Dr hemagglutinin structural subunit precursor. {GENE: Name=draA} - Escherichia coli

[HA17_CLOBO \(P46083\)](#)

Hemagglutinin component HA-17 (HA 17 kDa subunit). {GENE: Name=HA-17; Synonyms=antP-17} - Clostridium botulinum

[HA33_CLOBO \(P46084\)](#)

Main hemagglutinin component (HA 33 kDa subunit) (HA1). {GENE: Name=HA-33; Synonyms=antP-33, ha1} - Clostridium botulinum

Search in UniProt Knowledgebase (Swiss-Prot and TrEMBL) for: hemagglutinin influenza

UniProtKB/Swiss-Prot Release 48.9 of 24-Jan-2006

UniProtKB/TrEMBL Release 31.9 of 24-Jan-2006

- Number of sequences found in [UniProt Knowledgebase \(Swiss-Prot\)](#)⁽¹⁷⁷⁾ and [TrEMBL](#)⁽⁵⁴⁰⁶⁾: **5583**
- Note that the selected sequences can be saved to a file to be later retrieved; to do so, go to the [bottom](#) of this page.
- For more directed searches, you can use the Sequence Retrieval System [SRS](#).

~~Search in UniProtKB/Swiss-Prot:~~ There are matches to 177 out of 206586 entries

[HEMA_JAIC \(P03437\)](#)

Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain]. {GENE: Name=HA} - Influenza A virus (strain A/Aichi/2/68 H3N2)

[HEMA_JABAN \(P03441\)](#)

Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain] (Fragment). {GENE: Name=HA} - Influenza A virus (strain A/Bangkok/1/79 H3N2)

[HEMA_JABUD \(P19694\)](#)

Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain]. {GENE: Name=HA} - Influenza A virus (strain A/Budgerigar/Hokkaido/1/77 H4N6)

[HEMA_JACAO \(P26142\)](#)

Hemagglutinin (Fragment). {GENE: Name=HA} - Influenza A virus (strain A/Camel/Mongolia/82 H1N1)

[HEMA_JACKA \(P19695\)](#)

Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 chain]. {GENE: Name=HA} - Influenza A virus (strain A/Chicken/Alabama/1/75)

[ExpASY Home page](#)[Site Map](#)[Search ExpASY](#)[Contact us](#)[Swiss-Prot](#)Search for

UniProtKB/Swiss-Prot entry **P03437**

[Printer-friendly view](#)[Submit update](#)[Quick BlastP search](#)

[\[Entry info\]](#) [\[Name and origin\]](#) [\[References\]](#) [\[Comments\]](#) [\[Cross-references\]](#) [\[Keywords\]](#) [\[Features\]](#)
[\[Sequence\]](#) [\[Tools\]](#)

Note: most headings are clickable, even if they don't appear as links. They link to the user manual or other documents.

Entry information

Entry name	HEMA_IAAIC
Primary accession number	P03437
Secondary accession numbers	None
Entered in Swiss-Prot in	Release 01, July 1986
Sequence was last modified in	Release 01, July 1986
Annotations were last modified in	Release 49, February 2006

Name and origin of the protein

Cross-references

Sequence databases

EMBL	J02090; AAA43178.1; -, Genomic_RNA. [EMBL / GenBank / DDBJ] [CoDingSequence] V01085; CAA24269.1; -, Genomic_RNA. [EMBL / GenBank / DDBJ] [CoDingSequence]
PIR	A93231; HMIVHA .

3D structure databases

PDB	1EO8; X-ray, A=17-344, B=346-520. [ExPASy / RCSB / EBI]
	1HA0; X-ray, A=25-518. [ExPASy / RCSB / EBI]
	1HGD; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGE; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGF; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGG; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGH; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGI; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HGJ; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1HTM; X-ray, A/C/E=17-43, B/D/F=383-520. [ExPASy / RCSB / EBI]
	1J8H; X-ray, C=322-334. [ExPASy / RCSB / EBI]
	1KEN; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
	1QFU; X-ray, A=17-344, B=346-520. [ExPASy / RCSB / EBI]
	1QU1; X-ray, A/B/C/D/E/F=376-530. [ExPASy / RCSB / EBI]
	2HMG; X-ray, A/C/E=17-344, B/D/F=346-520. [ExPASy / RCSB / EBI]
2VIR; X-ray, C=44-325. [ExPASy / RCSB / EBI]	

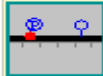
Family and domain databases

InterPro	IPR001364 ; Hemagglutn. IPR000149 ; Hemagglutn_1. Graphical view of domain structure.
Pfam	PF00509 ; Hemagglutinin; 1. Pfam graphical view of domain structure.
PRINTS	PR00330 ; HEMAGGLUTN1. PR00329 ; HEMAGGLUTN12.
ProDom	PD000225 ; Hemagglutn; 1. [Domain structure / List of seq. sharing at least 1 domain]
BLOCKS	P03437 .

Other

[UniProt](#) [P03437](#)

Features



Feature table viewer



Feature aligner

Key	From	To	Length	Description	FTId
SIGNAL	1	16	16		
CHAIN	17	344	328	Hemagglutinin HA1 chain.	PRO_0000038885
CHAIN	346	566	221	Hemagglutinin HA2 chain.	PRO_0000038886
TOPO_DOM	17	530	514	Extracellular (<i>Potential</i>).	
TRANSMEM	531	551	21	<i>Potential</i> .	
TOPO_DOM	552	566	15	Cytoplasmic (<i>Potential</i>).	
LIPID	555	555		S-palmitoyl cysteine (by host) (<i>By similarity</i>).	
LIPID	562	562		S-palmitoyl cysteine (by host) (<i>By similarity</i>).	
LIPID	565	565		S-palmitoyl cysteine (by host) (<i>By similarity</i>).	
CARBOHYD	499	499		N-linked (GlcNAc...).	
DISULFID	30	482		Interchain (between HA1 and HA2 chains).	
DISULFID	68	293			
DISULFID	80	92			
DISULFID	113	155			
DISULFID	297	321			
DISULFID	489	493			
STRAND	27	35	9		
STRAND	37	38	2		
STRAND	40	42	3		
STRAND	45	46	2		
STRAND	48	48	1		
STRAND	50	53	4		
STRAND	55	57	3		
STRAND	59	60	2		
STRAND	65	72	8		
STRAND	74	76	3		
TURN	78	79	2		
HELIX	82	87	6		
TURN	88	88	1		
HELIX	90	95	6		
TURN	96	97	2		
STRAND	99	99	1		
STRAND	101	105	5		

Sequence information

Length: **566 AA** [This is the length of the unprocessed precursor]

Molecular weight: **63416 Da** [This is the MW of the unprocessed precursor]

CRC64: **E395659C23CAFECA** [This is a checksum on the sequence]

```
      10      20      30      40      50      60
MKTIIALSYI FCLALGQDLP GNDNSTATLC LGHHAVPNGT LVKTIITDDQI EVTNATELVQ

      70      80      90     100     110     120
SSSTGKICNN PHRILDGIDC TLIDALLGDP HCDVFNQETW DLFVERSKAF SNCYPYDVPD

     130     140     150     160     170     180
YASLRSLVAS SGTLEFITEG FTWTGVTQNG GSNACKRGPQ SGFFSRLNWL TKSGSTYPVL

     190     200     210     220     230     240
NVTMPNNDNF DKLYIWIHH PSTNQEQTSL YVQASGRVTV STRRSQQTII PNIGSRPWVR

     250     260     270     280     290     300
GLSSRISIW TIVKPGDVLV INSNGLIAP RGYFKMRTGK SSIMRSDAPI DTCISECITP

     310     320     330     340     350     360
NGSIPNDKPF QNVNKITYGA CPKYVKQNTL KLATGMRNVP EKQTRGLFGA IAGFIENGWE

     370     380     390     400     410     420
GMIDGWYGR HONSEGTOQA ADLKSTQAAI DQINGKLNRV IEKTNEKFHQ IEKEFSEVEG
```

Formato FASTA
formato di interscambio
tra software

```
>sp|P03437|HEMA_IAAIC Hemagglutinin precursor [Contains: Hemagglutinin HA1 chain; Hemagglutinin HA2 c
MKTIIALSYIFCLALGQDLPGNDNSTATLCLGHHAVPNGTLVKTIITDDQIEVTNATELVQ
SSSTGKICNNPHRILDGIDCTLIDALLGDPHCDVFNQETWDLFVERSKAFSNCYPYDVPD
YASLRSLVASSGTLEFITEGFTWTGVTQNGGSSNACKRGPQSGFFSRLNWLTKSGSTYPVL
NVTMPNNDNFDKLYIWIHHPSSTNQEQTSLYVQASGRVTVSTRRSQQTIIIPNIGSRPWVR
GLSSRISIWTVIVKPGDVLVINSNGLIAPRGYFKMRTGKSSIMRSDAPIIDTCISECITP
NGSIPNDKPFQNVNKITYGACPKYVKQNTLKLATGMRNVPKQTRGLFGAIIAGFIENGWE
GMIDGWYGRHONSEGTOQAADLKSTQAAIDQINGKLNRVIEKTNEKFHQIEKEFSEVEG
RIQDLEKYVEDTKIDLWSYNAELLVALENQHTIDLTDSEMKNLFEKTRRQLRENAEEMGN
GCFKIYHRCDNACIESIRNGTYDHDVYRDEALNNRFQIKGVELKSGYKDWILWISFAISC
FLLCVLLGFIMWACQRGNIRCNICI
```

- [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.

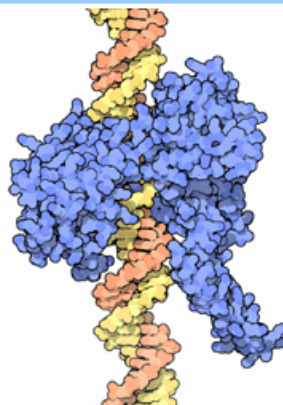
This site offers tools for browsing, searching, and reporting that utilize the data resulting from ongoing efforts to create a more consistent and comprehensive archive.

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: Topoisomerases



Each of your cells contains about 2 meters of DNA, all folded into the tiny space inside the nucleus, which is a million times smaller. As you might imagine, these long, thin strands can get tangled very easily in the busy environment of the nucleus. To make things even more complicated, DNA is a double helix, which must be unwound to access the genetic information.

- [More ...](#)
- [Previous Features](#)

NEWS

- [Complete News](#)
- [Newsletter](#)
- [Discussion Forum](#)

24-Jan-2006

Montville Township High School Places First at the NJ Science Olympiad Protein Modeling Trial Event

Several high school teams competed in the protein modeling event at the New Jersey Northern Regional **Science Olympiad** that was held January 12, 2006 at Montclair State University.



■ [Full Story ...](#)

17-Jan-2006

Time-stamped Copies of PDB Archive Available via FTP

10-Jan-2006

Structural Genomics Tools and Portal

Results (1-10 of 87)

[Refine this Search](#)

[1 Structures Awaiting Release](#)

[Select All](#)

[Deselect All](#)

[Download Selected](#)

[▶ Tabulate](#)

[▶ Narrow Query](#)

[▶ Sort Results](#)

[▶ Results per Page](#)

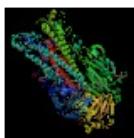
[■ Show Query Details](#)

[Results Help](#)

1RUZ



1918 H1 Hemagglutinin



Characteristics

Release Date: 30-Mar-2004 **Exp. Method:** X Ray Diffraction

Classification

Resolution: 2.90 Å

Virus/viral Protein

Compound

Mol. Id: 1 **Molecule:** Hemagglutinin **Mol. Id:** 2

Molecule: Hemagglutinin

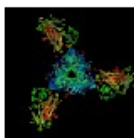
Authors

Gamblin, S.J., Haire, L.F., Russell, R.J., Stevens, D.J., Xiao, B., Ha, Y., Vasisht, N., Steinhauer, D.A., Daniels, R.S., Elliot, A., Wiley, D.C., Skehel, J.J.

1EOB



INFLUENZA VIRUS HEMAGGLUTININ COMPLEXED WITH A NEUTRALIZING ANTIBODY



Characteristics

Release Date: 12-Apr-2000 **Exp. Method:** X Ray Diffraction

Classification

Resolution: 2.80 Å

Virus/viral Protein

Compound

Mol. Id: 1 **Molecule:** Hemagglutinin (ha1 Chain) **Fragment:** Bromelain

Released Fragment **Mol. Id:** 2 **Molecule:** Hemagglutinin (ha2 Chain)

Fragment: Bromelain Released Fragment **Mol. Id:** 3 **Molecule:** Antibody

(light Chain) **Fragment:** Fab Fragment of Antibody Bh151 **Mol. Id:** 4

Molecule: Antibody (heavy Chain) **Fragment:** Fab Fragment of Antibody Bh151

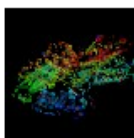
Authors

Fleury, D., Daniels, R.S., Skehel, J.J., Knossow, M., Bizebard, T.

1FLC



X-RAY STRUCTURE OF THE HAEMAGGLUTININ-ESTERASE-FUSION GLYCOPROTEIN OF INFLUENZA C VIRUS



Characteristics

Release Date: 01-Mar-2000 **Exp. Method:** X Ray Diffraction

Classification

Resolution: 3.20 Å

Hydrolase

Compound

Mol. Id: 1 **Molecule:** Haemagglutinin Esterase Fusion Glycoprotein

Fragment: Hef1 **Mol. Id:** 2 **Molecule:** Haemagglutinin Esterase Fusion

Glycoprotein **Fragment:** Hef2

Authors

Rosenthal, P.B., Zhang, X.

1FYT



CRYSTAL STRUCTURE OF A COMPLEX OF A HUMAN ALPHA/BETA-T CELL RECEPTOR, INFLUENZA HA ANTIGEN PEPTIDE, AND MHC CLASS II MOLECULE HLA-DR1

- 1RUZ
- Download Files**
- FASTA Sequence
- Display Files
- Display Molecule
- Structural Reports
- Structure Analysis
- Help

1RUZ

Images and Visualization

Biological Molecule / Asymmetric Unit



Display Options

- KING
- Jmol
- WebMol
- Protein Workshop
- QuickPDB
- All Images

Title 1918 H1 Hemagglutinin

Authors Skehel, J.J., Gamblin, S.J., Haire, L.F., Russell, R.J., Stevens, D.J., Xiao, B., Ha, Y., Vasisht, N., Steinhauer, D.A., Daniels, R.S.

Primary Citation Gamblin, S.J., Haire, L.F., Russell, R.J., Stevens, D.J., Xiao, B., Ha, Y., Vasisht, N., Steinhauer, D.A., Daniels, R.S., Elliot, A., Wiley, D.C., Skehel, J.J. The structure and receptor binding properties of the 1918 influenza hemagglutinin. *Science* v303 pp.1838-1842, 2004
[\[Abstract\]](#)

History Deposition 2003-12-12 Release 2004-03-30

Experimental Method Type X-RAY DIFFRACTION Data N/A

Parameters	Resolution[Å]	R-Value	R-Free	Space Group
	2.90	0.248 (obs.)	0.289	P 4 ₁ 2 ₁ 2

Unit Cell	Length [Å]	a	b	c	Angles [°]	alpha	beta	gamma
		171.46	171.46	153.45	90.00	90.00	90.00	90.00

Molecular Description Polymer: 1 Molecule: hemagglutinin Chains: H, J, L; Other Details: Hemagglutinin HA1 chain Polymer: 2 Molecule: hemagglutinin Chains: I, K, M; Other Details: Hemagglutinin HA2 chain

Functional Class Virus/viral Protein

Source Polymer: 1 Scientific Name: Influenza a virus Common Name: Virus Polymer: 2 Scientific Name: Influenza a virus Common Name: Virus

Related PDB Entries	Id	Details
	1RVZ	1934 H1 Hemagglutinin in complex with LSTC
	1RVX	1934 H1 Hemagglutinin in complex with LSTA
	1RVT	1930 H1 Hemagglutinin in complex with LSTC
	1RV0	1930 Swine H1 Hemagglutinin complexed with LSTA
	1RUY	1930 H1 Hemagglutinin
	1RU7	1934 H1 Hemagglutinin

Chemical Component	Identifier	Name	Formula	Drug Similarity	Ligand Structure	Ligand Interaction

La qualità e la quantità dei dati hanno spinto gli scienziati a puntare verso traguardi proporzionatamente ambiziosi:

- Comprendere gli aspetti integrativi della biologia degli organismi osservati come sistemi complessi coerenti.
- Correlare sequenza, struttura, interazioni e funzionalità di singole proteine, acidi nucleici e complessi, tra di essi
- Usare dati riguardanti organismi contemporanei come base per viaggiare avanti e indietro nel tempo; per dedurre eventi della storia evolutiva.
- Fornire un supporto ad applicazioni nel campo della medicina, agricoltura ecc.

Scenario futuro:

Nuovo Virus altamente pericoloso > pandemie sul pianeta

Sequenziamento del suo genoma > confronto con la sua sequenza con quella dei virus presenti in Banca Dati > Individuazione di proteine patogene per l'uomo responsabili dell'infettività e virulenza > Determinazione della struttura tridimensionale di tali proteine tramite Homology modeling o Modeling predittivo > confronto tramite Banca Dati con struttura di altre proteine e di conseguenza loro funzione > Drug Design per ricerca di nuovi farmaci o individuazioni di Anticorpi che possano neutralizzare il Virus.

Tempo necessario pochi giorni > pandemia evitata

Confronto di sequenze

Il confronto tra sequenze è in biologia computazionale la base per

- o misurare la “similarità” tra le sequenze
 - ❖ allineamento
- o misurare la “diversità” tra le sequenze
 - ❖ distanza di edit
- o trovare parti comuni alle sequenze
 - ❖ pattern discovery
 - ❖ allineamento locale

INTRODUZIONE

**DATABASE DI
SEQUENZE**



RICERCA



TESTUALE

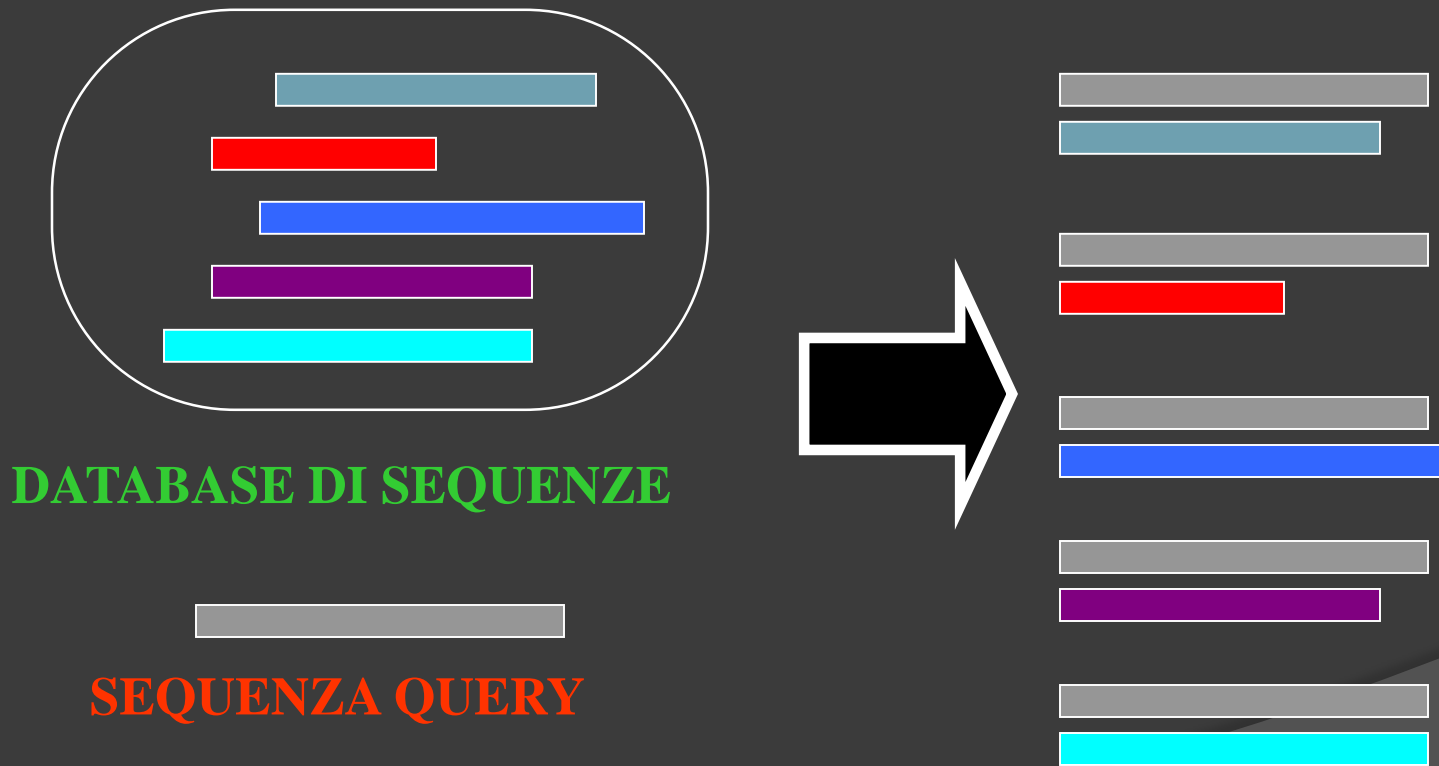
SIMILARITA'

Ricerca dei record i cui campi soddisfano determinati criteri (hanno certi valori)

Ricerca dei record che hanno le sequenze più “simili” ad una sequenza fornita come query

RICERCA PER SIMILARITA'

- La ricerca per similarità di una sequenza contro un database di sequenze richiede che sia possibile valutare la similarità della sequenza query contro ciascuna delle sequenze del database. Quindi il problema da risolvere è quello della ricerca delle similarità tra due sequenze

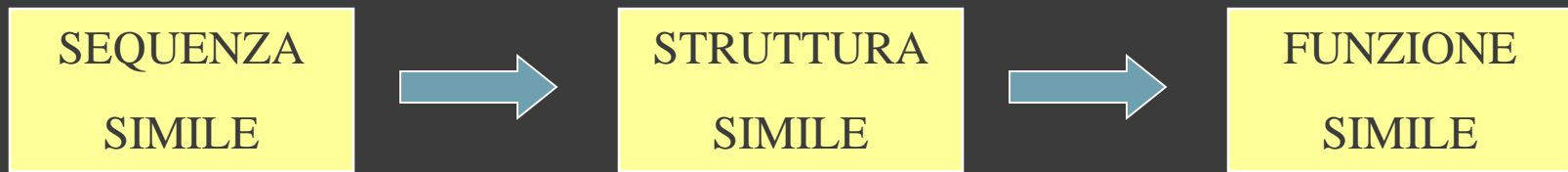


PERCHE' CERCARE SEQUENZE SIMILI?

- Quando si ottiene (in qualche modo) una sequenza di DNA o Aminoacidi si è interessati a capire cos'è quella sequenza (è già nota?) e a scoprire la sua funzione.
- Potrebbe anche capitare che la sequenza stessa sia presente nei database e già annotata (descritta la sua funzione)... Nel caso invece non si trovasse nei database esattamente la stessa sequenza, un modo semplice di ipotizzare (è comunque una predizione, che dovrà poi essere confermata sperimentalmente) la funzione della mia sequenza query è quello di cercare sequenze simili che invece siano già state annotate.
- In base al grado di similarità trovato diventa possibile fare delle ipotesi più o meno probabili sulla funzione della sequenza query semplicemente “trasferendo” ad essa la funzione delle sequenze target simili ad essa identificate .

QUANDO SUPPORRE LA FUNZIONE

- Se le sequenze di due proteine (DNA) sono molto simili allora lo saranno anche le strutture e le funzioni



- Non vale il viceversa! (Funzioni e strutture simili non implicano sequenze simili)
- Ci possono essere proteine con la stessa funzione, ma con struttura e soprattutto sequenza diversa. Es. mutazioni silenziose, che interessano la terza base di un codone. L'aminoacido rimane lo stesso ma è cambiato il DNA!

SIMILARITA' E OMOLOGIA

- Spesso si fa confusione tra similarità ed omologia!
- La similarità è un aspetto quantitativo che indica (fissato un criterio comparativo, % identità, % mutazioni conservative...) un livello di somiglianza tra le sequenze.
- L'omologia è un aspetto qualitativo che riguarda più propriamente la “funzione” delle sequenze ed indica un'origine **filogenetica comune**

Dobbiamo ricorrere a criteri statistici per giudicare la significatività delle similitudini e delle differenze

Proteine omologhe: proteine che si sono evolute da un comune antecore, nell'evoluzione la similarità di sequenza è meno preservata rispetto alla struttura terziaria

Si possono avere proteine omologhe con un'identità di sequenza fino al 20%

Come è possibile ciò?

La maggior parte delle mutazioni avviene sulla superficie della proteina mentre gli amminoacidi del core sono maggiormente conservati in modo da consentire il medesimo folding alle proteine.

I problemi principali nella deduzione dei rapporti filogenetici mediante confronto tra sequenze sono:

- ⦿ Ampio ambito di variabilità della similarità, che può scendere al di sotto della significatività statistica
- ⦿ Effetto delle diverse velocità evolutive lungo i diversi rami dell'albero evolutivo

The Genetic Code

Second letter

First letter

	U	C	A	G	
U	UUU Phenylalanine UUC Phenylalanine UUA Leucine UUG Leucine	UCU Serine UCC Serine UCA Serine UCG Serine	UAU Tyrosine UAC Tyrosine UAA Stop codon UAG Stop codon	UGU Cysteine UGC Cysteine UGA Stop codon UGG Tryptophan	U C A G
C	CUU Leucine CUC Leucine CUA Leucine CUG Leucine	CCU Proline CCC Proline CCA Proline CCG Proline	CAU Histidine CAC Histidine CAA Glutamine CAG Glutamine	CGU Arginine CGC Arginine CGA Arginine CGG Arginine	U C A G
A	AUU Isoleucine AUC Isoleucine AUA Methionine; initiation codon AUG Methionine; initiation codon	ACU Threonine ACC Threonine ACA Threonine ACG Threonine	AAU Asparagine AAC Asparagine AAA Lysine AAG Lysine	AGU Serine AGC Serine AGA Arginine AGG Arginine	U C A G
G	GUU Valine GUC Valine GUA Valine GUG Valine	GCU Alanine GCC Alanine GCA Alanine GCG Alanine	GAU Aspartic acid GAC Aspartic acid GAA Glutamic acid GAG Glutamic acid	GGU Glycine GGC Glycine GGA Glycine GGG Glycine	U C A G

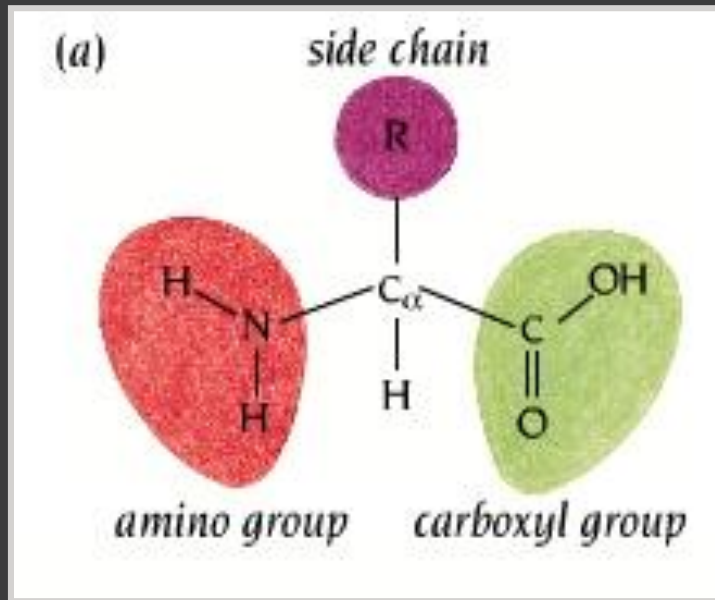
- The genetic code - Each amino acid is coded by 3 nucleotides, named codon.
- Code redundancy - Most amino acids are coded by several codons.
 - 64 triplets code for 20 amino acids & 3 stop codons.

PRINCIPI DELLA STRUTTURA DELLE PROTEINE

- ◉ Struttura primaria: gli amminoacidi e il legame peptidico
- ◉ Struttura secondaria
 - ✓ Strutture supersecondarie o motivi
- ◉ Struttura terziaria
 - ✓ Domini
- ◉ Struttura quaternaria
- ◉ Classificazione del *folding* proteico (SCOP e CATH)
- ◉ Il Protein Data Bank (PDB)

STRUTTURA PRIMARIA

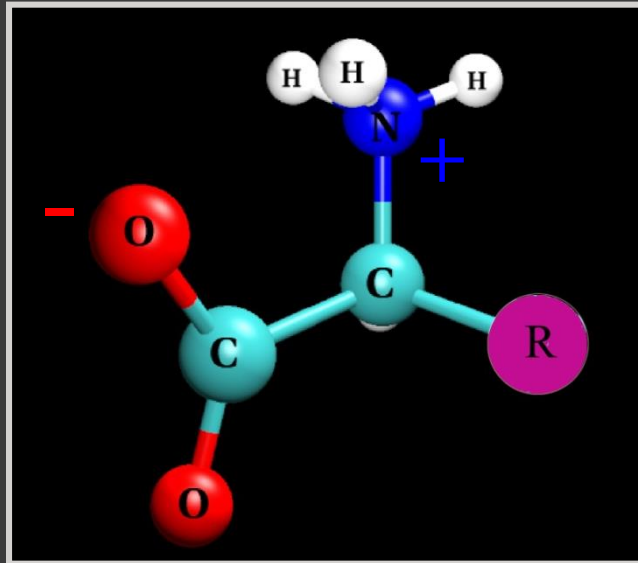
Il codice genetico specifica 20 diverse catene laterali, mentre altre possono essere prodotte dalle prime ad opera di enzimi (modificazioni post-traduzionali)



20 α -amminoacidi
standard

α -amminoacido: gruppo amminico legato al C (α) adiacente al gruppo carbossilico

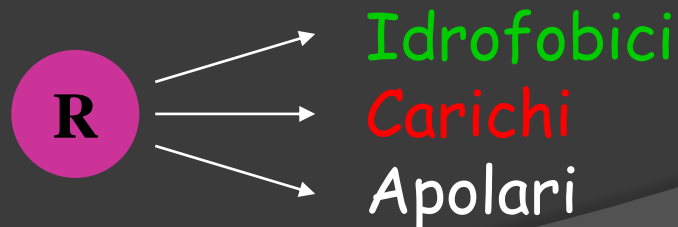
STRUTTURA PRIMARIA



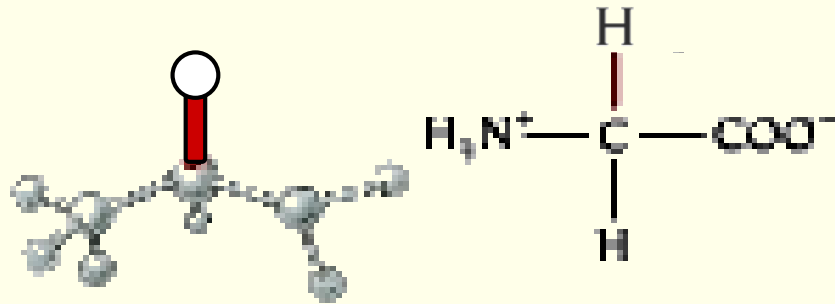
ph=7 AA anfoteri o
Zwitterioni

proprietà sia basiche che
acide

Classificazione: in base alla loro catena laterale

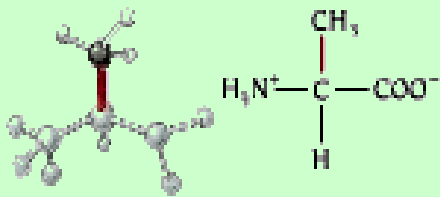
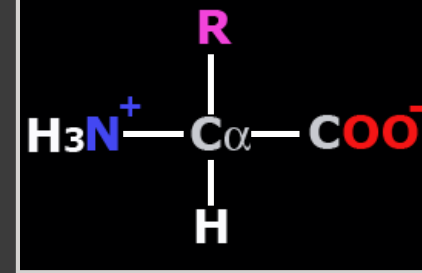


GLICINA

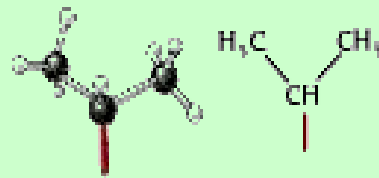


G Gly, Glicina

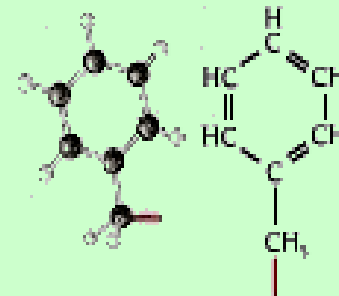
Amminoacidi IDROFOBICI



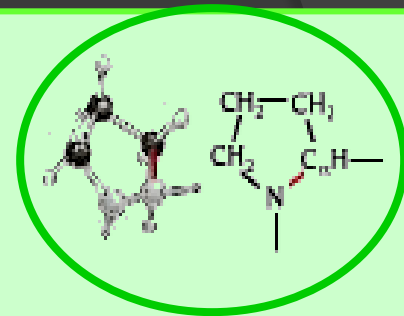
A Ala, Alanina



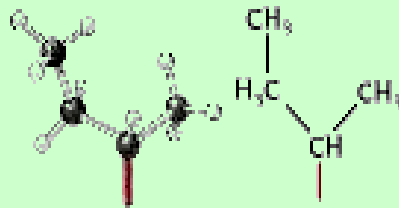
V Val, Valina



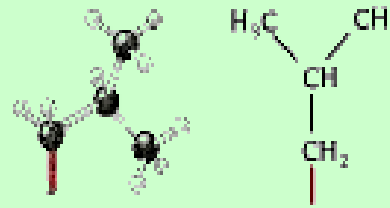
F Phe, Fenilalanina



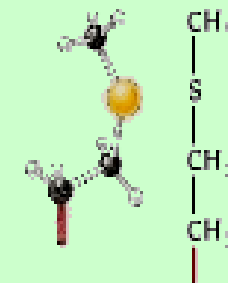
P Pro, Prolina



I Ile, Isoleucina

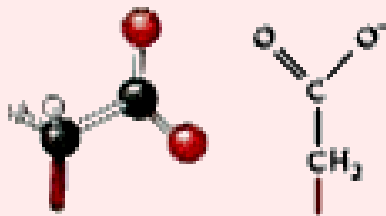
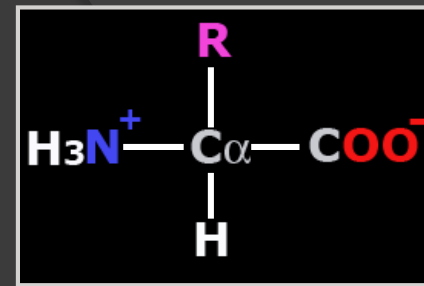


L Leu, Leucina

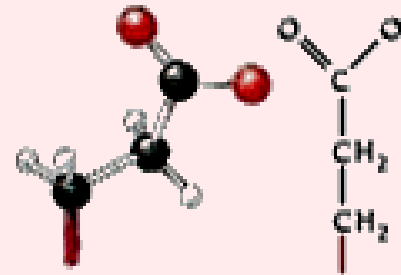


M Met, Metionina

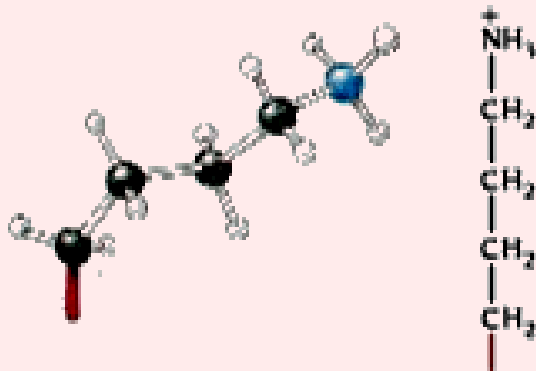
Amminoacidi CARICHI



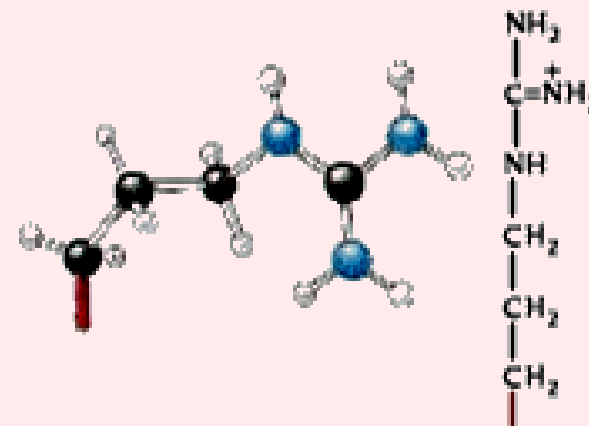
D Asp, Acido Aspartico (-)
(Aspartato)



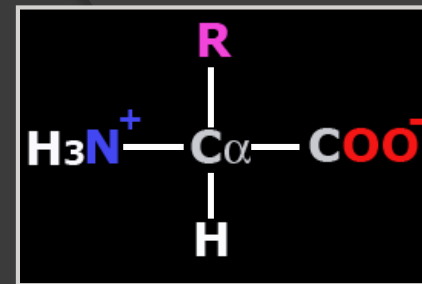
E Glu, Acido Glutammico (-)
(Glutammato)



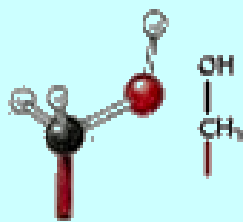
K Lys, Lisina (+)



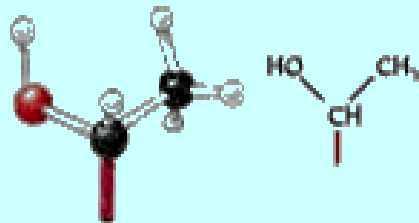
R Arg, Arginina (+)



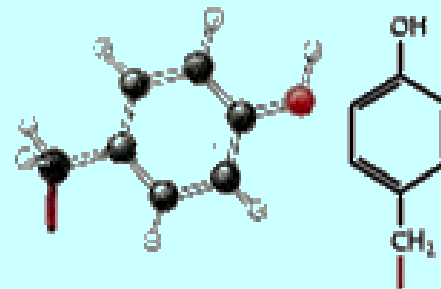
Amminoacidi POLARI



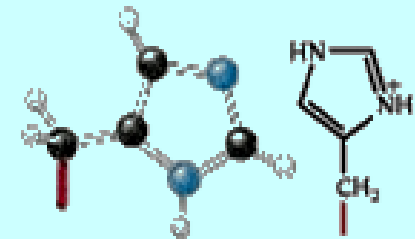
S Ser, Serina



T Thr, Treonina



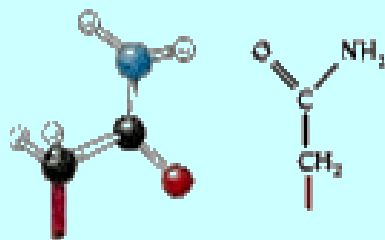
Y Tyr, Tirosina



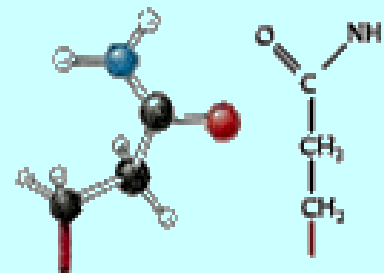
H His, Istidina



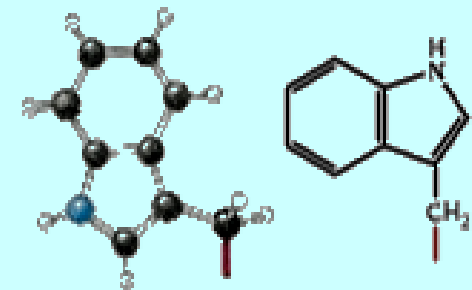
C Cys, Cisteina



N Asn, Asparagina

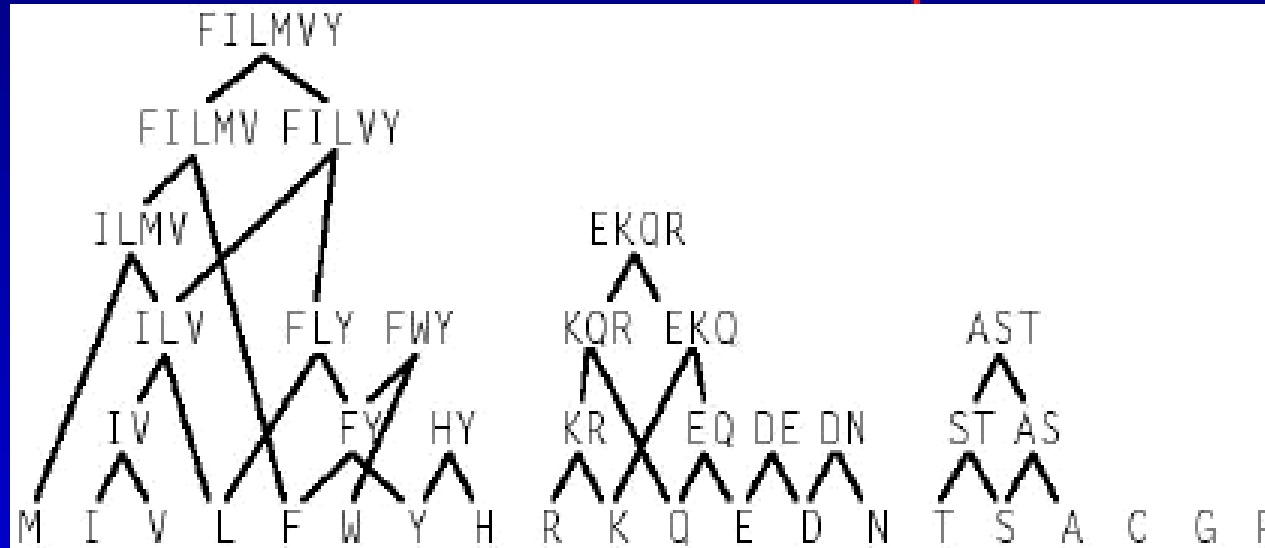


Q Gln, Glutammina



W Trp, Triptofano

Allowable Amino Acid Substitution Groups



http://www.imb-jena.de/IMAGE_AA.html

amino acid		
Alanine	ALA	A
Arginine	ARG	R
Aspartic Acid	ASP	D
Asparagine	ASN	N
Cysteine	CYS	C
Glutamic Acid	GLU	E
Glutamine	GLN	Q
Glycine	GLY	G
Histidine	HIS	H
Isoleucine	ILE	I
Leucine	LEU	L
Lysine	LYS	K
Methionine	MET	M
Phenylalanine	PHE	F
Proline	PRO	P
Serine	SER	S
Threonine	THR	T
Tryptophan	TRP	W
Tyrosine	TYR	Y
Valine	VAL	V

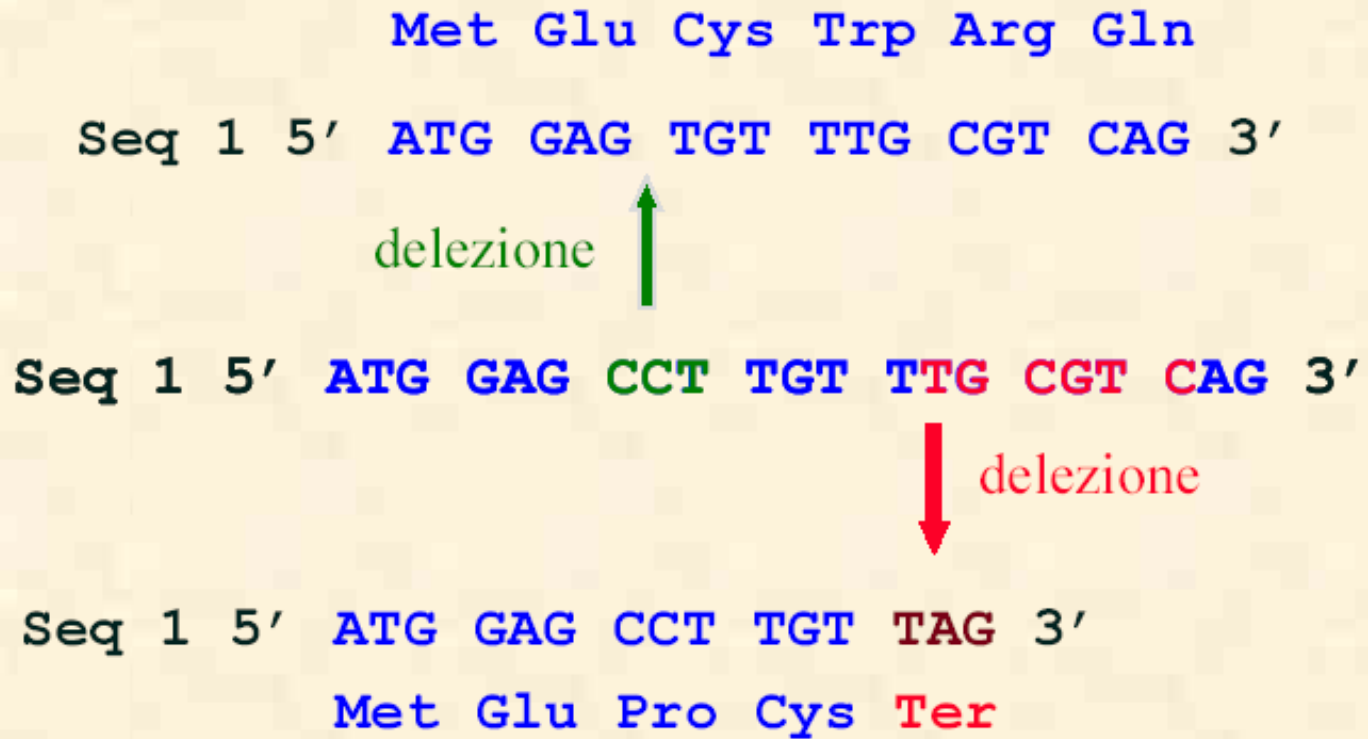
EVOLUZIONE DEI GENOMI

- Vari sono i meccanismi responsabili della variabilità genetica che oggi possiamo osservare:
 - Mutazioni puntiformi
 - Delezioni
 - Inserzioni
 - Inversioni

MUTAZIONI

		Met	Glu	Pro	Cys	Trp	Arg	Gln	
Seq 1	5'	ATG	GAG	CCT	TGT	TTG	CGT	CAG	3'
			↓		↓		↓		
		transizione			transversione			transizione	
Seq 2	5'	ATG	GAA	CCT	TCT	TTG	CGT	TAG	3'
		Met	Glu	Pro	Ser	Trp	Arg	Ter	

DELEZIONI



INSERZIONI

Seq 1a 5' Met Glu Pro His Cys Trp Arg Gln
ATG GAG CCT CAC TGT TTG CGT CAG 3'

inserzione



Seq 1 5' ATG GAG CCT TGT TTG CGT CAG 3'

inserzione



Seq 1b 5' ATG GAG CCT TGA TTT GCG TCA G 3'
Met Glu Pro Ter Phe Ala Ser

INVERSIONI

Seq 1 5' ATG GAG CCT TGT TTG CGT CAG 3'

inversione

Seq 1 5' ATG GAG ACA AGG TTG CGT CAG 3'

Met Glu Thr Arg Trp Arg Gln

Significato dell'allineamento

Qual è la corrispondenza fra gli aminoacidi delle due sequenze che più probabilmente rispecchia l'evoluzione delle due proteine?

L'allineamento tra due sequenze biologiche è utile per scoprire informazione funzionale, strutturale ed evolutiva

Cosa vuol dire allineare due sequenze?

scrivere due sequenze orizzontalmente in modo da avere il maggior numero di simboli identici o simili in registro verticale anche introducendo intervalli (gaps – inserzioni/delezioni – *indels*)

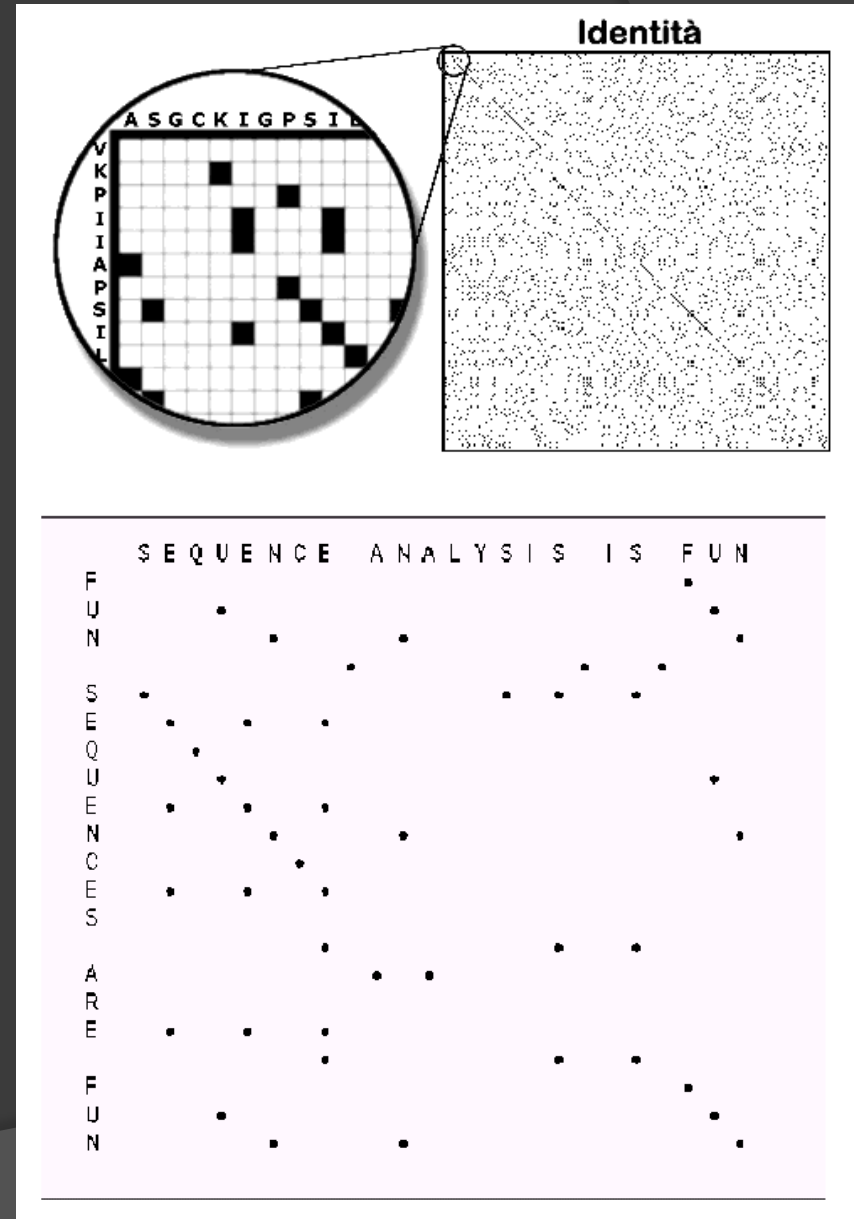
Pertanto si vogliono minimizzare le differenze o massimizzare le similarità

Metodi di allineamento

- ⦿ Analisi della matrice a punti (**dot matrix**)
- ⦿ programmazione dinamica (**dynamic programming**)
- ⦿ metodo (Fasta, Blast)

DOT MATRIX

- Il primo semplice sistema di visualizzazione di allineamenti (1970).
- Le due sequenze da confrontare sono ai margini di una matrice. Se le due lettere corrispondenti ad una casella sono uguali allora la casella viene colorata di nero ed apparirà come un punto (dot) all'interno della matrice.
- Gli allineamenti di una certa lunghezza appaiono come segmenti diagonali e saranno immediatamente distinguibili visivamente.
- I gap appaiono come salti in diagonale.
- Le sequenze ripetute appaiono come segmenti diagonali paralleli.



LE MATRICI DI SOSTITUZIONE

- Nel caso dell'allineamento di aminoacidi è opportuno applicare dei criteri di similarità che non si limitino a verificare l'identità assoluta ma tengano conto del fatto che gli aminoacidi possano essere più o meno simili tra loro. Aminoacidi molto simili possono essere indifferentemente sostituiti in una proteina senza alcuna variazione apprezzabile nella struttura della proteina.
- Per esempio acido aspartico (D) e acido glutammico (E) sono molto simili e molto spesso nel corso dell'evoluzione prendono il posto l'uno dell'altro nelle proteine. Al contrario acido aspartico (D) e triptofano (W) sono molto diversi e non sono assolutamente interscambiabili. E' quindi ragionevole valutare differenzialmente la sostituzione (in generale il confronto) di D con E e di D con W.
- Ciò viene descritto in matrici quadrate di 20*20 caselle in cui si attribuisce un punteggio ad ogni possibile coppia di aminoacidi. Quanto più alto è il punteggio tanto più interscambiabili sono gli aminoacidi. Punteggi negativi penalizzano invece aminoacidi molto differenti

metodi per la valutazione del punteggio

proposta: gli allineamenti e il calcolo della similarità potrebbero essere notevolmente migliorati dall'introduzione di schemi di punteggio diversi da 0 e da 1 per l'appaiamento di residui amminoacidici

si potrebbero per esempio prevedere punteggi alti per l'identità tra coppie di residui, punteggi un po' più bassi ma >0 per residui simili dal punto di vista chimico-fisico

punteggi invece negativi (o uguali a 0) per residui diversi o molto diversi dal punto di vista chimico-fisico

ATTENZIONE

non bisogna confondere le **matrici di punti** con le **matrici di sostituzione!**

le matrici di punti sono **grafici** che consentono di mettere in evidenza zone di identità tra sequenze diverse. Se una sequenza è lunga **m** caratteri e l'altra sequenza è lunga **n** caratteri, la matrice di punti sarà **rettangolare** e di dimensione $m \times n$

le matrici di sostituzione associano un **punteggio** ad ogni coppia di residui, sono matrici **quadrate e simmetriche**, che contengono $20 \times 20 = 400$ valori, parzialmente ridondanti (il valore relativo alla coppia RK è uguale a quello della coppia KR)

ma come si calcolano i valori di una matrice di sostituzione?

Due aminoacidi i e j frequenza della loro coppia “ f_{ij} ”, frequenza nella sequenza di i “ f_i ”
frequenza di j “ f_j ”

Rapporto “ f_{ij} ”/” f_i ” x “ f_j ” misurerà quanto spesso i e j appaiono in posizioni corrispondenti non dovute al caso
 \log_2 del rapporto valore della matrice

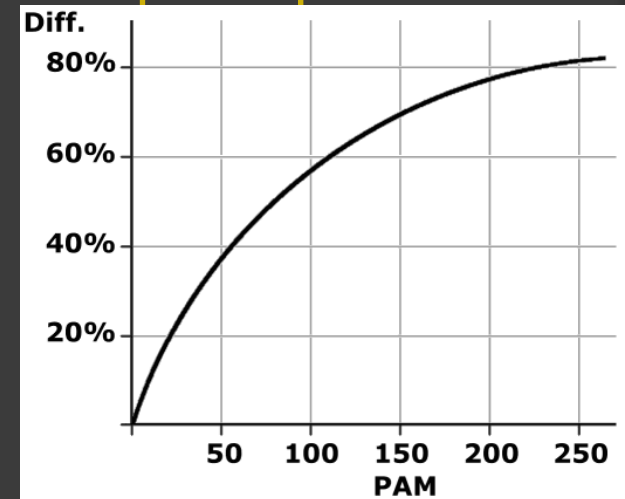
matrici di sostituzione

PAM	0	30	80	110	200	250
% identità	100	75	60	50	25	20

vediamo ora due tra le più usate matrici di sostituzione di tipo PAM: la **PAM120** e la **PAM250**, che si utilizzano per ottimizzare allineamenti tra sequenze che abbiano circa il **50%** o il **20%** di identità di sequenza

PAM

- PAM1 (con i punteggi e non con le frequenze) è molto simile alla matrice Identità (valori quasi sempre 1 sulla diagonale e 0 altrove)
- PAM2 è calcolata da PAM1 ipotizzando un altro passo evolutivo e così via...
- PAM_n è ottenuta da PAM_{n-1}
- PAM100 quindi rappresenta 100 passi evolutivi in ciascuno dei quali si è avuto un 1% di sostituzioni rispetto al passo precedente.



BLOSUM

- Introdotte da Henikoff & Henikoff nel 1992.
- A differenza delle PAM generate iterativamente, queste sono invece basate su una banca dati (BLOCKS) di allineamenti multipli di segmenti proteici senza GAP.
- Il numero associato alle matrici rappresenta la percentuale di aminoacidi identici in un certo blocco

**sequenze lontane
filogeneticamente**

**sequenze vicine
filogeneticamente**



BLOSUM35

BLOSUM62

RICERCA DELLE SIMILARITA' TRA 2 SEQUENZE

- Per determinare la similarità tra due sequenze è necessario considerare due aspetti:
 - 1- ALGORITMO DI ALLINEAMENTO
 - 2- CRITERIO DI SIMILARITA'

ALLINEAMENTI GLOBALI E LOCALI

Consideriamo i seguenti due differenti allineamenti delle stesse sequenze

Allineamento globale:

```
LTGARDWEDIPLWTDWDIEQESDFKTRAFGTANCHK
||.  | | | .|      .|  ||  || | ||
TGIPLWTDWDLEQESDNSCNTDHYTREWGTMNAHKA
```

Allineamento locale:

```
LTGARDWEDIPLWTDWDIEQESDFKTRAFGTANCHK
      ||||| | | | | | | | | | | | |
TGIPLWTDWDLEQESDNSCNTDHYTREWGTMNAHKA
```

- Nel primo caso si hanno 14 identità (evidenziate dalle linee verticali che uniscono aminoacidi uguali) e tre sostituzioni conservative (es. A-I, S-T) distribuite su tutta la lunghezza della sequenza. Nel secondo caso si hanno 13 identità ed una sostituzione conservativa su una regione di 14 aminoacidi.
- Quale dei due allineamenti è da considerarsi migliore?
- (N.B. L'allineamento globale, secondo alcuni, dovrebbe comprendere l'intera lunghezza di entrambe le sequenze; nell'esempio dovrebbero essere quindi aggiunti un gap all'inizio ed uno alla fine)

ALLINEAMENTI GLOBALI O LOCALI?

- Dal punto di vista biologico generalmente vengono privilegiati gli allineamenti locali, che riguardano regioni limitate delle proteine o di acidi nucleici.
- In Biologia Molecolare avrete sicuramente sentito parlare di domini delle proteine o anche degli acidi nucleici. Se ad esempio siamo interessati a trovare tutte le sequenze di proteine di una banca dati che contengono un certo dominio, allora sicuramente si cercheranno similarità locali.
- Gli allineamenti globali vengono applicati quando si vogliono confrontare accuratamente due sequenze in cui la similarità sia estesa per tutta la lunghezza
- N.B. Un allineamento locale non è necessariamente limitato ad una piccola regione della sequenza, ma potrebbe estendersi anche all'intera lunghezza della sequenza.

Programmazione dinamica

- ⦿ Fornisce l'allineamento ottimale tra due sequenze
- ⦿ semplici variazioni dell'algoritmo producono allineamento globali o locali
- ⦿ l'allineamento calcolato dipende dalla scelta di alcuni parametri

Allineamento globale o locale?

1) scegliamo il miglior allineamento dal punto di vista biologico, e poi...

2) cerchiamo il modo di privilegiarlo dal punto di vista computazionale

spesso gli allineamenti locali hanno una migliore rispondenza con la realtà funzionale

gli **allineamenti globali** possono comunque essere utilizzati per confrontare accuratamente due sequenze la cui **similarità sia estesa per tutta la lunghezza**

ESEMPIO

Proviamo a confrontare le due diverse proteine derivanti dallo stesso gene “subunità 1A di rubisco” di *Arabidopsis thaliana*

The screenshot shows the NCBI Entrez Gene search results for the gene *At1g67090*. The search criteria are "Gene" for "rubisco AND 1A AND arabidopsis[organism]". The results display the gene name, its full name, other aliases, chromosome, and GeneID.

NCBI Entrez Gene

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books OMIM

Search Gene for rubisco AND 1A AND arabidopsis[organism] Go Clear current records only

Limits Preview/Index History Clipboard Details

Display Summary Show: 20 Send to Text

1: [At1g67090](#) Links

ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A) (ATS1A) [*Arabidopsis thaliana*]

Other Aliases: At1g67090, F1O19.14

Chromosome: 1

GeneID: 843029

ESEMPIO

Accediamo al record del gene e attraverso i link ai due diversi record delle proteine (hanno i codici NP_176880 e NP_974098)

Display Graphics Show: 5 Send to Text

1: **At1g67090** ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A) [Links](#)
(ATSLA) [*Arabidopsis thaliana*]
GeneID: 843029 Locus tag: [At1g67090](#) updated 01-Apr-2004

Transcripts and products: (shown on reverse complement genome) [RefSeq below](#)

NC_003070

◀ 25052940] [25051821 ▶

[NM_105379](#)
[NM_202369](#)

■ - coding region ■ - untranslated region

[NP_176880](#) ribulose biphosphate carbo
[NP_974098](#) ribulose biphosphate carbo

Genomic context: chromosome: 1, map: unknown, clone: CHR1v01212004

[25055956 ▶ [25069080 ▶

At1g67070 At1g67080 At1g67090 At1g67100 At1g67110

Gene type: protein coding
RefSeq status: Provisional
Organism: [Arabidopsis thaliana](#) (ecotype: Columbia)

Link ai record dei mRNA
corrispondenti (nel database GENBANK)

Link ai record delle proteine
corrispondenti (nel database PROTEIN)

ESEMPIO

Codici **GENBANK** : 15219826 e 42572015

Codici **REFSEQ** (database composto di sequenze “pulite e verificate” di mRNA e **PROTEINE**: NP_176880 e NP_974098

```
>gi|15219826|ref|NP_176880.1| ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A)  
MASMLESATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNGGRVNCMQVWPPIGKKKFETL  
SYLPDLTDSELAKEVDYLIRNKWIPCVEFELEHGFVYREHGNSPGYYDGRYWTMWKLPFGCTDSAQVLK  
EVEECKKEYPNAFIRIIGFDNTRQVCISFIAYKPPSFTG
```

180 AA

```
>gi|42572015|ref|NP_974098.1| ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A)  
MASMLESATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNGGRVNCMQVWPPIGKKKFETL  
SYLPDLTDSELAKEVDYLIRNKWIPCVEFDLCTVSTVTHPDTMMDGTGQCGSFPCSVAPTPLKC
```

136 AA

ESEMPIO

The screenshot shows the Needle web interface. On the left is a navigation menu with categories: 'cons', 'megamerger', 'merger', 'ALIGNMENT DIFFERENCES', 'diffseq', 'ALIGNMENT DOT PLOTS', 'dotmatcher', 'dotpath', 'dottup', 'polydot', 'ALIGNMENT GLOBAL', 'alignwrap', 'est2genome', 'needle' (highlighted with a red box), 'stretcher', 'ALIGNMENT LOCAL', 'matcher', 'seqmatchall', 'supermatcher', 'water', 'wordmatch', and 'ALIGNMENT'. The main area contains a text box with a sequence: `>gi|15219826|ref|NP_176880.1| ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A) (ATS1A) [Arabidopsis thaliana] MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG GRVNCMQVWPPIGKKKFETL SYLPDLTDELAKVDYLRNKWIPCVEFELEHGFFVYREHGNSPGYDGR YWTMWKLPFLFGCTDSAQVLK`. Below this is the instruction 'Select a set of sequences. Use one of the following three fields:' followed by three numbered options. Option 1 is 'To access a sequence from a database, enter the USA path here: (dbnam)'. Option 2 is 'Or, upload a sequence file from your local computer here:' with a text box and a 'Browse...' button. Option 3 is 'Or enter the sequence data manually here:'. A second text box shows a similar sequence: `>gi|42572015|ref|NP_974098.1| ribulose biphosphate carboxylase small chain 1A / RuBisCO small subunit 1A (RBCS-1A) (ATS1A) [Arabidopsis thaliana] MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG GRVNCMQVWPPIGKKKFETL`. At the bottom, there is a 'Display percent identity and similarity' checkbox, an 'Alignment format' dropdown menu set to 'SRS format' (highlighted with a red box), a '2. SUBMIT TO NEEDLE...' button, and a 'run needle' button.

cons
megamerger
merger

ALIGNMENT DIFFERENCES

diffseq

ALIGNMENT DOT PLOTS

dotmatcher
dotpath
dottup
polydot

ALIGNMENT GLOBAL

alignwrap
est2genome
needle
stretcher

ALIGNMENT LOCAL

matcher
seqmatchall
supermatcher
water
wordmatch

ALIGNMENT

```
>gi|15219826|ref|NP_176880.1| ribulose
biphosphate carboxylase small chain 1A / RuBisCO
small subunit 1A (RBCS-1A) (ATS1A) [Arabidopsis
thaliana]
MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG
GRVNCMQVWPPIGKKKFETL
SYLPDLTDELAKVDYLRNKWIPCVEFELEHGFFVYREHGNSPGYDGR
YWTMWKLPFLFGCTDSAQVLK
```

Select a set of sequences.

Use one of the following three fields:

1. To access a sequence from a database, enter the USA path here: (dbnam)
2. Or, upload a sequence file from your local computer here:
3. Or enter the sequence data manually here:

```
>gi|42572015|ref|NP_974098.1| ribulose
biphosphate carboxylase small chain 1A / RuBisCO
small subunit 1A (RBCS-1A) (ATS1A) [Arabidopsis
thaliana]
MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG
GRVNCMQVWPPIGKKKFETL
```

Display percent identity and similarity:

Alignment format: **SRS format**

2. SUBMIT TO NEEDLE...

run needle

ESEMPIO

```
# Aligned_sequences: 2
# 1: NP_176880.1
# 2: NP_974098.1
# Matrix: EBLOSUM62
# Gap_penalty: 10.0
# Extend_penalty: 0.5
#
# Length: 181
# Identity:      108/181 (59.7%) | Aminoacidi conservati
# Similarity:   115/181 (63.5%) | Aminoacidi sostituiti in modo conservativo
# Gaps:         46/181 (25.4%)
# Score: 499.5
#
#
#=====
NP_176880.1      1  MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG      50
                |||
NP_974098.1      1  MASSMLSSATMVASPAQATMVAPFNGLKSSAAFPATRKANNDITSITSNG      50

NP_176880.1     51  GRVNCMQVWPPIGKKKFETLSYLPDLTSELAKEVDYLIRNKWIPCVEFE      100
                |||
NP_974098.1     51  GRVNCMQVWPPIGKKKFETLSYLPDLTSELAKEVDYLIRNKWIPCVEFD      100

NP_176880.1    101  LEHGFVYR-EHGNSPGYDGRYWTMWKLPLFGCTDSAQVLKEVEECKKEY      149
                :
NP_974098.1    101  TDLCTVSTVTHPDT--MMDG----TGQCGSFPCSVAPTPLK----C      136
                |
                I

NP_176880.1    150  PNAFIRIIGFDNTRQVQCISFIAYKPPSFTG      180

NP_974098.1    137
```

Quando e perché l'allineamento locale?

- 4 Confronto sequenze DNA “anonimo”, per individuare sottostringhe collegate
- 4 Individuazione subunità strutturali comuni a proteine diverse
- 4 ...

1. Ricerca di omologhe in banche dati.

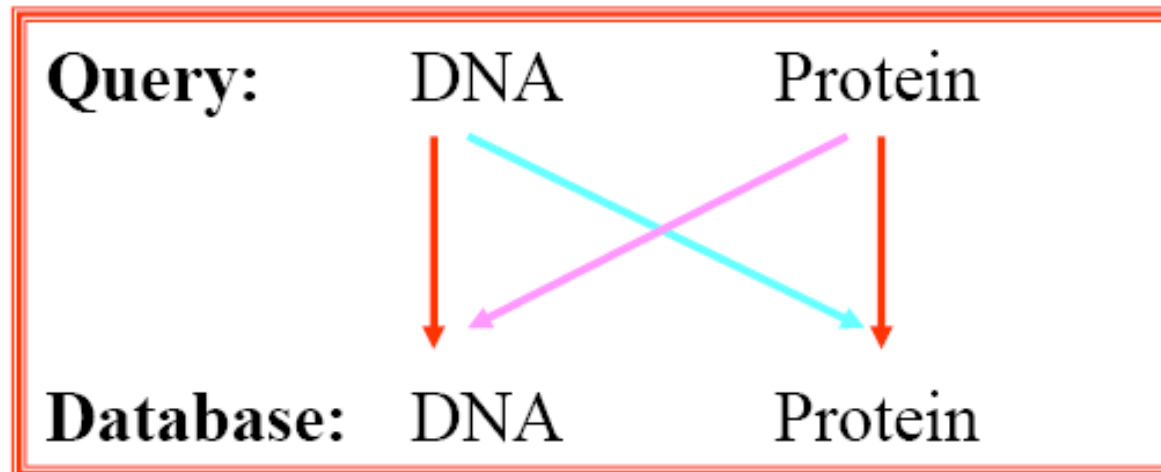
2. Programmi per la ricerca:

FASTA

BLAST

Ricerca di omologhe in banche dati

- Proteina vs. proteine
- Gene (traduzione in aa) vs. proteine
- Gene vs. geni
- Proteina vs. traduzione in aa di sequenze nucleotidiche (tutti i moduli)



Quando confrontiamo sequenze proteiche cerchiamo la migliore corrispondenza per **20** diversi **amminoacidi**

Quando confrontiamo sequenze nucleotidiche cerchiamo la migliore corrispondenza per sole **4 basi nucleotidiche**

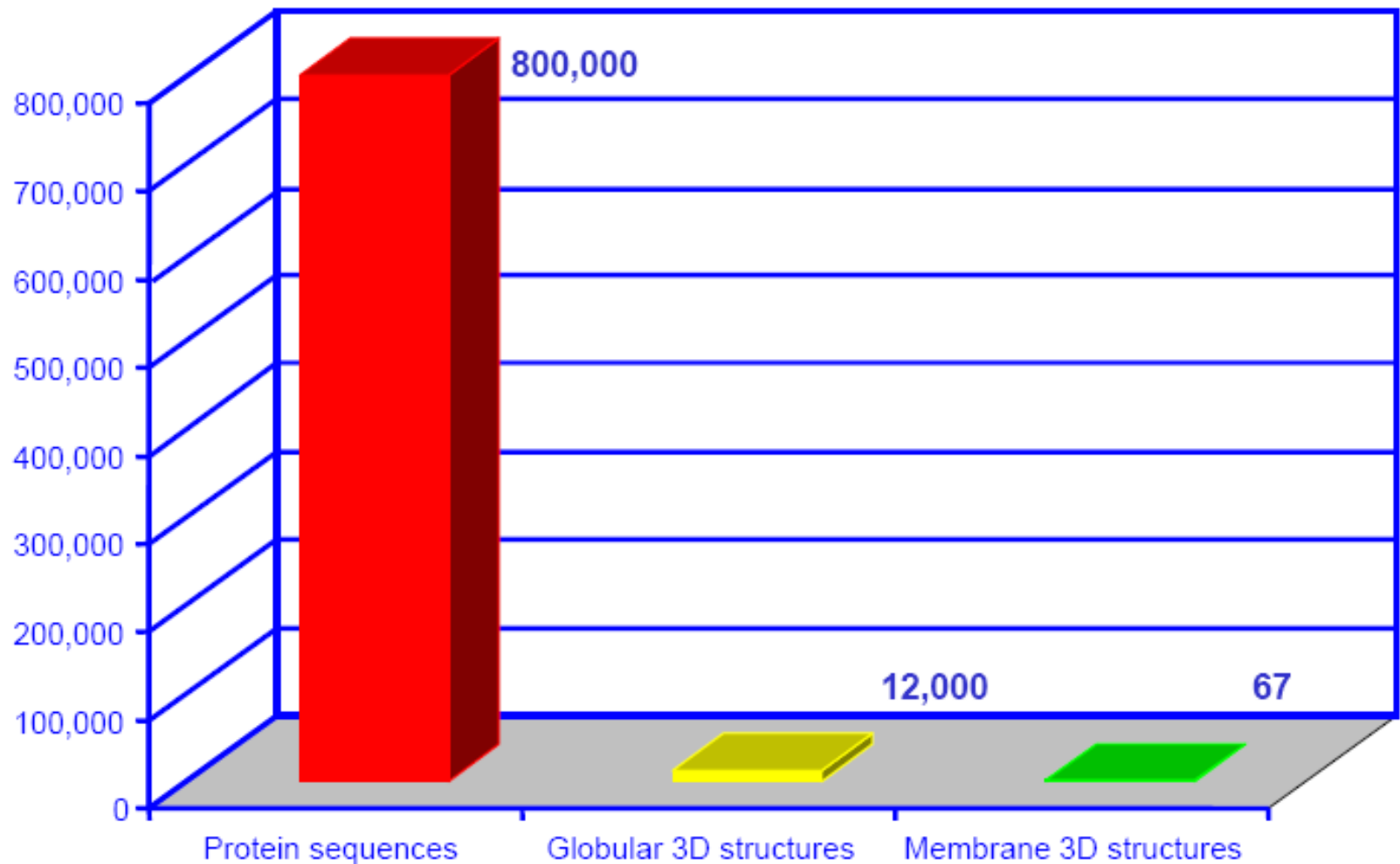
La probabilita' di trovare una buona corrispondenza (allineamento con punteggio alto) **per caso** è più alta per le sequenze nucleotidiche che per quelle proteiche

Inoltre, quando confrontiamo sequenze proteiche possiamo tener conto della **similarita'** tra i diversi amminoacidi



Quando è possibile, è preferibile confrontare sequenze proteiche !

Quante sequenze di proteine nelle banche dati ?



Come possiamo “pescare” dai databases di sequenze potenziali omologhe?



Algoritmi esatti (Smith-Waterman)

Lezione precedente

Esatto, garantisce di trovare il/i **miglior** allineamento/i per una coppia di sequenze.

Per 2 sequenze: **A** di lunghezza **n** and **B** di lunghezza **m**, Smith-Waterman impiega **$n*m$** passi computazionali.

Cerchiamo omologhe della sequenza query **A** (**$n=200$ aa**)

Cerchiamo nel DB (**10^6 sequenze di $m=200$ aa**)

Numero di passi computazionali = **$10^6 \times 200 \times 200 = \sim 10^{10}$**

10^3 passi al sec = 10^7 secs = 120 giorni = 4 mesi!

Necessità di algoritmi approssimati

Algoritmi esatti (Smith-Waterman)

Lezione precedente

Esatto, garantisce di trovare il/i **miglior** allineamento/i per una coppia di sequenze.

Per 2 sequenze: **A** di lunghezza **n** and **B** di lunghezza **m**, Smith-Waterman impiega **$n*m$** passi computazionali.

Come scartiamo gli allineamenti irrilevanti ?



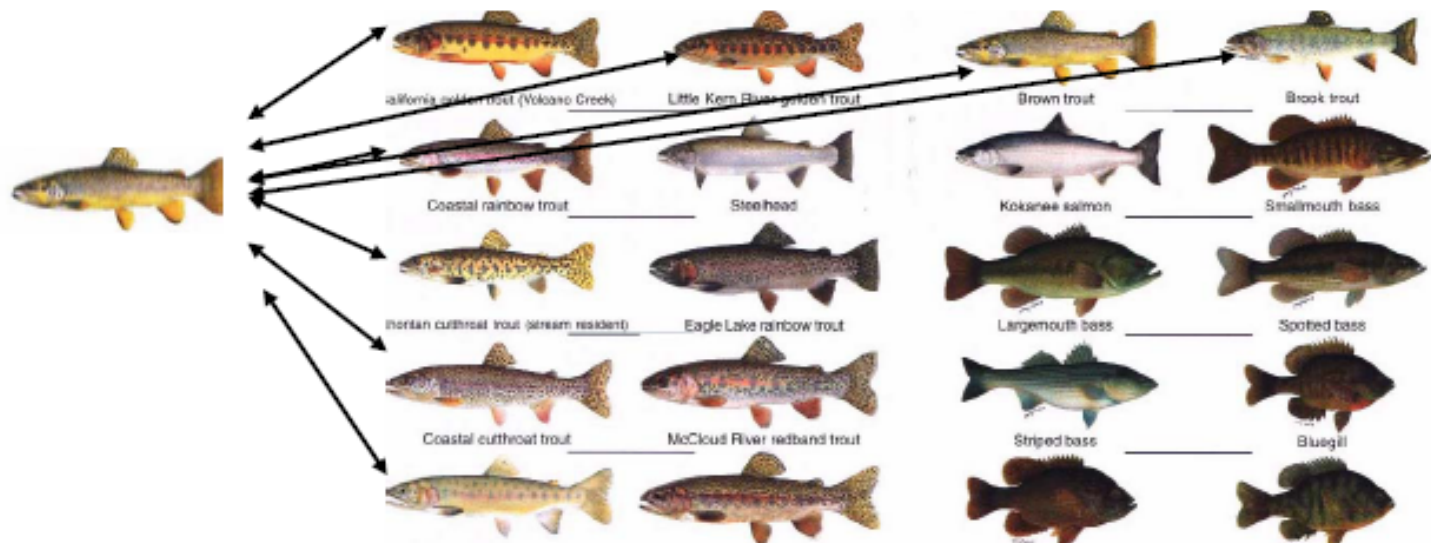
Gli algoritmi euristici (BLAST, FASTA) servono a scartare la gran parte degli allineamenti irrilevanti.

Programmi quali **FASTA** e **BLAST**, partendo da una sequenza query:



prima “pescano” dalle banche dati un sottoinsieme di sequenze che sono potenziali omologhe

poi allineano al meglio ciascuna sequenza di questo sottoinsieme alla sequenza query



Spesso, per risolvere un problema reale, si devono affrontare problemi di ottimizzazione NP-difficili, e siccome è normale che i problemi reali siano di dimensione elevata, la risoluzione di tali problemi di ottimizzazione può richiedere tempi di calcolo proibitivi. In questi casi, siccome bisogna comunque trovare una soluzione, non resta che affidarsi a procedure che non garantiscono l'ottimalità, ma che sono veloci e forniscono una soluzione spesso accettabile. Algoritmi che operano in questo modo vengono detti *algoritmi euristici* o più semplicemente euristiche.

Algoritmi EURISTICI di allineamento

Sono nati insieme alle banche dati, con lo scopo di permettere una ricerca per similarità rapida anche se meno accurata contro le migliaia di sequenze depositate.

Attualmente i programmi più utilizzati sono:

FASTA: Lipman & Pearson (1985)

BLAST: Altshul (1990)

Si tratta di sequenze omologhe?

Valutazione della significatività dell'allineamento

OPT grande

Z-score molto grande

E-value < 0.01 $1.9e-110$ significa 10 elevato alla -110 quindi un valore molto prossimo allo 0!

Bit-score grande

BLAST

Search

Ne

Search The Web

Find a Web page containing

Search

Brought to you by MSN Search

Search for other items: Files or Folders Computers People

© 2006 Microsoft MSN Privacy



Nucleotide

Protein

Translations

Retrieve results for an RID

protein-protein **BLAST**

[Search](#)

```
>gi|4504111|ref|NP_002077.1| growth factor receptor-bound
protein 2 isoform 1 [Homo sapiens]
MEAIARYDFKATADDELSFKRGDILKVLNEECDQNUYKAE LNKGEDGFI PKNVIEMKPHPW
FFGKIPRAKA
EEMLSKQRHDGAF LIRESSESAPGDFLSVKFGNDVQHFKVLRDGA GKYFLWVVKFNLSLNE
```

[Set subsequence](#) From: To:

[Choose database](#) nr

- nr
- refseq
- swissprot
- pat
- pdb
- env_nr
- month

Now:

or

Options for advanced blasting

[Limit by entries query](#) or select from:

[Compositional adjustments](#)

[Choose filter](#) Low-complexity Mask-Spaceman Mask-SPACEMAN

BLAST

Search X

No >>

Search
The
WebFind
a
Web
page
containi

Search

Brought
to
you
by
MSN
SearchSearch
for
other
items:
[Files or
Folders](#)
[Computer
People](#)©
2006
Micros
MSN
Privacy**Options** for advanced blasting

[Limit by enter query](#) or select from

[Compositional adjustments](#)

[Choose filter](#) Low complexity Mask for lookup table only Mask lower case

[Expect](#)

[Word Size](#)

[Matrix](#) Gap Costs

[PSSM](#)

- PAM30
- PAM70
- BLOSUM80
- BLOSUM62**
- BLOSUM45

[Other advanced](#)

[PHI pattern](#)

Format

Query = g|4504111|ref|NP_002077.1| growth factor receptor-bound protein 2 isoform 1 [Homo sapiens] (217 letters)

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



The request ID is 1142860123-17210-152380642238.BLAST04

Format! or **Reset all**

The results are estimated to be ready in 9 seconds but may be done sooner.

Please press "FORMAT!" when you wish to check your results. You may change the formatting options for your result via the form below and press "FORMAT!" again. You may also request results of a different search by entering any other valid request ID to see other recent jobs.

Format

Show Graphical Overview Linkout Sequence Retrieval NCBI-gi Alignment in HTML format

CDS feature

Masking Character Lower Case **Masking Color** Grey

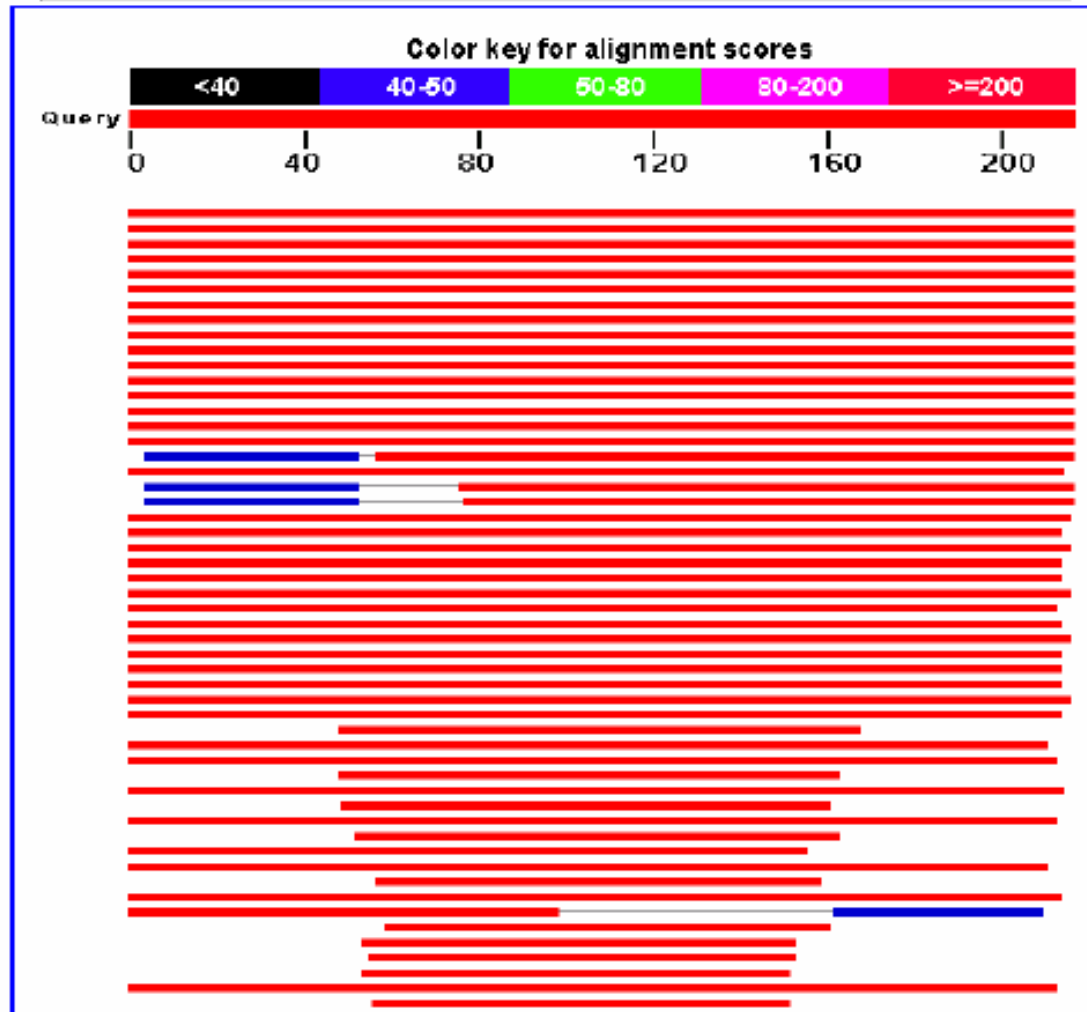
Number of: Descriptions 500 Alignments 250

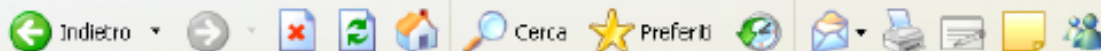
Alignment view Pairwise

I tratti di sequenza "mascherati" perché a bassa complessità sono riportati in lettere minuscole colorate in grigio

Distribution of 985 Blast Hits on the Query Sequence

Mouse-over to show define and scores, click to show alignments





Affine Phospho Peptide

[gi|20150610|pdb|1JYQ|A](#) Chain A, Xray Structure Of Grb2 Sh2 Domain Complexed With A Highly Affine Phospho Peptide
Length=96

Score = 189 bits (481), Expect = 5e-47

Identities = 92/92 (100%), Positives = 92/92 (100%), Gaps = 0/92 (0%)

```
Query 60  WFFGKIPRAKAEEMLSKQRHGDGAFLIRESESAPGDFSLSVKFGNDVQHFKVLRDGAGKYF 119
          WFFGKIPRAKAEEMLSKQRHGDGAFLIRESESAPGDFSLSVKFGNDVQHFKVLRDGAGKYF
Sbjct 5   WFFGKIPRAKAEEMLSKQRHGDGAFLIRESESAPGDFSLSVKFGNDVQHFKVLRDGAGKYF 64

Query 120 LWVVKFNSLNELVYHRSTSVSRNQQIFLRDI 151
          LWVVKFNSLNELVYHRSTSVSRNQQIFLRDI
Sbjct 65  LWVVKFNSLNELVYHRSTSVSRNQQIFLRDI 96
```

> [gi|72081594|ref|XP_788430.1](#) PREDICTED: similar to CG6033-PA, isoform A [Strongylocentrotus purpuratus]
Length=177

Score = 176 bits (451), Expect = 1e-43

Identities = 103/214 (48%), Positives = 126/214 (59%), Gaps = 43/214 (20%)

```
Query 2   EAIKAVDFKATADDELSFKRGDILKVLNEECDQNGYKAELNCKDGFIPKVIENKPHPVF 61
          EA AK+DF + ELSFK+ ILKV +DG
Sbjct 3   EATAKHDFNGQEESELSFKKMSILKVT-----RDG----- 32

Query 62  FGIKIPRAKAEEMLSKQRHGDGAFLIRESESAPGDFSLSVKFGNDVQHFKVLRDGAGKYFLW 121
          AEE+L K DGAFIRESE PGD+SLSVKF + VQHFKVLRDGAGKYFLW
Sbjct 33  -----AEELL-KNDGDGAFLIRESEGTGPGDYSLSVKFVDGVQHFKVLRDGAGKYFLW 63

Query 122 VVKFNSLNELVYHRSTSVSRNQQIFLRDIEQVPQQPTYVQALFDFDPQEDGELCFRRGD 181
          VVKFNSLN+LV+YHR++SVSR+Q I+L+D + + V AL+DF E+GEL F++GD
Sbjct 84  VVKFNSLNQLVEYHRTSSVSRSTIYLKD--RKSES IHLVLAALYDFTAGEEGELSFKKGD 141

Query 152 FIVHVDNNDPNWVKG--ACHGOTGHFPRNTVTFV 213
```

Per ricerche in banche dati nucleotidiche, l'indicizzazione in w-mers ha poca rilevanza.

Inoltre il valore w di default di Blast per i nucleotidi è 11, il che lo porta a non riconoscere sequenze che condividano in modo esatto meno di 11 basi, è questo è un limite grosso.

Fasta è molto più tollerante per sequenze che presentano gaps, visto che già nelle prime fasi prevede il loro inserimento, mentre Blast li inserisce solo in fase di allungamento.

⇒ FASTA è più adatto a ricerche in banche dati nucleotidiche

⇒ BLAST è più adatto a ricerche in banche dati proteiche

Anche se questa regola è un po' troppo arbitraria...

Perché l'allineamento multiplo?

Il modo più efficace per decifrare un testo scritto in lingua sconosciuta è quello di confrontarlo con testi equivalenti scritti in linguaggi noti (Stele di Rosetta 196 a.C)

Allineamento multiplo: confrontare più sequenze ricerca di omologhe

Omologia (carattere qualitativo): geni o proteine che hanno un comune progenitore

Ortologhi: divergenza di sequenza dovuta ad un evento di speciazione (Beta-globina ratto e umana)

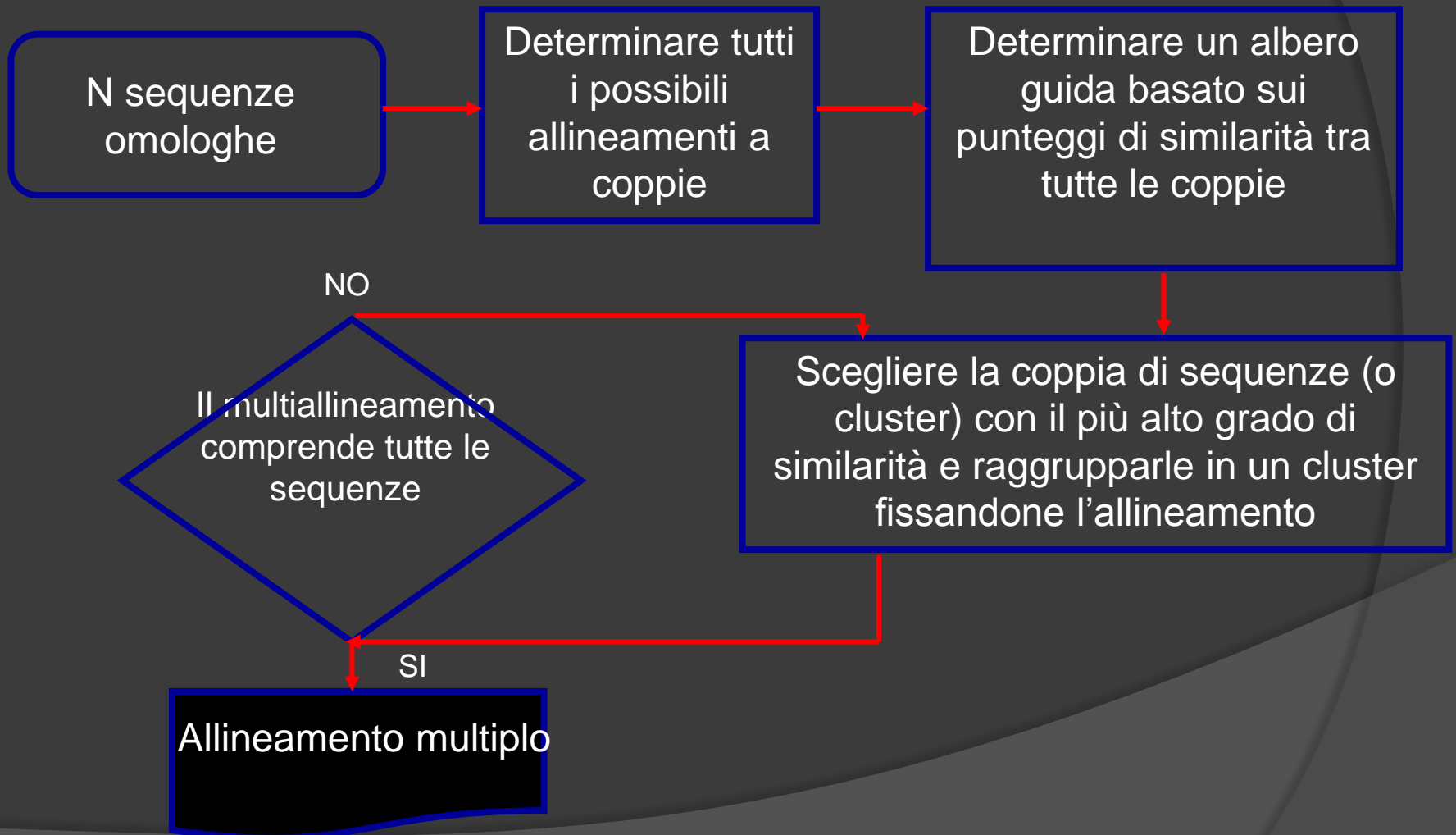
Paraloghi: divergenza in caso di duplicazione genica (alfa e beta globine umana)

Per potere affermare che due sequenze siano omologhe ci si affida al criterio di similarità. Se due sequenze possiedono un elevato grado di similarità sarà improbabile che ciò sia dovuto esclusivamente al caso.

Non si può escludere il contrario

Organismi molto lontani evolutivamente

Similarità per convergenza adattativa





Pôle BioInformatique Lyonnais Network Protein Sequence Analysis

NPS@ is the [IBCP](#) contribution to [PBIL](#) in Lyon, France

[\[HOME\]](#) [\[NPS@\]](#) [\[SRS\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#) [\[MPSA\]](#) [\[ANTHEPROT\]](#) [\[Geno3D\]](#) [\[SuMo\]](#) [\[Positions\]](#) [\[PBIL\]](#)

Wednesday, December 28th 2005: NPS@ server will not be available from January 9th to 13th. We apologize for any inconveniences.

Monday, September 26th 2005 : fixed secondary structure prediction insertion in Multalin. ([see news](#))
When sending automatic requests on NPS@, please use HTTP POST method not GET.

CLUSTALW

[\[Abstract\]](#) [\[NPS@ help\]](#) [\[Original server\]](#)

Paste a protein sequence databank in Pearson/Fasta format below : [help](#)

```
>sp|P09345|HEMA_IACKS Hemagglutinin precursor [Contains:  
Hemagglutinin HA1 chain; Hemagglutinin HA2 chain] -  
Influenza A virus (strain A/Chicken/Scotland/59 H5N1).  
MERIVLLLAIVSLVKSDQICIGYHANKSTKQVDTIMEKNTVTVTHAQDILERTHNGKLCSL  
NGVKPLILRDCSVAGWLLGNPMCDEFNLNLPWLYIVEKDNPI NSLCYPGDFNDYEELKYL  
LSSSTNHF EKIRIIPRSSWSNHDASSGVSSACPYIGRSSFLRNVVWLIKKNNTYPTIKRSY  
NNTNQEDLLILWGIHHPNDAAEQTKLYQNPTTYVSVGTSTLNQRSIPEIATRPKVNGQSG  
RMEFFWTILKPNDAINFESNGNFIAPRYAYKIVKKGDSAIMKSGLAYGNCDTKCQTPVGE  
INSSMPFHMHPHTIGCEPKYVKS DRLVLATGLRNVPQRKKRGLFGAIAGFIEGGWQGMV
```

All sequence names must be different !

SUBMIT

CLEAR

Output width :

Alignment

Hide Colors

View Alignment File

CLUSTAL W (1.82) multiple sequence alignment

```
sp|P01922|HBA_HUMAN      -VLSPADKTNVKAAWGKVGAGHAGEYGAELERMFLSFPTTKTYFFPHF-DLS-----HGSA 53
sp|P01958|HBA_HORSE     -VLSAADKTNVKAAWSKVGGHAGEYGAELERMFLGFPTTKTYFFPHF-DLS-----HGSA 53
sp|P02023|HBB_HUMAN     VHLTPEEKSAVTALWGKVNV--VDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNP 58
sp|P02062|HBB_HORSE     VQLSGEEKAAVLAALWDKVNV--EEEVGGEALGRLLVVYPWTQRFFDSFGDLSNPGAVMGNP 58
sp|P02179|MYG_BALAC     -VLSDAEWHLVLNIAKVEADVAGHQDILIRLFKGHPELEKFDKFKHLKTEAEMKASE 59
      *:  :  *  *  **      *  :  *  ** :  . * *  *  *  . * . . . .

sp|P01922|HBA_HUMAN     QVKGHGKQVADALTNAVAVHVDMPNALSALSDLHAHKLKRVDPVNFKLLSHCLLVTLAAHL 113
sp|P01958|HBA_HORSE     QVKAHGKQVGDALTLAVGHLDDLPGALSNSDLHAHKLKRVDPVNFKLLSHCLLSTLAVHL 113
sp|P02023|HBB_HUMAN     KVKAHGKQVLAGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLCVLAHHF 118
sp|P02062|HBB_HORSE     KVKAHGKQVLSHSGEGVHHLDNLKGTFAAALSELHCDKLHVDPENFRLLGNVLVVLAHF 118
sp|P02179|MYG_BALAC     DLKKHGNTVLTALGGILKKGHHEAELKPLAQSHATKHKIPIKYLEFISDAIHVLSHRH 119
      . : *  ** : . *  : :  :  :  . .  :  * : : * . *  : :  : : . . . : : . *  :

sp|P01922|HBA_HUMAN     PAEFTPAVHASLDKFLASVSTVLTISKYR----- 141
sp|P01958|HBA_HORSE     PNDFTPAVHASLDKFLSSVSTVLTISKYR----- 141
sp|P02023|HBB_HUMAN     GKEFTPPVQAAVQKVVAGVANALAHKYH----- 146
sp|P02062|HBB_HORSE     GKDFTEPELQASYQKVVAGVANALAHKYH----- 146
sp|P02179|MYG_BALAC     PAEFGADAQAAMNKALELFRKDIAAKYKELGFQG 153
      : * .  : ** : *  : . . : : ** :
```

PLEASE NOTE: Showing colors on large alignments is slow.

Hide Colors

View Alignment File

Alignment data :

Alignment length : 575

Identity (*) : 361 is 62.78 %

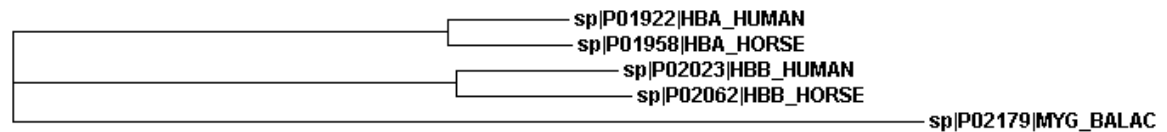
Strongly similar (:): 86 is 14.96 %

Weakly similar (.) : 32 is 5.57 %

Different : 96 is 16.70 %

```
(  
(  
sp|P01922|HBA_HUMAN:0.05940,  
sp|P01958|HBA_HORSE:0.06117)  
:0.21271,  
(  
sp|P02023|HBB_HUMAN:0.07873,  
sp|P02062|HBB_HORSE:0.08566)  
:0.21751,  
sp|P02179|MYG_BALAC:0.44687);
```

Phylogram



Show as Cladogram Tree

Show Distances

View DND File

Right-click on the above tree to see display options.

Problems printing? Read [how to print a Phylogram or Cladogram](#).

Quando non ricorrano tali condizioni, l'allineamento richiede una ulteriore fase di aggiustamento, che tenga conto di informazioni addizionali quali:

- Conservazioni di siti funzionali e catalitici noti a priori
- Predizioni di struttura secondaria
- Predizioni di struttura terziaria

La struttura di PSI-BLAST

PSI-BLAST: usa una ricerca iterativa per cui le sequenze trovate a ogni ciclo sono

Usate per costruire un modello di punteggio per il ciclo successivo.

Tutta l'informazione contenuta nel multiallineamento può essere contenuta nel suo "profilo" Tra le principali applicazioni dei profili troviamo la ricerca di sequenze omologhe molto divergenti.

1° iterazione normale ricerca di similarità con l'algoritmo BLAST

2° iterazione individuato il profilo, seconda ricerca

Risultati inaccurati se si inseriscono nell'allineamento per il profilo sequenza intruse

ESEMPIO: ferritina umana

Sequences with E-value BETTER than threshold

[Related Structures](#)

Sequences producing significant alignments:				Score (Bits)	E Value	
NEW	<input checked="" type="checkbox"/>	gi 47125326 gb AAH70494.1	FTH1 protein [Homo sapiens]	360	2e-98	U G
NEW	<input checked="" type="checkbox"/>	gi 114637918 ref XP_001140124.1	PREDICTED: similar to Ferrit...	357	2e-97	G
NEW	<input checked="" type="checkbox"/>	gi 76779199 gb AAI05803.1	FTH1 protein [Homo sapiens]	354	1e-96	U G
NEW	<input checked="" type="checkbox"/>	gi 114649149 ref XP_509574.2	PREDICTED: similar to FTH1 protein	337	1e-91	G
NEW	<input checked="" type="checkbox"/>	gi 56682959 ref NP_002023.2	ferritin, heavy polypeptide 1 [H...	335	5e-91	U G
NEW	<input checked="" type="checkbox"/>	gi 50927649 gb AAH78892.1	Fth1 protein [Rattus norvegicus]	334	1e-90	U G
NEW	<input checked="" type="checkbox"/>	gi 83404987 gb AAI11079.1	Fth1 protein [Rattus norvegicus]	334	1e-90	U G
NEW	<input checked="" type="checkbox"/>	gi 58477732 gb AAH89817.1	Fth1 protein [Rattus norvegicus] >...	334	1e-90	U G
NEW	<input checked="" type="checkbox"/>	gi 229918 pdb 1FHA 	Chain , Ferritin (H-Chain) Mutant (Lys ...	333	2e-90	S
NEW	<input checked="" type="checkbox"/>	gi 62900172 sp Q5R8J7 FRIH_PONPY	Ferritin heavy chain (Ferrit...	333	2e-90	
NEW	<input checked="" type="checkbox"/>	gi 117558589 gb AAI27508.1	Fth1 protein [Rattus norvegicus]	333	2e-90	G
NEW	<input checked="" type="checkbox"/>	gi 51859472 gb AAH81845.1	Fth1 protein [Rattus norvegicus]	333	2e-90	U G
NEW	<input checked="" type="checkbox"/>	gi 28435 emb CAA25086.1	unnamed protein product [Homo sapiens]	330	2e-89	U G
NEW	<input checked="" type="checkbox"/>	gi 42490866 gb AAH66341.1	FTH1 protein [Homo sapiens]	330	2e-89	U G
NEW	<input checked="" type="checkbox"/>	gi 109085791 ref XP_001104405.1	PREDICTED: similar to Ferrit...	329	3e-89	G
NEW	<input checked="" type="checkbox"/>	gi 17367250 sp Q9XT73 FRIH_TRIVU	Ferritin heavy chain (Ferrit...	323	2e-87	
NEW	<input checked="" type="checkbox"/>	gi 50978756 ref NP_001003080.1	ferritin, heavy polypeptide 1...	322	4e-87	U G
NEW	<input checked="" type="checkbox"/>	gi 114326408 ref NP_001041616.1	ferritin heavy chain [Felis ...	320	2e-86	G
NEW	<input checked="" type="checkbox"/>	gi 16416389 dbj BAB70615.1	ferritin heavy chain [Cavia porcellu	318	5e-86	

XP_043424 Ferritina umana

NEW	<input checked="" type="checkbox"/>	gi 73983772 ref XP_854977.1 	PREDICTED: similar to Ferritin 1...	58.9	9e-08	G
NEW	<input checked="" type="checkbox"/>	gi 106895228 ref ZP_01362330.1 	Ferritin and Dps [Clostridium...	58.5	1e-07	
NEW	<input checked="" type="checkbox"/>	gi 120528 sp P18686 FRIL SHEEP	Ferritin light chain (Ferritin L	58.5	1e-07	
NEW	<input checked="" type="checkbox"/>	gi 45358722 ref NP_988279.1 	Ferritin [Methanococcus maripalu...	58.5	2e-07	G
NEW	<input checked="" type="checkbox"/>	gi 288834 emb CAA47982.1 	ferritin 1 [Vigna unguiculata]	58.2	2e-07	
NEW	<input checked="" type="checkbox"/>	gi 54113875 gb AAV29571.1 	NT02FT1049 [synthetic construct]	57.4	3e-07	
NEW	<input checked="" type="checkbox"/>	gi 56707774 ref YP_169670.1 	Ferritin-like protein [Francisel...	57.4	3e-07	G
NEW	<input checked="" type="checkbox"/>	gi 90590871 ref ZP_01246517.1 	Ferritin and Dps [Flavobacteri...	57.0	4e-07	
NEW	<input checked="" type="checkbox"/>	gi 19553724 ref NP_601726.1 	ferritin-like protein [Corynebac...	57.0	4e-07	G
NEW	<input checked="" type="checkbox"/>	gi 30023087 ref NP_834718.1 	Ferritin [Bacillus cereus ATCC 1...	57.0	4e-07	G
NEW	<input checked="" type="checkbox"/>	gi 118170246 gb ABK71142.1 	ferritin family protein [Mycobacteri	57.0	4e-07	
NEW	<input checked="" type="checkbox"/>	gi 485051 pir PQ0614	ferritin 2 - cowpea (fragment)	57.0	4e-07	
NEW	<input checked="" type="checkbox"/>	gi 84517979 ref ZP_01005328.1 	ferritin [Prochlorococcus mari...	57.0	4e-07	
NEW	<input checked="" type="checkbox"/>	gi 89209389 ref ZP_01187815.1 	Ferritin and Dps [Bacillus wei...	56.6	5e-07	
NEW	<input checked="" type="checkbox"/>	gi 116074582 ref ZP_01471843.1 	ferritin [Synechococcus sp. R...	56.6	5e-07	
NEW	<input checked="" type="checkbox"/>	gi 30265097 ref NP_847474.1 	ferritin [Bacillus anthracis str...	56.6	6e-07	G
NEW	<input checked="" type="checkbox"/>	gi 53714336 ref YP_100328.1 	ferritin A [Bacteroides fragilis...	56.6	6e-07	G
NEW	<input checked="" type="checkbox"/>	gi 62263156 gb AAX78137.1 	unknown protein [synthetic construct]	56.2	6e-07	
NEW	<input checked="" type="checkbox"/>	gi 18148456 dbj BAB83264.1 	ferritin-like protein [Clostridium s	56.2	6e-07	
NEW	<input checked="" type="checkbox"/>	gi 29346517 ref NP_810020.1 	ferritin A [Bacteroides thetaiot...	55.8	9e-07	G
NEW	<input checked="" type="checkbox"/>	gi 116671801 ref YP_832734.1 	Ferritin, Dps family protein [A...	55.5	1e-06	G
NEW	<input checked="" type="checkbox"/>	gi 67938798 ref ZP_00531317.1 	Ferritin and Dps [Chlorobium p...	55.5	1e-06	
NEW	<input checked="" type="checkbox"/>	gi 39996409 ref NP_952360.1 	ferritin [Geobacter sulfurreduce...	55.5	1e-06	G
NEW	<input checked="" type="checkbox"/>	gi 92908944 ref ZP_01277722.1 	Ferritin and Dps [Mycobacteriu...	55.5	1e-06	
NEW	<input checked="" type="checkbox"/>	gi 42525965 ref NP_971063.1 	ferritin, putative [Treponema de...	55.5	1e-06	G
NEW	<input checked="" type="checkbox"/>	gi 52140475 ref YP_086354.1 	ferritin [Bacillus cereus E33L] ...	55.1	1e-06	G
NEW	<input checked="" type="checkbox"/>	gi 4756992 ref ZP_00240593.1 	SA1709 [Bacillus cereus G9241]...	55.1	2e-06	

Run PSI-Blast iteration 2

Run PSI-Blast iteration 3

Hit list size

[Distance tree of results](#) ^{NEW}

Sequences with E-value BETTER than threshold

[Related Structures](#)

Sequences producing significant alignments:

		Score	E	
		(Bits)	Value	

<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 47125326 gb AAH70494.1	FTH1 protein [Homo sapiens]	334	9e-91	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 114637918 ref XP_001140124.1	PREDICTED: similar to Ferrit...	332	4e-90	G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 76779199 gb AAI05803.1	FTH1 protein [Homo sapiens]	328	7e-89	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 83404987 gb AAI11079.1	Fth1 protein [Rattus norvegicus]	323	2e-87	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 50927649 gb AAH78892.1	Fth1 protein [Rattus norvegicus]	323	2e-87	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 58477732 gb AAH89817.1	Fth1 protein [Rattus norvegicus] >...	323	3e-87	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 51859472 gb AAH81845.1	Fth1 protein [Rattus norvegicus]	322	4e-87	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 117558589 gb AAI27508.1	Fth1 protein [Rattus norvegicus]	322	5e-87	G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 229918 pdb 1FHA 	Chain , Ferritin (H-Chain) Mutant (Lys ...	314	1e-84	S
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 56682959 ref NP_002023.2	ferritin, heavy polypeptide 1 [H...	313	2e-84	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 62900172 sp Q5R8J7 FRIH_PONPY	Ferritin heavy chain (Ferrit...	313	2e-84	
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 17367250 sp Q9XT73 FRIH_TRIVU	Ferritin heavy chain (Ferrit...	312	6e-84	
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	gi 50978756 ref NP_001003080.1	ferritin, heavy polypeptide 1...	311	9e-84	U G
<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>					G

	<input checked="" type="checkbox"/>	gi 30265097 ref NP_847474.1	ferritin [Bacillus anthracis str...	119	8e-26	G
	<input checked="" type="checkbox"/>	gi 76676681 ref XP_870469.1	PREDICTED: similar to EGF-like m...	118	1e-25	G
	<input checked="" type="checkbox"/>	gi 52140475 ref YP_086354.1	ferritin [Bacillus cereus E33L] ...	118	1e-25	G
NEW	<input checked="" type="checkbox"/>	gi 88711313 ref ZP_01105401.1	RsgA [Flavobacteriales bacteri...	118	1e-25	
NEW	<input checked="" type="checkbox"/>	gi 88804213 ref ZP_01119733.1	ferritin 1 [Robiginitalea bifo...	118	1e-25	
NEW	<input checked="" type="checkbox"/>	gi 47569927 ref ZP_00240593.1	SA1709 [Bacillus cereus G9241]...	117	2e-25	
	<input checked="" type="checkbox"/>	gi 67918558 ref ZP_00512155.1	Ferritin and Dps [Chlorobium l...	117	2e-25	
NEW	<input checked="" type="checkbox"/>	gi 113461064 ref YP_719132.1	ferritin like protein 2 [Haemop...	117	3e-25	G
NEW	<input checked="" type="checkbox"/>	gi 57241274 ref ZP_00369221.1	ferritin [Campylobacter lari R...	117	3e-25	
NEW	<input checked="" type="checkbox"/>	gi 15602532 ref NP_245604.1	RsgA [Pasteurella multocida subs...	117	3e-25	G
NEW	<input checked="" type="checkbox"/>	gi 89341529 ref ZP_01193772.1	Ferritin and Dps [Mycobacteriu...	117	3e-25	
NEW	<input checked="" type="checkbox"/>	gi 83016346 dbj BAE53405.1	ferritin like protein-2 [Actinobacil	116	4e-25	
NEW	<input checked="" type="checkbox"/>	gi 34540987 ref NP_905466.1	ferritin [Porphyromonas gingival...	116	4e-25	G
	<input checked="" type="checkbox"/>	gi 89209389 ref ZP_01187815.1	Ferritin and Dps [Bacillus wei...	116	5e-25	
NEW	<input checked="" type="checkbox"/>	gi 32030887 ref ZP_00133622.1	COG1528: Ferritin-like protein [H	116	6e-25	
NEW	<input checked="" type="checkbox"/>	gi 68250098 ref YP_249210.1	ferritin like protein 2 [Haemoph...	116	6e-25	G
NEW	<input checked="" type="checkbox"/>	gi 16273294 ref NP_439537.1	ferritin [Haemophilus influenzae...	115	6e-25	G
	<input checked="" type="checkbox"/>	gi 118170246 gb ABK71142.1	ferritin family protein [Mycobacteri	115	7e-25	
	<input checked="" type="checkbox"/>	gi 116671801 ref YP_832734.1	Ferritin, Dps family protein [A...	115	7e-25	G
NEW	<input checked="" type="checkbox"/>	gi 42629321 ref ZP_00154868.1	COG1528: Ferritin-like protein...	115	1e-24	
NEW	<input checked="" type="checkbox"/>	gi 113968478 ref YP_732271.1	Ferritin, Dps family protein [S...	115	1e-24	G
NEW	<input checked="" type="checkbox"/>	gi 117619373 ref YP_854579.1	ferritin [Aeromonas hydrophila ...	115	1e-24	G
	<input checked="" type="checkbox"/>	gi 108873403 gb EAT37628.1	ferritin subunit, putative [Aedes ae	114	1e-24	
NEW	<input checked="" type="checkbox"/>	gi 57242023 ref ZP_00369963.1	ferritin VC0078 [Campylobacter...	114	2e-24	
NEW	<input checked="" type="checkbox"/>	gi 89893918 ref YP_517405.1	hypothetical protein DSY1172 [De...	114	2e-24	G
NEW	<input checked="" type="checkbox"/>	gi 53732651 ref ZP_00154869.2	COG1528: Ferritin-like protein [H	114	2e-24	
NEW	<input checked="" type="checkbox"/>	gi 57867357 ref YP_189000.1	ferritin family protein [Staphy...	114	2e-24	G

Run PSI-Blast iteration 3



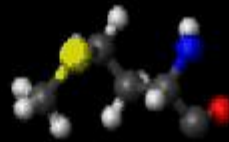
ALA



VAL



PHE



MET



PRO



LEU



ILE



ARG



LYS



ASP



GLU



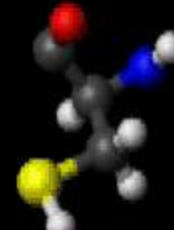
HIS



SER



THR



CYS



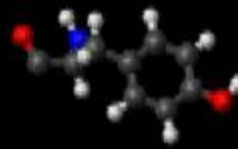
ASN



GLN



TRP



TYR

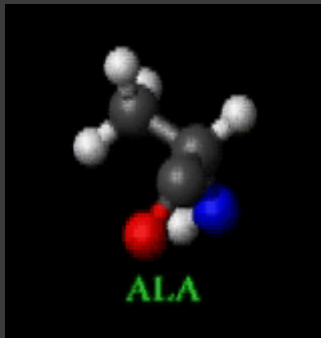


GLY

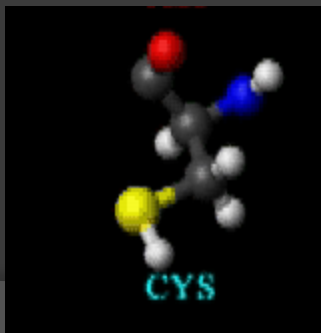
Ogni aa ha una sua particolare funzione sia a livello strutturale che di funzione 😊



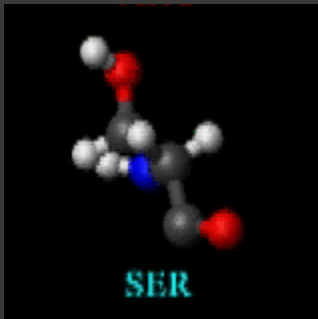
Molto flessibile (ϕ e ψ) 😊 si trova soprattutto in corrispondenza di turn, rara in eliche e sheet per l'assenza di legami H



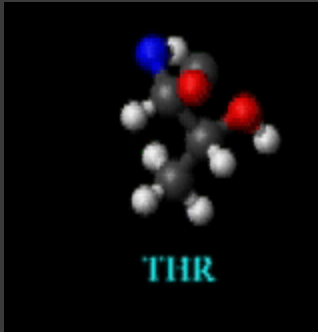
Idrofobico ma allo stesso tempo può trovarsi esposto al solvente pertanto tipico nelle zone di interazione tra proteine



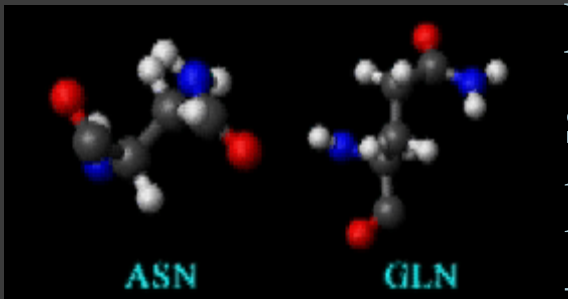
Ponti a disolfuro, coordinante metalli, attività redox (albumina)



Dentro e fuori le proteine può essere associato a deboli turns per la flessibilità e la capacità di formare legami H con il bb, può essere fosforilata...cascata di traduzione del segnale



Proprietà simili alla Ser ma essendo ramificato è meno flessibile difficile da posizionarsi nelle eliche comune negli sheet.

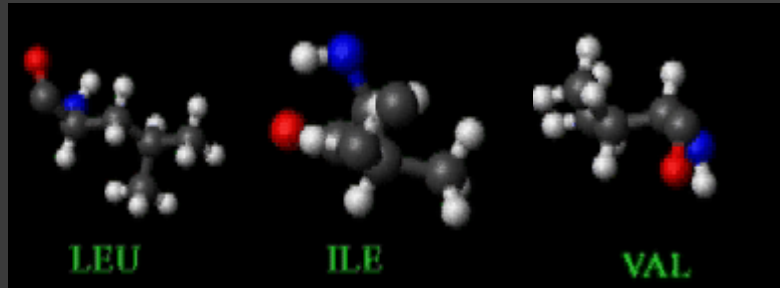


Entrambi polari si trovano di solito sulla superficie e siti attivi, Asn può formare legami H con il bb e indurre la formazione di eliche

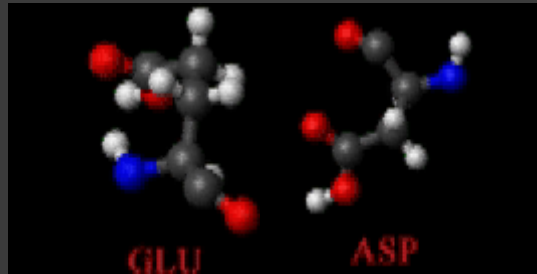


E' unico perché imminoacido, ciò restringe la sua flessibilità impedendogli di formare alcuni legami H induce un ripiegamento, nelle eliche...

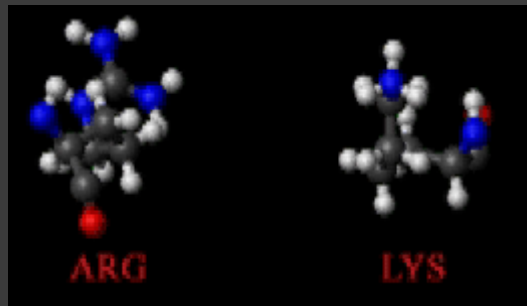
Idrofobici raramente in siti attivi
tipici degli sheet.



Carichi – di solito esposti se interni
formano ponti salini con aa carichi +, siti
attivi e coordinanti metalli (Zn)



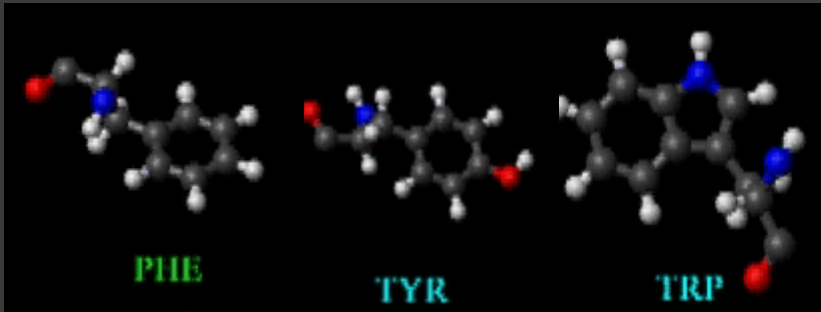
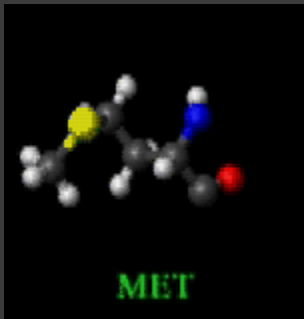
Carichi + hanno la prima parte della
catena idrofobica, possono formare ponti
salini e interagire con fosfati (ATP)



A pH fisiologico può agire come base o acido,
residuo ideale per centri funzionali, può inoltre
coordinare metalli (His Tag)



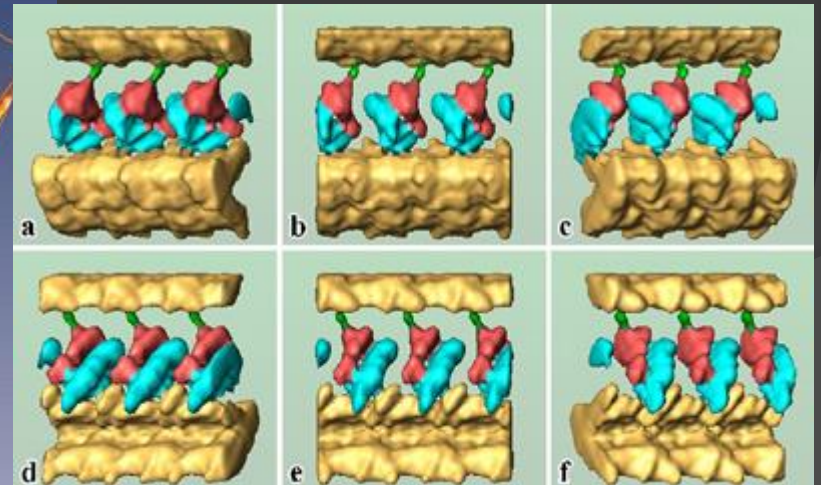
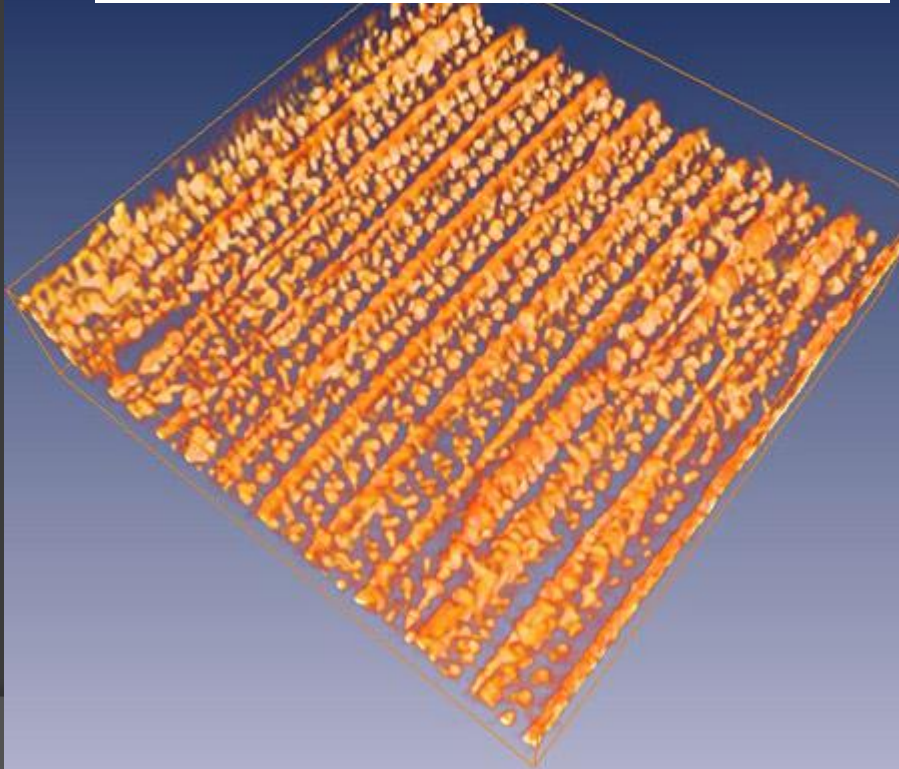
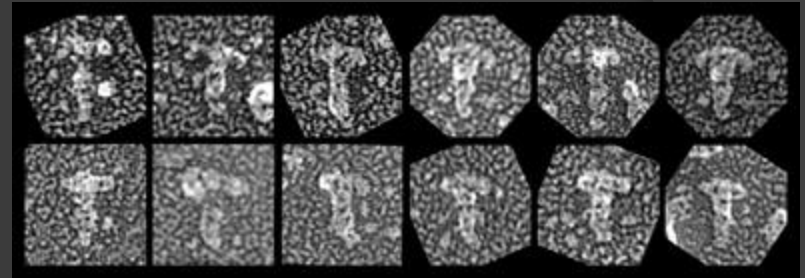
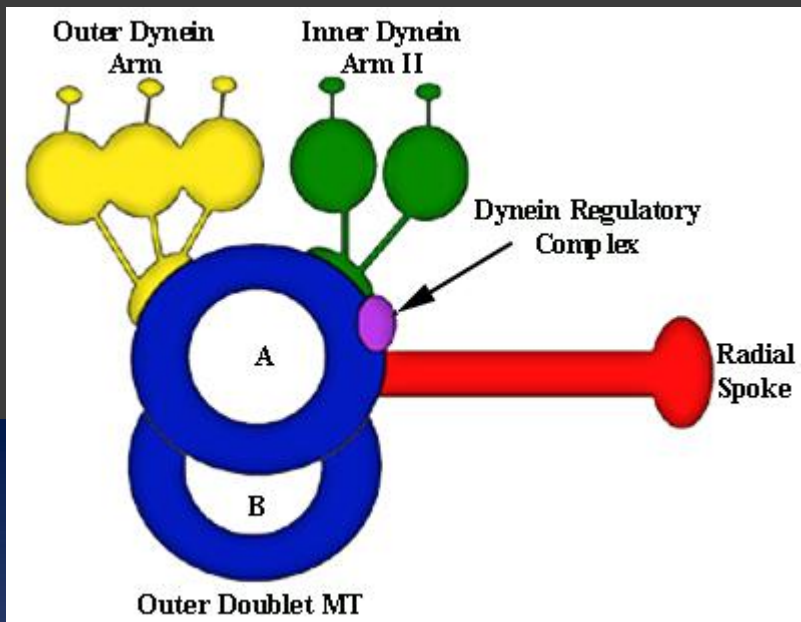
Lunga sc idrofobica e flessibile, dentro le proteine, a un atomo di S legato però ad un gruppo metile che lo rende meno reattivo.



Aromatici.

La struttura quaternaria di grandi molecole proteiche può essere studiata in 3D a partire da proiezioni ottenute al TEM lungo direzioni casuali.

Transmission Electron Microscopy -> La Microscopia Elettronica a Trasmissione (TEM) permette di ottenere, da un campione sufficientemente assottigliato (< 0.1 micron), immagini ad alta risoluzione ($< 10 \text{ \AA}$) prodotte da elettroni ad alta energia (100 KeV) trasmessi su uno schermo fluorescente o su una pellicola fotografica.



Il problema del protein folding...

$\Delta G = \Delta H - T\Delta S$ differenza tra lo stato foldato e non

Le informazioni contenute nella sequenza proteica sono sufficienti a specificare la struttura

La Structural Genomics Initiative si ripromette di determinare tra 10000 e 20000 nuove strutture proteiche nei prossimi anni e la maggior parte di esse saranno **MODELLATE**

Proteine omologhe potrebbero essere originate da duplicazione genica con differente evoluzione ed avere acquisito differenti funzioni

Alcuni ripiegamenti sono adottati da proteine che mostrano differenti funzioni

La proteina considerata potrebbe avere un ripiegamento nuovo non ancora osservato

Cosa possiamo imparare però dall'analisi della struttura?

Quali residui sono esposti e quali no

La struttura quaternaria

La struttura ottenuta sarà sicuramente molto simile a quella biologicamente attiva.

Per esempio helix-turn-helix motifs > lega il DNA

Due eliche attorcigliate per circa otto giri con Leu presenti ogni sette residui indicano un sito di dimerizzazione di molte proteine leganti il DNA 😊

Un motivo in cui lo Zn è legato a due cys e due his separate da dodici residui zink-finger

Nel caso di assenza di motivi funzionali possiamo comunque analizzare la presenza di crepe sulla superficie ed evidenziare possibili siti catalitici, non è però sempre possibile, a meno che non si conoscano altre proteine della stessa famiglia evolutiva.

Recentemente si è visto come molte proteine hanno più di una attività

Altre hanno mostrato di ripiegarsi correttamente solo quando si legano ad uno specifico ligando.

153320+1062164 sequenze proteiche conosciute
SwissProt + Trembl

26059 strutture conosciute

Tecniche lunghe per determinare la struttura delle proteine
(NMR, Cristallografia, TEM)

Incapacità di prevedere il folding di una proteina

Questo ha portato alla necessità di sviluppare una nuova metodologia.

Predizione struttura proteine:

Per farlo bisogna dividere il procedimento in step successivi

Struttura secondaria è predetta a partire dalla sequenza primaria

Gli elementi di struttura secondaria si riarrangiano per formare la struttura terziaria

Struttura secondaria

Perché è così complicato?

Per una sequenza di 100 aa, se assumiamo solo 2 possibili conformazioni per ogni residuo, ci sono $2^{100} \sim 10^{30}$ per l'intera catena

Dal momento che la sequenza determina la struttura bisogna determinare la struttura secondaria a partire dalla sequenza

Naturalmente la predizione non è così semplice

ci sono casi in cui:

Sequenze simili danno differenti strutture

Mutazioni puntiformi possono alterare la struttura

Differenti sequenze danno strutture simili

Globin fold

Predizione Struttura Secondaria

I metodi usati sono tre e si basano sulle informazioni raccolte dalle proteine la cui struttura terziaria è già risolta.

-Statistico di Chou e Fasman: i 20 aa mostrano preferenze significative per particolari strutture secondarie (A,R,Q,E,M,L,K eliche) (C,I,F,T,W,Y,V foglietti)
GOR attendibilità del 56%

-Stereochimico di Lim: tiene conto delle proprietà idrofobiche, idrofiliche ed elettrostatiche considerando il loro ruolo nel folding (alternanza di idrofilici e idrofobici, foglietti) utile per predire eliche anfipatiche e transmembrana. SOSUI, TMPRED, ecc.

-Neural Network: tiene conto di entrambe le precedenti e del processo evolutivo a partire dall'allineamento multiplo. PHD 70%

Chou-Fasman / GOR Method

Prediction: predicted sequence is scanned

- α helix prediction: when 4/6 a.a. have a probability > 1.03 to be α helix.
- β sheet prediction: when 3/5 a.a. have a probability > 1 to be β sheet.
- prediction is elongated until prediction values of 4 a.a. are < 1 .



PSORT.org provides links to the PSORT family of programs for subcellular localization prediction as well as other datasets and resources relevant to localization prediction. The page is currently hosted by the Brinkman Laboratory at Simon Fraser University, and our goal is to provide an open-source resource centre for researchers interested in subcellular localization prediction.

Please choose from the following PSORT programs for localization prediction:

- ◆ [PSORTb v.2.0](#) ([Gardy et al, 2004](#)) (v.1.0: [Gardy et al, 2003](#)) for **bacterial sequences**
- ◆ [WoLF PSORT](#) (Horton et al., to be published) is a **recently updated version of PSORT II** for the prediction of **eukaryotic sequences**
- ◆ [PSORT II](#) ([Nakai and Horton, 1997](#)) for **eukaryotic sequences**
- ◆ [PSORT](#) ([Nakai and Kanehisa, 1991](#)) for **plant sequences**
- ◆ [iPSORT](#) ([Bannai et al, 2002](#)) for **classification of eukaryotic N-terminal sorting signals**

See also [PSORTdb](#), our new database of bacterial protein subcellular localizations.

PSORTb and PSORTdb are maintained by the [Brinkman Laboratory](#), [Simon Fraser University](#), British Columbia, Canada. PSORT and PSORT II are maintained by [Kenta Nakai](#), at the [Human Genome Center](#), Institute for Medical Science, University of Tokyo, Japan. iPSORT is maintained by [Hideo Bannai](#) at the [Human Genome Center](#).



[Submit Sequences](#) | [Documentation](#) | [Resources](#) | [Contact](#) | [Updates](#)

PSORT-B Results

SeqID: ecoli

Analysis Report:

SCL-BLASTe- Periplasmic, CytoplasmicMembrane[matched 100% [129663](#): Cytoplasmic membrane associated periplasmic protein]

Localization Scores:

Cytoplasmic	0.00
CytoplasmicMembrane	10.00
Periplasmic	10.00
OuterMembrane	0.00
Extracellular	0.00

Final Prediction:

CytoplasmicMembrane (This protein may have multiple localization sites.) 10.00

TargetP 1.1 Server

Predice la localizzazione subcellulare di proteine eucariotiche .
Testato per la prima volta su Arabidopsis Thaliana (per le piante) e Homo Sapiens (per le non piante).

La localizzazione è basato sulla presenza delle presequenze N-terminali, cioè:

- *peptide targeting di cloroplasti (cTP)*, per tale predizione è utilizzato **ChloroP** che predice la presenza di sequenze segnale e la potenziale localizzazione del sito di taglio
- *peptide targeting di mitocondri (mTP)*
- *peptide segnale della via secretoria (SP)*, per tale predizione si usa **SignalP** che come **ChloroP** predice la presenza di sequenze segnale e la potenziale localizzazione del sito di taglio, ciò lo esegue in sequenze aminoacidiche di organismi diversi.

Come si usa TargetP???

1-Inserire i dati

Paste a single sequence or several sequences in [FASTA](#) format into the field below:

Submit a file in [FASTA](#) format directly from your local disk:

Sfoggia...

Organism group

Non-plant
 Plant

Prediction scope

Perform cleavage site predictions

Cutoffs

no cutoffs; winner-takes-all (default)
 specificity >0.95 (predefined set of cutoffs that yielded this specificity on the TargetP test sets)
 specificity >0.90 (predefined set of cutoffs that yielded this specificity on the TargetP test sets)
 define your own cutoffs (0.00 - 1.00): cTP: mTP: SP: other:

Submit Clear fields

Restrictions:
At most 2,000 sequences and 200,000 amino acids per submission; each sequence not more than 4,000 amino acids.



Codice ad una lettera



Percorso di directory



Paste a single sequence or several sequences in **FASTA** format into the field below:

Submit a file in **FASTA** format directly from your local disk:

Organism group

- Non-plant
- Plant

Prediction scope

- Perform cleavage site predictions

Cutoffs

- no cutoffs; winner-takes-all (default)
- specificity >0.95 (predefined set of cutoffs that yielded this specificity on the TargetP test sets)
- specificity >0.90 (predefined set of cutoffs that yielded this specificity on the TargetP test sets)
- define your own cutoffs (0.00 - 1.00): cTP: mTP: SP: other:

Restrictions:

At most 2,000 sequences and 200,000 amino acids per submission; each sequence not more than 4,000 amino acids.

Demanded specificity	Category	Cutoff	
		Plant	Non-plant
0.95	cTP	0.73	-
	mTP	0.86	0.78
	SP	0.43	0.00
	other	0.84	0.73
0.90	cTP	0.62	-
	mTP	0.76	0.65
	SP	0.00	0.00
	other	0.53	0.52

Si può definire anche il range di cutoffs per cTP, mTP e SP che va da 0,00 a 1,00 (range valori di output)

Risultati

Nome max 20 caratteri

Lunghezza sequenza

Predizione di localizzazione basata sugli score

C,M,S,-,*

Reliability class

```
### targetp v1.1 prediction results #####
Number of query sequences: 12
Cleavage site predictions included.
Using PLANT networks
```

Name	Len	cTP	mTP	SP	other	Loc	RC	TPlen
P11043_has_a_very_ve	516	0.873	0.012	0.004	0.320	C	3	65
P07505	266	0.330	0.047	0.004	0.444	-	5	-
P12360	246	0.580	0.119	0.210	0.089	C	4	42
P12352	97	0.397	0.555	0.014	0.150	M	5	40
Q01289	399	0.733	0.017	0.031	0.462	C	4	62
P08817	129	0.844	0.092	0.089	0.015	C	2	47
P07263	546	0.400	0.380	0.075	0.020	C	5	41
P07597	117	0.005	0.095	0.967	0.006	S	1	26
P48786	1088	0.199	0.070	0.067	0.822	-	2	-
Q01238	102	0.420	0.277	0.033	0.164	C	5	41
P35334	342	0.055	0.010	0.968	0.041	S	1	29
P13086	333	0.053	0.905	0.045	0.034	M	1	21
cutoff		0.000	0.000	0.000	0.000			

Score finale su cui è basata la predizione finale. La localizzazione con un alto score è quella più in accordo con TargetP

Predizione lunghezza presequenza, mostarta solo quando a TargetP si chiede la localizzazione del sito di taglio

Esempi

>P16096; 46 FRUCTOSE-BISPHOSPHATE
ALDOLASE, CHLOROPLAST PRECURSOR

```
### targetp v1.1 prediction results #####
Number of query sequences: 1
Cleavage site predictions not included.
Using PLANT networks.

Name          Len      cTP      mTP      SP  other  Loc  RC
-----
394            0    0.090    0.112    0.951  0.000
-----
cutoff          0.000    0.000    0.000    0.000
```

>P00508; 22 ASPARTATE AMINOTRANSFERASE,
MITOCHONDRIAL PRECURSOR

```
### targetp v1.1 prediction results #####
Number of query sequences: 1
Cleavage site predictions not included.
Using NON-PLANT networks.

Name          Len      mTP      SP  other  Loc  RC
-----
423            0        0.031    0.978  0.000
-----
cutoff          0.000    0.000    0.000
```

Predizioni di struttura **secondaria**

**CHOU &
FASMAN**

Residue	helix	β -sheet	turns
Glu	1.59	0.52	1.01
Ala	1.41	0.72	0.82
Leu	1.34	1.22	0.57
Met	1.30	1.14	0.52
Gln	1.27	0.98	0.84
Lys	1.23	0.69	1.07
Arg	1.21	0.84	0.90
His	1.01	0.80	0.81
Val	0.90	1.87	0.41
Ile	1.09	1.67	0.47
Tyr	0.74	1.45	0.76
Cys	0.66	1.40	0.54
Trp	1.02	1.35	0.65
Phe	1.16	1.33	0.59
Thr	0.76	1.17	0.90
Gly	0.43	0.58	1.77
Asn	0.76	0.48	1.34
Pro	0.34	0.31	1.32
Ser	0.57	0.96	1.22
Asp	0.99	0.39	1.24



Pôle BioInformatique Lyonnais

Network Protein Sequence Analysis

NPS@ is the [IBCP](#) contribution to [PBIL](#) in Lyon, France

[\[HOME\]](#) [\[NPS@\]](#) [\[SRS\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#) [\[MPSA\]](#) [\[ANTHEPROT\]](#) [\[Geno3D\]](#) [\[SuMo\]](#) [\[Positions\]](#) [\[PBIL\]](#)

Wednesday, December 28th 2005: NPS@ server will not be available from January 9th to 13th. We apologize for any inconveniences.

Monday, September 26th 2005 : fixed secondary structure prediction insertion in Multalin. ([see news](#))
When sending automatic requests on NPS@, please use HTTP POST method not GET.

GOR IV SECONDARY STRUCTURE PREDICTION METHOD

[\[Abstract\]](#) [\[NPS@ help\]](#) [\[Original server\]](#)

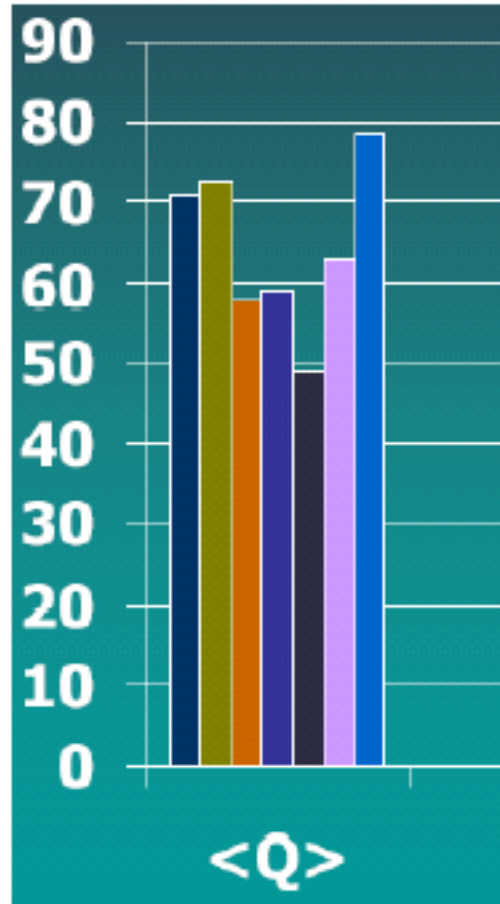
Sequence name (optional) :

Paste a protein sequence below : [help](#)

```
NGVKPLILRDCSVAGWLLGNPMCDEF INVPEWSYIVEKASPANDLCYPGNFNDYEELKHL
LSRINHFEKIQIIPKSSWSNHDASSGVSSACPYLGRSSFFRMVVWL IKKNSAYPTIKRSY
NNTNQEDLLVLWGVHHPNDAAEQTKLYQNPTTYISVGTSTLNQRLVPEIATRPKVNGQSG
RHEFFWTILKPNDAINFESNGNFI APEYAYKIVKKG DSTIMKSELEYGNCNTKCQTPMGA
INSSMPFHNIHPLTIGCEPKYVKS NRLVLTGLRNTPQRE RRRKKRGLFGA IAGFIEGGW
QGMVDGWYGYHHSNEQGSCYSADKESTQKAIDGVTNKVNSIINKMNTQFEAVGREFNLE
RRIENLNKKMEDGFLDVWVTYNAELLVLMENERTLDFHDSNVKNL YDKVRLQLRDNAKELG
NGCFEFYHKCDNECMESVKNGTYDYPQYSEE ARLNREEISGVKLESMGTYQILSIYSTVA
SSLALAIMVAGLSLWMCNSNGSLQCRICI
```

Output width :

Algoritmi predizione struttura secondaria



■ PhD (Rost, Sander)

■ PhD3 (Rost, Sander)

■ COMBINE

■ GORIII (Garnier, Osguthorpe, Robson)

■ Chou / Fassman

■ PrISM (Yang)


■ PSIPRED (Jones)

The PSIPRED protein structure prediction server - Microsoft Internet Explorer

File Modifica Visualizza Preferiti Strumenti ?

Indietro Avanti Termina Aggiorna Pagina iniziale Cerca

Indirizzo <http://bioinf.cs.ucl.ac.uk/psipred/> Vai

Bioinformatics Group 

University College of London

[Jones home>](#)
[McGuffin home>](#)
[Bryson home>](#)

The PSIPRED Protein Structure Prediction Server

David T. Jones, Liam J. McGuffin & Kevin Bryson


Description

[PSIPRED Server Help Page](#)

[PSIPRED Server History](#)

The PSIPRED protein structure prediction server allows you to submit prediction of your choice and receive the results of the prediction via three prediction methods to apply to your sequence. PSIPRED - a hidden Markov model based method, MEMSAT - our widely used transfer method and GenTHREADER - a sequence profile based fold recogni

PSIPRED server

Bioinformatics Group 

[PSIPRED home>](#)

The PSIPRED Protein Structure Prediction Server

Info

We suggest that you do not bookmark this page as it is liable to move. It is best to access the server via the [PSIPRED home page](#), which has more information about the methods and a full reference list.

Input Sequence

[Help](#)
Input sequence (single letter code)

Choose Prediction Method

- [Help](#) Predict Secondary Structure (PSIPRED v2.6)
- Predict Transmembrane Topology (MEMSAT3)
- Fold Recognition(GenTHREADER - quick)
- Fold Recognition (mGenTHREADER - with profiles and predicted secondary structure)

[CLICK HERE TO ACCESS THE SERVER](#)

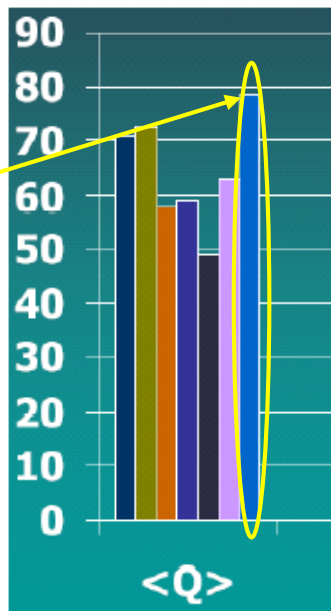
Predizione struttura secondaria

Il metodo PSIPRED tiene conto di:

- metodo statistico di Chou e Fasman;
- proprietà chimico-fisiche degli aa;
- allineamento multiplo effettuato da PSI-BLAST.

Residue	helix	β -sheet	turns
Glu	1.59	0.52	1.01
Ala	1.41	0.72	0.82
Leu	1.34	1.22	0.57
Met	1.30	1.14	0.52
Gln	1.27	0.98	0.84
Lys	1.23	0.69	1.07
Arg	1.21	0.84	0.90
His	1.01	0.80	0.81
Val	0.90	1.87	0.41
Ile	1.09	1.67	0.47
Tyr	0.74	1.45	0.76
Cys	0.66	1.40	0.54
Trp	1.02	1.35	0.65
Phe	1.16	1.33	0.59
Thr	0.76	1.17	0.80

Attendibilità dei risultati molto elevata



- PhD (Rost, Sander)
- PhD3 (Rost, Sander)
- COMBINE
- GORIII (Garnier, Osguthorpe, Robson)
- Chou / Fassman
- PrISM (Yang)
- PSIPRED (Jones)

PSIPRED Protein Structure Prediction Server - Microsoft Internet Explorer

Indietro Avanti Termina Aggiorna Pagina iniziale Cerca

Indirizzo <http://bioinf.cs.ucl.ac.uk/psipred/psiform.html>

Sequenza formato FASTA

Input Sequence

[Help](#)
Input sequence (single letter code)

```
DASEFVGVGSI AAGGRYDNLVNMVFSEASGKKSTQIPCVGISFGVERIFSLIKQRINSSTT
IKPTATQVFMVMAFGGKDWGTGYLPERMKVTKQLWDAGIEAEYVYKAKANPRKQFDAAEKA
GCHIAVILGKEEYLEGKLRVKRLGQEFADDDGELVSAADIVP IVQEKLSQIHEDGLNEVT
RLIKGL
```

Choose Prediction Method

[Help](#)

- Predict Secondary Structure (PSIPRED v2.6)
- Predict Transmembrane Topology (MEMSAT3)
- Fold Recognition(GenTHREADER - quick)
- Fold Recognition (mGenTHREADER - with profiles)

Filtering Options

[Help](#)

- Mask low complexity regions
- Mask transmembrane helices
- Mask coiled-coil regions

Warning: Turn off all filtering if you are running MEMSAT3

Submit Sequence

E-mail address [Help](#)

Password (only required for commercial e-mail addresses)

Short name for sequence [Help](#)

Operazione completata

The PSIPRED Server Help Page - Microsoft Internet Explorer

Input Sequence

Type your AMINO ACID sequence here. Please do not try to enter a nucleic acid sequence. We recommend that you enter your sequence as a plain single-letter string like this:

```
ALGSNLTNPVEQLHAALKAISQLSNLHTLUTTSFVKSEKPLGPQDQPDYVWAWAKIETELS
```

Alternatively, you can enter your sequence in FASTA format, but the description text will be ignored by the server.

Note that there is an upper limit to the length of sequences which can be submitted. For mGenTHREADER that limit is 1000 residues. For the other methods, the limit is 1500 residues. If your sequence is longer than this, try breaking it into likely domains before submitting it. Our [DomPred](#) server can help you in doing this.

Choose Prediction Method

Select which method you wish to use. See the [PSIPRED home page](#) for more details on the different methods (and references).

Filtering Options

To reduce the false positive rate of fold recognition methods, particularly when applied to long sequences, it is important that biased regions of the target sequence are filtered out before the prediction is carried out. The PSIPRED server uses the PFILT program to perform the masking and has 3 filtering options, which will filter out low complexity regions, likely transmembrane segments and coiled-coil regions. The default setting is for just low-complexity regions of the sequence to be masked out. Regions which are masked out will

Risultati predizione (formato di testo)

PSIPREDPREDICTION RESULTS

Key

Conf: Confidence (0=low, 9=high)
Pred: Predicted secondary structure (H=helix, E=strand, C=coil)
AA: Target sequence

PSIPRED HFORMAT (PSIPRED V2.6 by David Jones)

Conf: 95433564133311487702778884256899999985568789750334412689997
Pred:
CCCCCCEEEHHHHHHHCC CCCCCHHHHHHHHHHHHCC CCCCCHHHHHHCCCC
CC
AA:
MLSRSLNKVVTSIKSSIIIRMSATAAATSAPTANAANALKASKAPKKGKLVSLKTPKG
10 20 30 40 50 60

Conf: 4215999999999999999999849878016554508774444564521069999589
Pred:
CCCCCHHHHHHHHHHHHHHHHHHHHCC CCEEEBCC CCCCCHHHHCCCC CCHHEEEEEE
CC
AA:
TKDWADSDMVR EAFIS TLS GLFKKHGGVTIDTP VFELREILAGKYGEDSKLIYNLEDQG
70 80 90 100 110 120

Conf: 98897588562999999957997501698745684478754588635112517799648
Pred:
CCEEECC CCHHHHHHHHCC CCEEEBCC CCEBCC CCCC CCCC CEEEC EEEBCC
AA:
GELCSLR YDL TVPFARYVAMNNIQSIKRYHIAKVYRRDQ PAMTKGRMREFYQCDFDVAGT
130 140 150 160 170 180

Conf: 99658799999999999977998759997785668999984999999999999986
Pred:
CCCHHHHHHHHHHHHHHHHHHCC CCEEEBCC HHHHHHHHHHCC CCHHHHHHHHHH
HHH
AA: FESMVPDS ECL SILVEGLTSLGKDFKIKLNRKILDGFIQAGVKDEDVRKISSAVDKL
190 200 210 220 230 240

Conf: 40209999998641348788999998899862987899999984644345665899999
Pred:
CCCCHHHHHHHHHHHCC CCCCCHHHHHHHHHHHHCC CCHHHHHHHHHHCC CCCCCHHHH
HHHH

250 260 270 280 290 300

Conf: 99999999973998279986410158724586599999568888665420002566755
Pred:
HHHHHHHHHHHCC CCEEEBCC CCCC CCCC CCEEEBCC CCCC CCCC CCCC CCCC
AA:
DIATLMKYTEAFDIDSFISFDLSLARGLDYVTGLIYEVVTSASAPPENASELKKKAKSAE
310 320 330 340 350 360

Conf: 4556665525874745778999857000257997747999619999999975315776
Pred:
CCCCCCEEEBCC CCHHHHHHHHCC CCCC CCCC CCEEEBCC HHHHHHHHHHHHCC C
C
AA: DASEFVGVGSLAAGR YDNL VNMFSEAS GKKS TQIPC VGISF GVERIFS LIKQRINS STT
370 380 390 400 410 420

Conf: 678876299997670566777899999999997798099996888799999999987
Pred:
CCCCCCEEEBCC CCHHHHHHHHHHHHHHHHHHHHCC CCEEEBCC CCHHHHHHHHH
HC
AA:
IKPTATQVFMVAFGGKDWGTGYLPERMKVTKQLWDAGIEAEVYKAKANPRKQFDDAE
KA
430 440 450 460 470 480

Conf: 989999987647767969999897645668625875999999999999872357899
Pred:
CCEEEBCC CCHHC CEEEBECC CCCC CCEEBCC HHHHHHHHHHHHHHHHHHHHHH
H
AA:
GCHAVILGKEEYLEGKLRVKRLGQEFADDDGELVSAADIVPIVQEKLSQIHEDGLNEVT
490 500 510 520 530 540

Conf: 997329
Pred: HHHHCC
AA: RLIKGL

Calculate PostScript, PDF and JPEG graphical output for this result
using:
<http://bioinf.cs.ucl.ac.uk/cei-bin/psipred/era/nph-view2.cei?id=082262e93452d467.psi>

PSIPRED

Graphical output for

Liam J. ...

Bioinformatics
University College

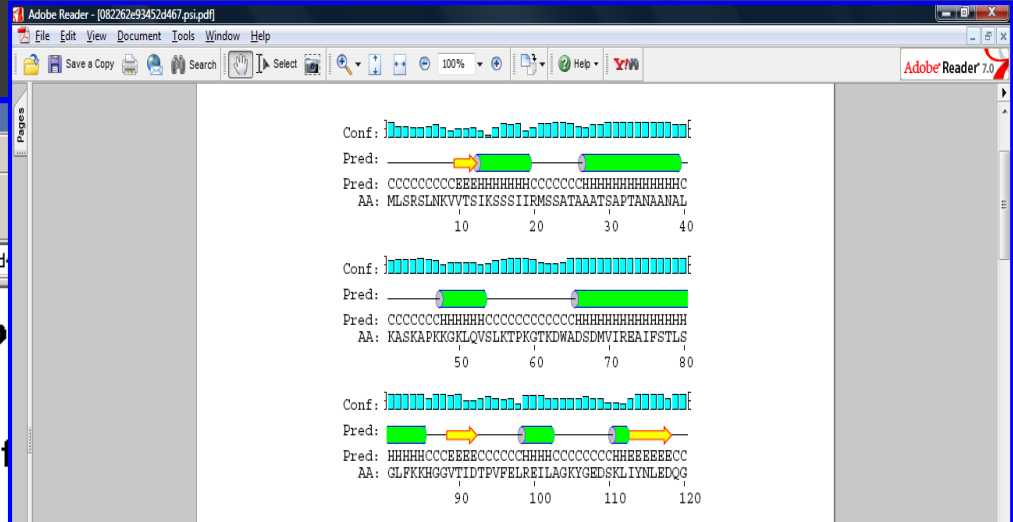
Click on t

Format

- RawScript File <http://b...>
- PDF File <http://b...>
- JPEG Page 1 <http://b...>
- JPEG Page 2 <http://b...>
- JPEG Page 3 <http://b...>

Jones DT (1999) Protein secondary structure prediction

McGuffin LJ, Bryson K, Jones DT (2000) The



Risultati predizione (versione grafica)

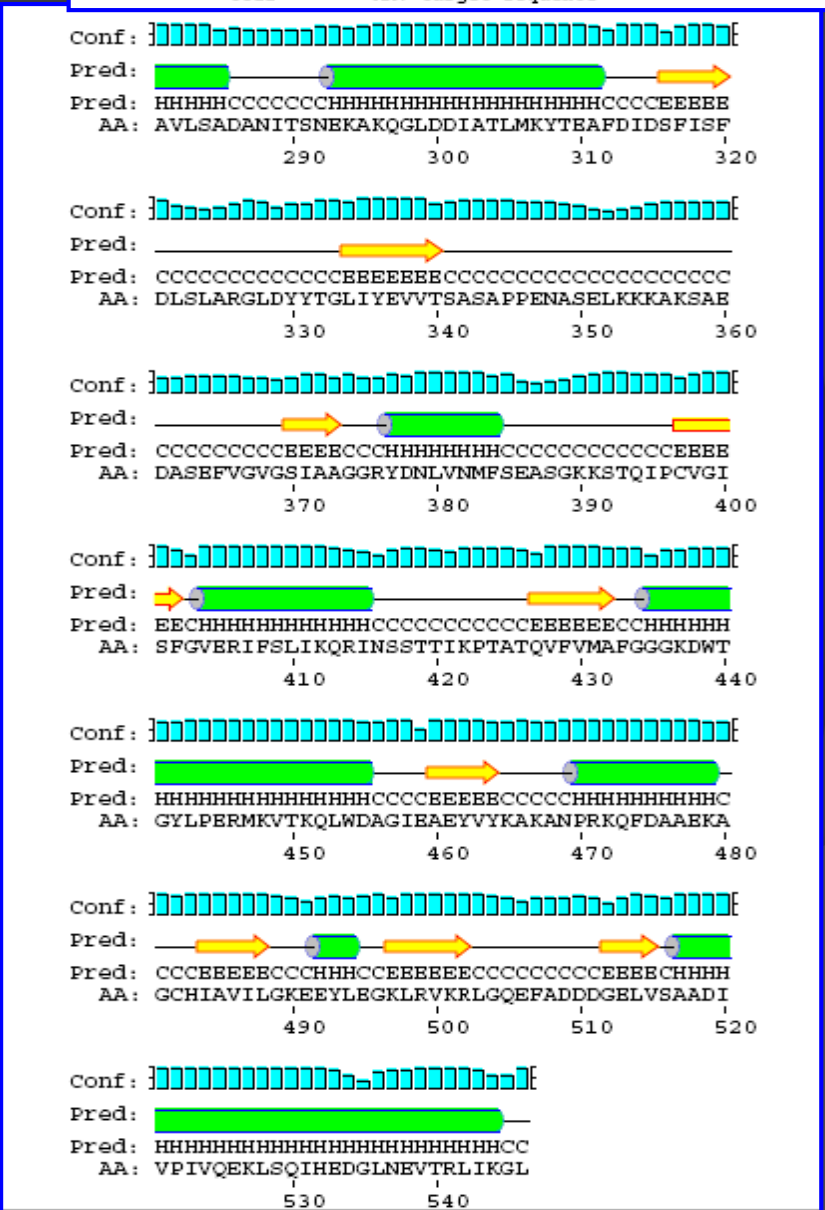
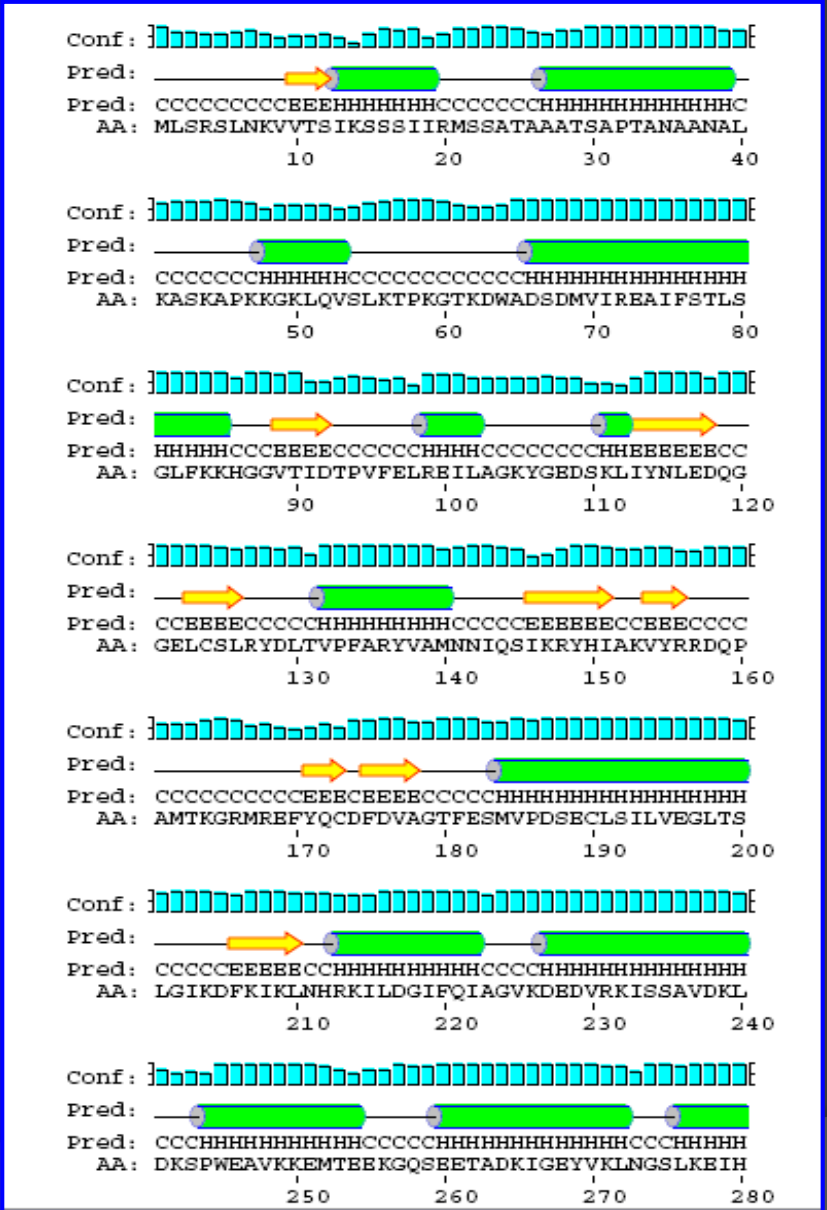
Legend:

- = helix
- = strand
- = coil

Conf: = confidence of prediction

Pred: = predicted secondary structure

AA: target sequence



Jpred 3

Incorporating Jnet

What is Jpred?

JPred is not a single computer program, but is a web system. In the original JPred, a range of different secondary structure prediction algorithms were run and the results combined.

JNet

How does Jpred work?

Development of the JNet algorithm showed that this was more accurate than JPred, so this is now the default and only algorithm that is run.

The server runs in two modes; single sequence and multiple sequence.



1. **Multiple sequence** If you already have a set of aligned sequences you may submit them as either **MSF format** or **BLC format**, and the predictions will run. Your alignment will be modified so that it does not contain gaps in the first sequence. The first sequence should therefore always be your target sequence.
2. **Single sequence** For single sequences a multiple alignment is constructed. It is created by the PSI-BLAST algorithm with 3 iterations. Redundant sequences are removed and gaps that have appeared in the query sequence are removed along with the aligned positions in the sequences. The prediction algorithms are then run.

4formati

Raw/FASTA

BLC

MSF

BATCH

6metodi di predizione

Six different prediction methods (DSC [King & Sternberg, 1996], PHD [Rost & Sander, 1993], NNSSP [Salamov & Solovyev, 1995], PREDATOR [Frishman & Argos, 1995], ZPRED [Zvelebil *et al.*, 1987] and MULPRED (Barton, 1988, unpublished) are then run, and the results from each method are combined into a simple file format.

OUTPUT

- ★ **Physico-chemical properties**
- ★ **Solvent accessibility**
- ★ **Prediction reliability**
- ★ **Conservation number values**



Publications: 2000

Cuff, J. A. and Barton, G. J. (2000), "Application of Multiple Sequence Alignment Profiles to Improve Protein Secondary Structure Prediction", *PROTEINS: Structure Function and Genetics*, 40:502-511. [[PubMed](#)] (EBI)

This paper describes the JNet secondary structure prediction algorithm. JNet gives over 76% accuracy in a comprehensive blind test and also predicts burial of residues and assigns confidence levels to predictions. The JNet software is available for

2007 → 81,5%

Advanced Jpred

fasta



Paste your sequence data here: [Help](#)

```
MTQENIPVTMSSSLSIILVILVSLRTALSELCNFPQDKQALLQIKKDLGNPTTLSSWLPPT  
DCCNRTWLGVLCDTDTQTYRVNNLDLSGHNLPKPYPIPSSLANLPYLNFLYIGGINNLVG  
RIPPAAIAKLTQLHYLYITHTNVSGAIPDFLSQIKTLVTLDFSYNALSGTLPPSISSLPNL  
GGITFDGNRISGAIPDSYGSFSKLFMTAMTISRNRLTGKIPPTFANLNLAFFVDSLNRMLEG  
LPSVLFSGDKNTKKIHLAKNSLAFDLGKVGLSKNLNLGLDLRNNRIYGTLPQGLTQLKPLQ  
SLNVSFNNLCCEIPQGGNLKRFDVSSYANNKCLCGSPLPST
```

Or upload a file: [Help](#)

Select type of input: [Help](#)


Single Sequence: Raw/Fasta Batch Mode
Multiple Alignment: MSF BLC Fasta

Skip searching PDB before prediction [Help](#)

Email address (optional): [Help](#)

Query name (optional): [Help](#)

Email



Match found in PDB

You might want to reconsider the accuracy and what you might gain from secondary structure prediction, if close sequence homologues exist in the structural database.

If you still want to carry out a Jpred prediction click

Hits found

PDB	Chain	Description	Blast E-value
1acx	A	ACTINOXANTHIN	3e-27
1j48	B	Apoprotein of C1027	4e-25
1j48	A	Apoprotein of C1027	4e-25
1hzl	A	C-1027 APOPROTEIN	4e-25
1hzk	A	C-1027 APOPROTEIN	4e-25
2g0l	A	NEOCARZINOSTATIN	8e-07
2g0k	A	Neocarzinostatin	8e-07
1o5p	A	Neocarzinostatin	8e-07
1noa	A	NEOCARZINOSTATIN	8e-07
1nco	B	HOLO-NEOCARZINOSTATIN	8e-07
1nco	A	HOLO-NEOCARZINOSTATIN	8e-07
1j5i	A	PROTEIN (Apo-Neocarzinostatin)	8e-07
1j5h	A	Apo-Neocarzinostatin	8e-07
2mcm	A	MACROMOMYCIN	2e-06
2cbt	B	NEOCARZINOSTATIN	8e-06
2cbt	A	NEOCARZINOSTATIN	8e-06

PDB last updated on: 2008-05-21

Alignment of PDB hits to your sequence

```
>1ack_A mol:protein length:100 ACTINOXANTHIN
      Length = 108

Score = 117 bits (293), Expect = 3e-27
Identities = 58/58 (100%), Positives = 58/58 (100%)

Query: 1 APAFSVSPASGASDQGSVSVVAAAGETTYIAQCAPVGGQDACPATATSFTTASGA 58
      APAFSVSPASG SDGQSVSVV AAAGETTYIAQCAPVGGQDACPATATSFTTASGA
Sbjct: 1 APAFSVSPASGASDQGSVSVVAAAGETTYIAQCAPVGGQDACPATATSFTTASGA 58

>1j48_B mol:protein length:110 Apoprotein of C1027
      Length = 110

Score = 110 bits (275), Expect = 4e-25

Query: 1 APAFSVSPASGASDQGSVSVV--AAAGETTYIAQCAPVGGQDACPATATSFTTASGA 58
      APAFSVSPASG SDGQSVSVV AAAGETTYIAQCAPVGGQDACPATATSFTTASGA
Sbjct: 1 APAFSVSPASGLSDGQSVSVVSGAAGETTYIAQCAPVGGQDACPATATSFTTASGA 60

>1hzi_A mol:protein length:110 C-1027 APOPROTEIN
      Length = 110

Score = 110 bits (275), Expect = 4e-25
Identities = 57/60 (95%), Positives = 57/60 (95%), Gaps = 2/60 (3%)

Query: 1 APAFSVSPASGASDQGSVSVV--AAAGETTYIAQCAPVGGQDACPATATSFTTASGA 58
      APAFSVSPASG SDGQSVSVV AAAGETTYIAQCAPVGGQDACPATATSFTTASGA
Sbjct: 1 APAFSVSPASGLSDGQSVSVVSGAAGETTYIAQCAPVGGQDACPATATSFTTASGA 60

>1hzk_A mol:protein length:110 C-1027 APOPROTEIN
      Length = 110

Score = 110 bits (275), Expect = 4e-25
Identities = 57/60 (95%), Positives = 57/60 (95%), Gaps = 2/60 (3%)

Query: 1 APAFSVSPASGASDQGSVSVV--AAAGETTYIAQCAPVGGQDACPATATSFTTASGA 58
      APAFSVSPASG SDGQSVSVV AAAGETTYIAQCAPVGGQDACPATATSFTTASGA
Sbjct: 1 APAFSVSPASGLSDGQSVSVVSGAAGETTYIAQCAPVGGQDACPATATSFTTASGA 60

>2g0l_A mol:protein length:122 NEOCARZINOSTATIN
      Length = 122

Score = 49.7 bits (117), Expect = 8e-07
Identities = 30/61 (49%), Positives = 38/61 (62%), Gaps = 3/61 (4%)

Query: 1 APAFSVSPASGASDQGSVSVVAA--AGETTYIAQCAPVGGQDACPATATSFTTASG 57
      AP +V+P+SG SDG V V+ A AG Y + QCA V G ACNPA +S T DA+G
Sbjct: 11 APTATVTPSSGLSDGTVVVKVAGAGLQAGTAYDVGGQCAWVDTOVLACNPADFSSVTADANG 70

Query: 58 A 58
      +
Sbjct: 71 S 71

>2g0k_A mol:protein length:122 Neocarzinostatin
      Length = 122
```

```
Score = 49.7 bits (117), Expect = 8e-07
Identities = 30/61 (49%), Positives = 38/61 (62%), Gaps = 3/61 (4%)

Query: 1 APAFSVSPASGASDQGSVSVVAA--AGETTYIAQCAPVGGQDACPATATSFTTASG 57
      AP +V+P+SG SDG V V+ A AG Y + QCA V G ACNPA +S T DA+G
Sbjct: 2 APTATVTPSSGLSDGTVVVKVAGAGLQAGTAYDVGGQCAWVDTOVLACNPADFSSVTADANG 61

Query: 58 A 58
      +
Sbjct: 62 S 62

>1nco_B mol:protein length:113 HOLO-NEOCARZINOSTATIN
      Length = 113

Score = 49.7 bits (117), Expect = 8e-07
Identities = 30/61 (49%), Positives = 38/61 (62%), Gaps = 3/61 (4%)

Query: 1 APAFSVSPASGASDQGSVSVVAA--AGETTYIAQCAPVGGQDACPATATSFTTASG 57
      AP +V+P+SG SDG V V+ A AG Y + QCA V G ACNPA +S T DA+G
Sbjct: 2 APTATVTPSSGLSDGTVVVKVAGAGLQAGTAYDVGGQCAWVDTOVLACNPADFSSVTADANG 61

Query: 58 A 58
      +
Sbjct: 62 S 62

>1j5i_A mol:protein length:122 PROTEIN (Apo-Neocarzinostatin)
      Length = 122

Score = 49.7 bits (117), Expect = 8e-07
Identities = 30/61 (49%), Positives = 38/61 (62%), Gaps = 3/61 (4%)

Query: 1 APAFSVSPASGASDQGSVSVVAA--AGETTYIAQCAPVGGQDACPATATSFTTASG 57
      AP +V+P+SG SDG V V+ A AG Y + QCA V G ACNPA +S T DA+G
Sbjct: 11 APTATVTPSSGLSDGTVVVKVAGAGLQAGTAYDVGGQCAWVDTOVLACNPADFSSVTADANG 70

Query: 58 A 58
      +
Sbjct: 71 S 71

>1j5h_A mol:protein length:122 Apo-Neocarzinostatin
      Length = 122

Score = 49.7 bits (117), Expect = 8e-07
Identities = 30/61 (49%), Positives = 38/61 (62%), Gaps = 3/61 (4%)

Query: 1 APAFSVSPASGASDQGSVSVVAA--AGETTYIAQCAPVGGQDACPATATSFTTASG 57
      AP +V+P+SG SDG V V+ A AG Y + QCA V G ACNPA +S T DA+G
Sbjct: 11 APTATVTPSSGLSDGTVVVKVAGAGLQAGTAYDVGGQCAWVDTOVLACNPADFSSVTADANG 70

Query: 58 A 58
```




[Top page](#)

[Protein](#) [Clefs](#) [Links](#)

Antibacterial protein

PDB id

1acx



Main view



Bottom view



Right view



Contents

Description

- [Header details](#)
- [Header records](#)
- [References](#)
- [PROCHECK](#)

Protein chain

[108 a.a.](#)

Tools

[Image Generation](#)

[AstexViewer™@MSD-EBI](#)

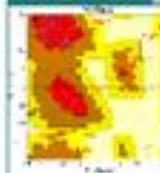
[Run PROCHECK](#)

[Clefs Calculation](#)

Quick links

- [RCSB](#)
- [MSD](#)
- [SRS](#)
- [EMDB](#)
- [JmolLib](#)
- [OCA](#)
- [CATH](#)
- [SCOP](#)
- [ESSP](#)
- [HSSP](#)
- [PQS](#)
- [ProSAT](#)
- [Whatcheck](#)

Procheck



Clefs



Surface



PDB id: **1acx**

Name: **Antibacterial protein**

Title: Actinoxanthin structure at the atomic level (russian)

Structure: Actinoxanthin. Chain: a. Engineered: yes

Source: Streptomyces globisporus

UniProt: [P01551](#) (ATXA_STROLY) [\[Pfam\]](#)

Seq: 143 a.a.

Struc: 108 a.a.*

Key: PfamA domain Secondary structure

* PDB and UniProt seqs differ at 5 residue positions (black crosses)

Resolution: 2.00Å

R-factor: not given

Authors: V.Z.Pletnev,A.P.Kuzin

Key ref: [\[PubMed id: \]](#) v.z.pletnev et al. (1982). Actinoxanthin Structure at the Atomic Level (Russian). *Bioorg.Khim.*, 8, 1637.

Date: 17-Dec-82

Release date: 09-Mar-83

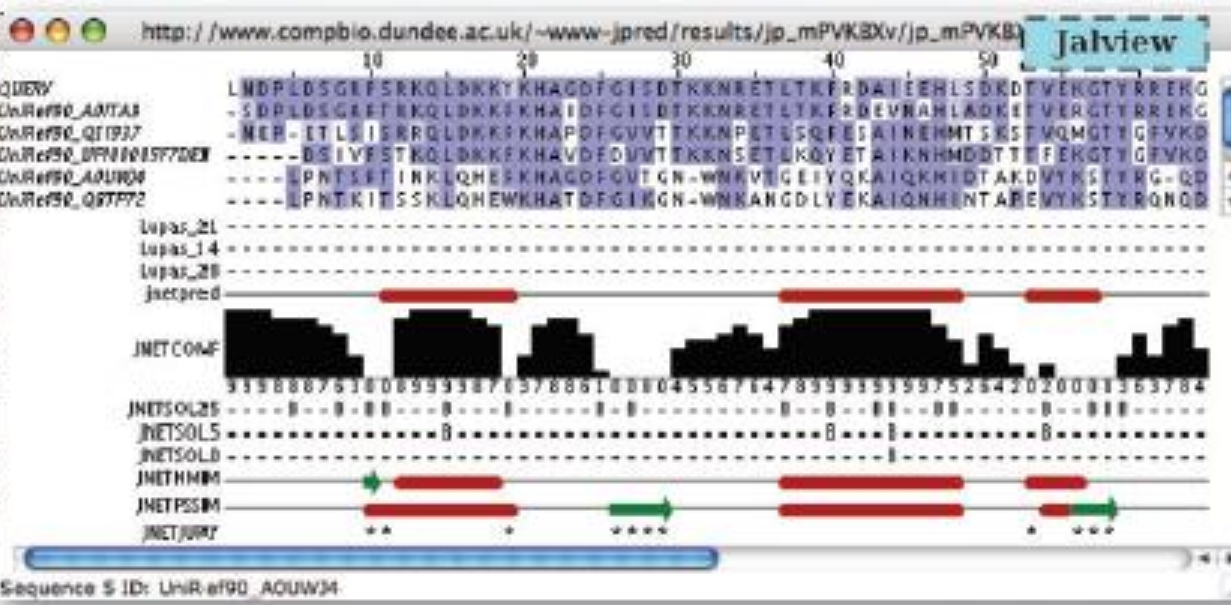
QUERY
 UniRef90_A01TA9
 UniRef90_Q11937
 UniRef90_UPI00005F7DE1
 UniRef90_A0UWJ4
 UniRef90_Q87P72
 UniRef90_A1DD43

OrigSeq
 1: LNDPLDSGRFSRKLQDKKTKF

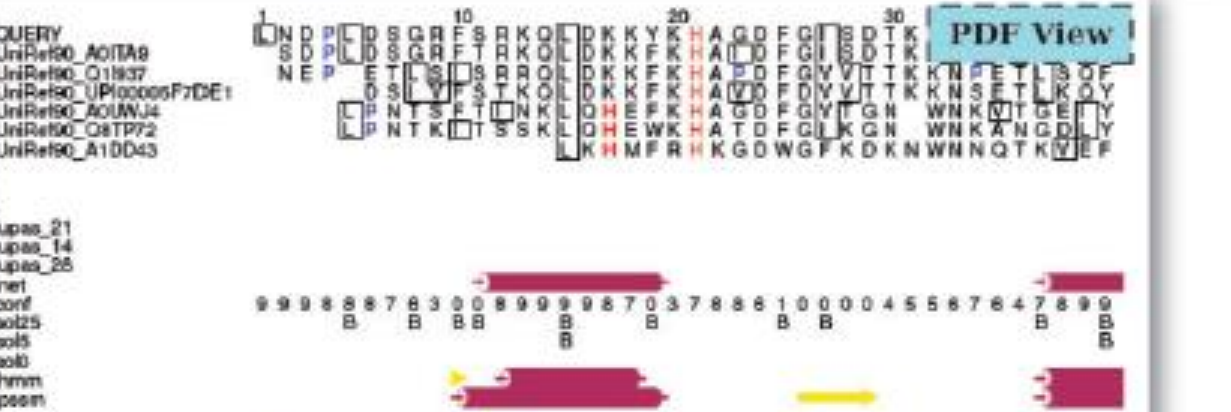
Jnet
 jhm
 jpsm

Lupas 14
 Lupas 21
 Lupas 28

Jnet_25
 Jnet_5
 Jnet_0
 Jnet_Rel
 1: 9998887630089998703



Notes
 Key:
 Colour code for alignment:
 Blue - Complete identity at a position
 Shades of red - The more red a position is, the more conservation of chemical properties
 Jnet - Final secondary structure prediction
 jalign - Jnet alignment prediction
 jhm - Jnet hhm profile prediction
 jpsm - Jnet PSIBLAST psam profile prediction
 Lupas - Lupas Coil prediction (window size)
 Note on coiled coil predictions
 - = less than
 c = between
 C = greater
 Jnet_25 - Jnet prediction of burial, 25%
 Jnet_5 - Jnet prediction of burial, 5%
 Jnet_0 - Jnet prediction of burial, 0%
 Jnet_Rel - Jnet reliability of prediction



conclusion

They represent extended (E), helical (H) and other (-) types of secondary structure respectively. In the solvent accessibility predictions they represent buried (B) and exposed (-) for each of the 0%, 5% and 25% solvent accessibility cut-offs.

What do the colours mean in the Postscript/PDF output from Jpred?

Character	Property
p	conserved polar
h	conserved hydrophobic
+	conserved positive charge
s	conserved small residues

Colour	Meaning
Pale blue	hydrophobic
Pale green	conserved polar
Small letters	small residues
Red	fully conserved
Blue text	Proline
Red text	Histidine
Boxed	Aliphatic (L, I or V)
Yellow	Cystine



Classification and Secondary Structure Prediction of Membrane Proteins

Mitaku Group
Department of Biotechnology
Tokyo University of Agriculture and Technology
[mailto: sosui@proteome.bio.tuat.ac.jp](mailto:sosui@proteome.bio.tuat.ac.jp)
[Naogoya/ TUAT]

[Contents]

1. SOSUI system
 - ◊ [SOSUI](#)
 - ◊ [SOSUI\(Batch\)](#)
 - ◊ [SOSUISignal](#)
 - ◊ [SOSUIDumbbell](#)
 - ◊ [SOSUIbreaker](#)
2. Database
 - ◊ Membrane Protein Database
 - ◊ [Database of Predicted Secretory Proteins in Prokaryota: SOSUIDBsecretory](#)

TMpred output for unknown

[ISREC-Server] Date: Tue Nov 28 12:46:27 Europe/Zurich 2006

Sequence: MNG...APA, length: 348

Prediction parameters: TM-helix length between 17 and 33

1.) Possible transmembrane helices

The sequence positions in brackets denominate the core region.

Only scores above 500 are considered significant.

Inside to outside helices : 6 found

	from		to	score	center
	37 (41)		59 (59)	2429	50
	74 (74)		99 (96)	1998	86
	117 (117)		139 (134)	1008	126
	153 (153)		175 (173)	2552	164
	205 (205)		222 (222)	2341	214
	253 (253)		274 (274)	3298	264

Outside to inside helices : 7 found

	from		to	score	center
	37 (39)		63 (56)	2497	47
	75 (78)		99 (95)	1617	86
	115 (118)		140 (140)	1441	127
	153 (157)		176 (173)	1808	165
	203 (203)		221 (221)	2960	213
	253 (253)		273 (271)	2684	263
	286 (290)		309 (309)	1023	298

TMpred output for unknown

[ISREC-Server] Date: Tue Nov 28 12:46:27 Europe/Zurich 2006

Sequence: MNG...APA, length: 348

Prediction parameters: TM-helix length between 17 and 33

1.) Possible transmembrane helices

The sequence positions in brackets denominate the core region.

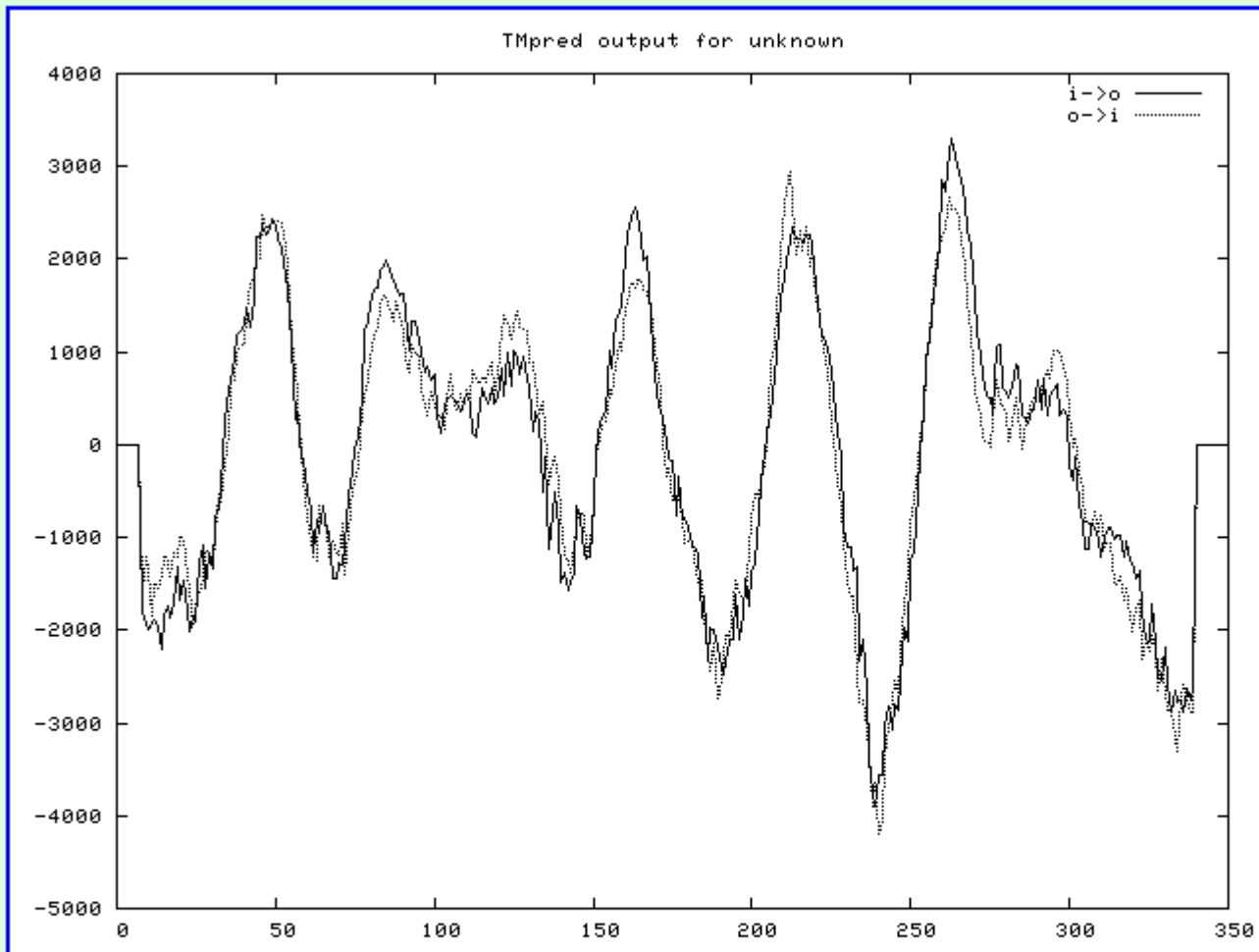
Only scores above 500 are considered significant.

Inside to outside helices : 6 found

	from		to	score	center
	37 (41)		59 (59)	2429	50
	74 (74)		99 (96)	1998	86
	117 (117)		139 (134)	1008	126
	153 (153)		175 (173)	2552	164
	205 (205)		222 (222)	2341	214
	253 (253)		274 (274)	3298	264

Outside to inside helices : 7 found

	from		to	score	center
	37 (39)		63 (56)	2497	47
	75 (78)		99 (95)	1617	86
	115 (118)		140 (140)	1441	127
	153 (157)		176 (173)	1808	165
	203 (203)		221 (221)	2960	213
	253 (253)		273 (271)	2684	263
	286 (290)		309 (309)	1023	298



You can get the prediction graphics shown above in one of the following formats:

- ◆ [GIF-format](#)
- ◆ [Postscript-format](#)
- ◆ [numerical format](#)

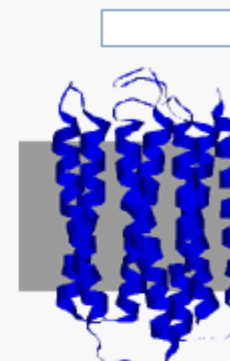
[CBS](#) >> [CBS Prediction Servers](#) >> [TMHMM](#)

TMHMM Server v. 2.0

Prediction of transmembrane helices in proteins

Update Nov. 29 2001: Minor change to the html output.

NOTE: You can submit many proteins at once in one fasta file. Please limit each submission to at most 4000 proteins. Please tick the 'One line per protein' option. Please leave time between each large submission.



Instructions

SUBMISSION

Submission of a local file in **FASTA** format (HTML 3.0 or higher)

OR by pasting sequence(s) in **FASTA** format:

```
MNGTEGPNFYVPPFSNKTGVVRSPPFEAPQYYLAEPWQFSMLAAYMFLILVIGFPINFLTLYVTVQHKKLRT
PLNYILLNLAVADLFMVFGGFTTTLTYTSLHGYFVFGPTGCNLEGFFATLGGEIALWSLVVLAIERVYVVC
KPMNSNFRFGENHAIMGVAFTWVMALACAAPPLVGWSRYIPQGMQCSCGALYFTLKPEINNESFVIYMFVV
HFSIPLIVIFFCYGQLVFTVKEAAAQQQESATTQKAEKEVTRMVIIMVIAFLICWLPYAGVAFYIFTHQG
SDFGPIFMTIPAFFAKSSSVYNPVIYIMMNKQFRNCMLTTLCCGKNPLGDDEASTTVSKTETSQVAPA
```

Output format:

- Extensive, with graphics
- Extensive, no graphics
- One line per protein

Other options:

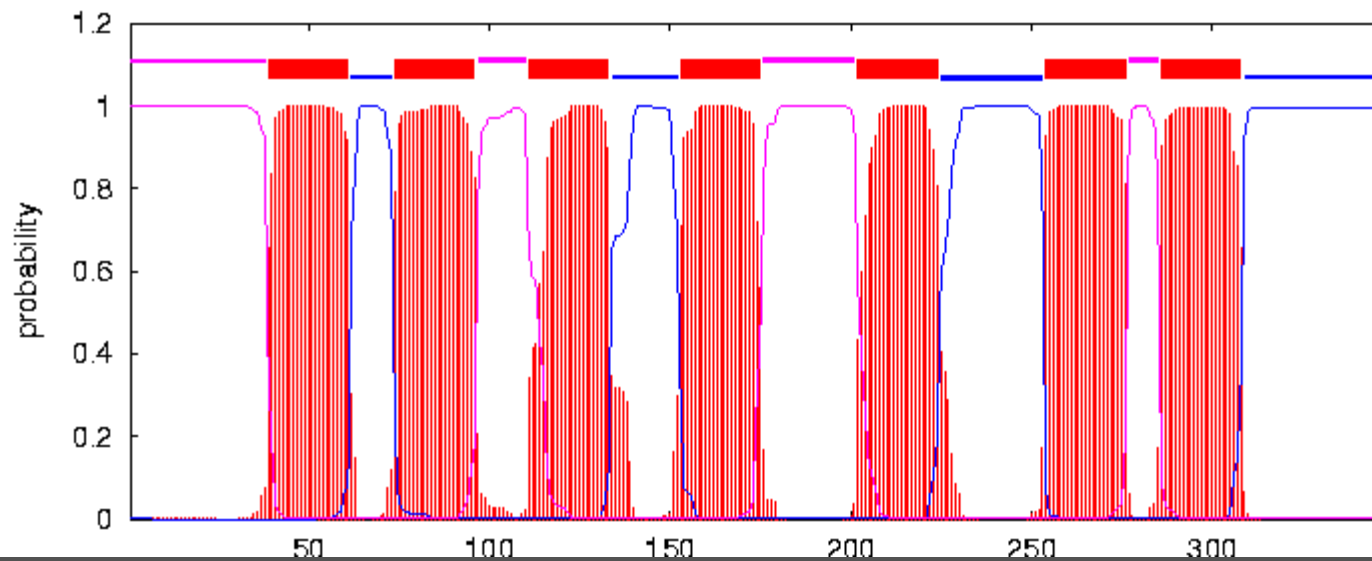
- Use old model (version 1)

```

# Sequence Length: 348
# Sequence Number of predicted TMHs: 7
# Sequence Exp number of AAs in TMHs: 157.8625
# Sequence Exp number, first 60 AAs: 21.49827
# Sequence Total prob of N-in: 0.00024
# Sequence POSSIBLE N-term signal sequence
Sequence TMHMM2.0 outside 1 38
Sequence TMHMM2.0 TMhelix 39 61
Sequence TMHMM2.0 inside 62 73
Sequence TMHMM2.0 TMhelix 74 96
Sequence TMHMM2.0 outside 97 110
Sequence TMHMM2.0 TMhelix 111 133
Sequence TMHMM2.0 inside 134 152
Sequence TMHMM2.0 TMhelix 153 175
Sequence TMHMM2.0 outside 176 201
Sequence TMHMM2.0 TMhelix 202 224
Sequence TMHMM2.0 inside 225 253
Sequence TMHMM2.0 TMhelix 254 276
Sequence TMHMM2.0 outside 277 285
Sequence TMHMM2.0 TMhelix 286 308
Sequence TMHMM2.0 inside 309 348

```

TMHMM posterior probabilities for Sequence





Search OPM

PDB ID or protein name

HOME

ABOUT OPM

DOWNLOAD OPM FILES

CONTACT US

Protein Classification

Types (3 types)

Classes (10 classes)

Superfamilies (168 superfamilies)

Families (238 families)

Species (237 species)

Localization (22 types)

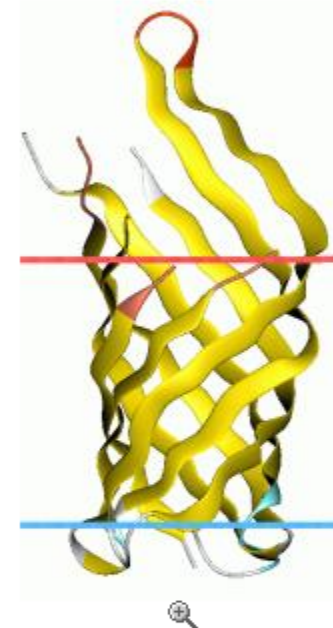
All proteins in OPM (629 entries)

Protein Links

[PDB Sum](#), [PDB](#), [MSD](#),
[SCOP](#), [Swiss-Prot](#), [Pfam](#),
[OCA](#), [MMDB](#), [FSSP](#), [HSSP](#)

1qjp » Outer membrane protein A (OMPA)

- **Type:** [1. Transmembrane](#) (2 classes)
- **Class:** [1.2. Beta-barrel transmembrane](#) (13 superfamilies)
- **Superfamily:** [1.2.01. OMPA-like \(n=8,S=10\)](#) (2 families) [1.B.6 \(TCDB\)](#)
- **Family:** [1.2.01.01. OMPA-like proteins](#) (4 proteins) [1.B.6 \(TCDB\)](#)
- **Species:** [Escherichia coli](#) (60 proteins)
- **Localization:** [Bacterial gram-negative outer membrane](#) (42 proteins)



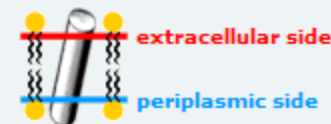
1qjp » Outer membrane protein A (OMPA)

Hydrophobic Thickness	24.9 ± 2.2 Å
Tilt Angle	8 ± 6°
ΔG_{transfer}	-36.7 kcal/mol
Links to 1qjp	PDB Sum , SCOP , MSD , OCA , MMDB , Dali , HSSP
Topology	subunit A (N-terminus periplasmic)
Resolution	1.65 Å
Related PDB Sum	1bxw , 1q90 , 2qed

3D view in [Chime](#), [Jmol](#) or [Webmol](#)

[Download Coordinates](#)

Topology in *Bacterial gram-negative outer membrane*





ProtScale

ProtScale [[Reference](#) / [Documentation](#)] allows you to compute and represent the profile produced by any amino acid scale on a selected protein.

An **amino acid scale** is defined by a numerical value assigned to each type of amino acid. The most frequently used scales are the hydrophobicity or hydrophilicity scales and the secondary structure conformational parameters scales, but many other scales exist which are based on different chemical and physical properties of the amino acids. This program provides 55 predefined scales entered from the literature.

Enter a [UniProtKB/Swiss-Prot](#) or [UniProtKB/TrEMBL](#) accession number (AC) (e.g. **P05130**) or a sequence identifier (ID) (e.g. **KPC1_DROME**):

Or you can paste your own sequence in the box below:

```
LSRINHFEDIQIIPKSSWSNHDASSGVSSACPYLGRSSFFRNVVWL IKKN
NNTNQEDLLVLWGVHHPNDAAEQTKLYQNPTTYSVGTSTLNQRLVPEIA
RMEFFWTILKPNDAINFESNGNFIAPYAYKIVKKGDSTIMKSELEYGNCI
INSSMPFHNHPLTIGCEPKYVKSRLVLA TGLRNTPQRERRRKRGLFG
QGMVDGWYGYHHSNEQGCYSADKESTQK AIDGVTNKVNS IINKMNTQFE
RRIENLNKKMEDGFLDVWTYNAELLVLMENERTLDFHDSNVKNLYDKVRL
NGCFEFYHKCDNECMESVKNGTYDYPQYSEEARLNREEISGVKLESMGTY
SSLALAIMVAGLSLWMCNGLQCRIC
```

Perché interessa il profilo idrofobico/idrofilico: per predire anse tra elementi di Struttura secondaria i residui esposti e quelli sepolti regioni che attraversano La membrana i siti antigenici

SEQUENCE LENGTH: 633

Using the scale [Hphob. / Kyte & Doolittle](#), the individual values for the 20 amino acids are:

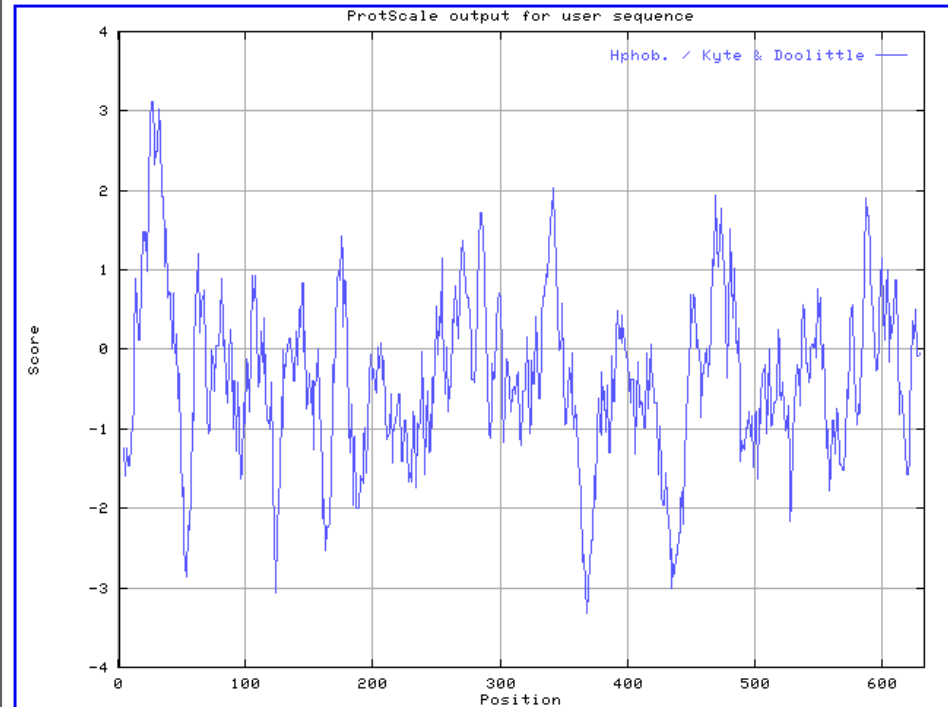
Ala: 1.800	Arg: -4.500	Asn: -3.500	Asp: -3.500	Cys: 2.500	Gln: -3.500
Glu: -3.500	Gly: -0.400	His: -3.200	Ile: 4.500	Leu: 3.800	Lys: -3.900
Met: 1.900	Phe: 2.800	Pro: -1.600	Ser: -0.800	Thr: -0.700	Trp: -0.900
Tyr: -1.300	Val: 4.200	Asx: -3.500	Glx: -3.500	Xaa: -0.490	

Weights for window positions 1,...,9, using **linear weight variation model**:

1	2	3	4	5	6	7	8	9
1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00
edge				center				edge

Profilo idropatia ProtScale

MIN: -3.322
MAX: 3.122





The ScanProsite tool [[Help](#) / [Commercial users](#)] allows to scan protein sequence(s) (either from [UniProt Knowledgebase \(Swiss-Prot/TrEMBL\)](#) or PDB or provided by the user) for the occurrence of patterns, profiles and rules (motifs) stored in the [PROSITE](#) database, or to search protein database(s) for hits by specific motif(s) [[Reference](#) / [Download ps_scan, the standalone version](#)]. The program [PRATT](#) can be used to generate your own patterns. You may either:

- ◆ Enter one or more PROSITE accession numbers and/or patterns [1 by line] to search the UniProt Knowledgebase (Swiss-Prot/TrEMBL) and/or PDB databases, **OR**
- ◆ Enter one or more sequences [raw, Swiss-Prot or fasta format] and/or UniProt Knowledgebase (Swiss-Prot/TrEMBL) accession numbers and/or PDB accession numbers [1 by line] to be scanned with all patterns, profiles, rules in PROSITE, **OR**
- ◆ Fill in both fields to find all occurrences of specified motifs in specified sequences.

Protein(s) to be scanned:

Enter one or more Swiss-Prot/TrEMBL accession number(s) [AC] (e.g. **P00747**) and/or sequence identifier(s) [ID] (e.g. **ENTK_HUMAN**), and/or PDB identifier, and/or paste **your own protein sequence(s)** in the box below:
(leave this box blank to scan PROSITE entrie(s) against selected protein databases)

PROSITE pattern(s)/profile(s) to scan for:

Enter one or more PROSITE accession number(s) (e.g. **PS50240**), and/or identifier(s) (e.g. **CHEB**), and/or type **your pattern(s)** in [PROSITE format](#) in the box below:
(leave this box blank to scan sequence(s) against the entire PROSITE database)

Se si vuole testare la propria sequenza contro tutti i motivi di PROSITE

Oppure un codice di Swiss-Prot o TREMBL/PDB contro tutti i motivi di PROSITE

LATO GIALLO

Se si vuole testare il proprio PROSITE pattern o Profilo sul database

LATO AZZURRO

Si possono provare più sequenze contemporaneamente 8 o 16, stessa cosa per i motivi

CLASSIFICAZIONE delle proteine

➤ Biochimica

- Globulari
- di Membrana
- Fibrose

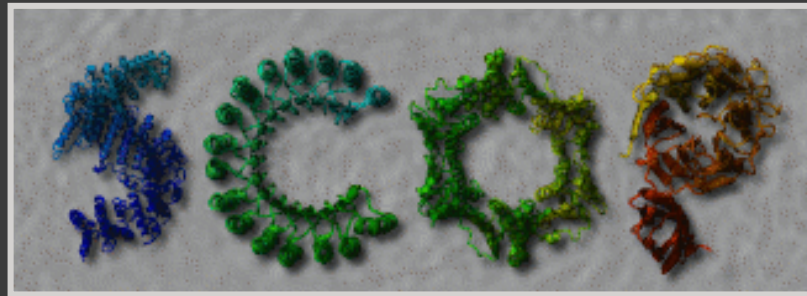
➤ Strutturale

- Substrutture:
BANCHE DATI
- SCOP
 - CATH

SCOP:

Structural Classification Of Proteins

(<http://scop.mrc-lmb.cam.ac.uk/scop>)



Classificazione e descrizione delle
relazioni strutturali ed evolutive tra tutte
le strutture proteiche conosciute

SCOP:

Structural Classification Of Proteins

(<http://scop.mrc-lmb.cam.ac.uk/scop>)

- Organizzato secondo dei livelli gerarchici:
 - Classe (strutturale)
 - Ripiegamento (strutturale)
 - Superfamiglia (evoluzionistico)
 - Famiglia (evoluzionistico)

SCO
P

- Unità di categorizzazione: **DOMINIO** (*domain*)

SCOP: FAMIGLIA


★ Insiemi di domini omologhi  FAMIGLIA
(*family*)




Proteine con:

- Identità di sequenza $\geq 30\%$
- Funzioni e strutture molto simili

SCOP: FAMIGLIA

★ Insiemi di domini omologhi  FAMIGLIA
(family)


Proteine con:

Origine
evolutiva
comune

- Identità di sequenza $\geq 30\%$
- Funzioni e strutture molto simili

SCOP: SUPERFAMIGLIA

- ★ Famiglie di proteine con:
 - Bassa identità di sequenza
 - Funzioni e strutture simili



SUPERFAMIGLIA
(*superfamily*)

SCOP: SUPERFAMIGLIA

- ★ Famiglie di proteine con:
- Bassa identità di sequenza
 - Funzioni e strutture simili

Origine
evolutiva
comune
probabile



SUPERFAMIGLIA
(*superfamily*)

SCOP: RIPIEGAMENTO

- ★ Due o più superfamiglie con una comune topologia di ripiegamento per un'ampia frazione della struttura



RIPIEGAMENTO
(*fold*)

SCOP: RIPIEGAMENTO

- ★ Due o più superfamiglie con una comune topologia di ripiegamento per un'ampia frazione della struttura

Nessuna
origine
evolutiva
comune



RIPIEGAMENTO
(*fold*)

SCOP: RIPIEGAMENTO

- ★ Due o più superfamiglie con una comune topologia di ripiegamento per un'ampia frazione della struttura

Stesso
ripiegamento:
Ragioni
chimico-fisiche



RIPIEGAMENTO
(*fold*)

SCOP: CLASSE

★ Ogni ripiegamento fa parte di una delle CLASSI (Class) generali.

➤ Classi SCOP:

- Tutto α
- Tutto β
- $\alpha+\beta$
- α/β
- *Multidomain protein*
- *Membrane and Cell Surface Protein*
- *Small protein*

SCOP: Structural Classification Of Proteins

Scop Classification Statistics

SCOP: Structural Classification of Proteins. 1.71 release
27599 PDB Entries (18 Jan 2005). 75930 Domains. 1 Literature Reference
(excluding nucleic acids and theoretical models)

Class	Number of folds	Number of superfamilies	Number of families
All alpha proteins	226	392	645
All beta proteins	149	300	594
Alpha and beta proteins (a/b)	134	221	661
Alpha and beta proteins (a+b)	286	424	753
Multi-domain proteins	48	48	64
Membrane and cell surface proteins	49	90	101
Small proteins	79	114	186
Total	971	1589	3004

SCOP

ESEMPIO:

Classificazione
dell'emoglobina
umana (catena α)

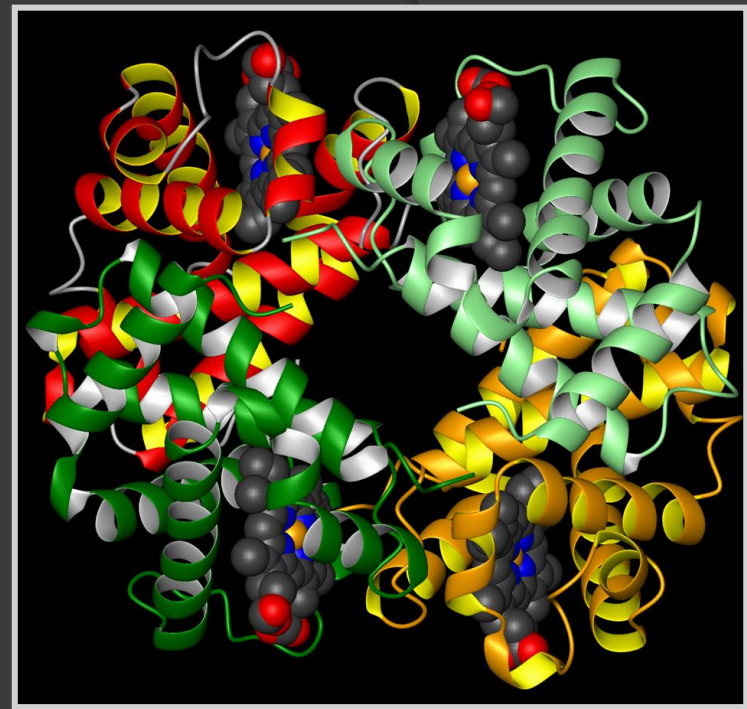
Protein: Hemoglobin, alpha-chain
from Human (*Homo sapiens*)

Root: scop

- **Class:** All alpha proteins
- **Fold:** Globin-like core: 6 helices; folded leaf, partly opened
- **Superfamily:** Globin-like
- **Family:** Globins Heme-binding protein

Protein: Hemoglobin, alpha-chain

Species: Human (*Homo sapiens*)



CATH:

Class Architecture Topology and Homologous

superfamily

(<http://www.biochem.ucl.ac.uk/bsm/cath/>)

- **Database di Classificazione** delle strutture proteiche presenti nella Protein Data Bank (PDB)

CATH:

Class Architecture Topology and Homologous
superfamily

(<http://www.biochem.ucl.ac.uk/bsm/cath/>)

- Organizzato secondo dei livelli gerarchici:
 - Classe (class)
 - Architettura (architecture)
 - Topologia (topology)
 - Superfamiglia Omologa (homologous superfamily)

CATH: CLASSE

- ★ Determinata dalla **composizione della struttura secondaria e dell'impaccamento**



CLASSE (*C-level*)

CATH: CLASSI

- *mainly* α
- *mainly* β
- $\alpha\beta$: α/β e $\alpha+\beta$
- Basso contenuto di struttura secondaria

CATH:

ARCHITETTURA

(A-level)

- ★ Descrive la forma complessiva del dominio
- ★ E' determinata dall'orientamento delle strutture secondarie nello spazio 3D, senza tener conto delle loro connessioni

Esempio: struttura a botte

CATH: TOPOLOGIA
o famiglia di ripiegamento
(*T-level*)

- ★ Descrive il ripiegamento proteico tenendo conto dell'orientazione delle strutture secondarie e delle connessioni tra esse

CATH: SUPERFAMIGLIA OMOLOGA (*H-level*)

- ★ Proteine con relazioni evolucionistiche (antenato comune) = omologhe:
 - Similarità strutturale e/o funzionale
 - Alta identità di sequenza

Motif Scan

user:
anonymous

Protein Sequence Input
Enter a protein sequence in RAW or FASTA or Swiss-Prot format or a db:AC or db:ID identifier

```
MERTVLLLATVSLVKSDQICIGYHANNSTEQVDTIMEKNVTVTHAQDILEI  
NGVKPLILRDCSVAGWLLGNPMCDEFINVPEWSYIVEKASPANDLCYPGNI  
LSRINHFEDIQIIPKSSWSNHDASSGVSSACPYLGRSSFFRNVVWLIKKN:  
NNTNQEDLLVLWGVHHPNDAAEQTKLYQNPTTYISVGTSTLNQRLVPEIA'  
RMEFFWTILKPNDAINFESNGNFIAPFYAYKIVKKG DSTIMKSELEYGNCI  
INSSMPFHNIHPLTIGECPKYVKS NRLVLATGLRNTPQRE RRRRKRGLFG  
QGMVDGWYGYHHSNEQGSCYSADKESTQKAIDGVTNKVNSIINKMNTQFE.  
RR IENLNKKMEDGFLDVWVTYNAELLVLMENERTLDFHDSNVK NLYDKVRLC
```

Motif scanning means finding all known motifs that occur in a sequence. This form lets you paste a protein sequence, select the collections of motifs to scan for, and launch the search. Some general [documentation](#) is available about the Prosite and Pfam collections of motifs. Another [document](#) deals with the interpretation of the match scores. You should consult the home pages of [Prosite](#) on ExPASy, [Pfam](#) and [InterPro](#) for additional information.

Warning: The scan might take a few minutes, thus if your proteins of interest are already in the sequence databases (see [list](#)), the [Query by Protein](#) form is much faster, and the [Protein Hub](#) provides a collection of tools that you might find useful.

Parameters

Database of motifs
([db description](#))

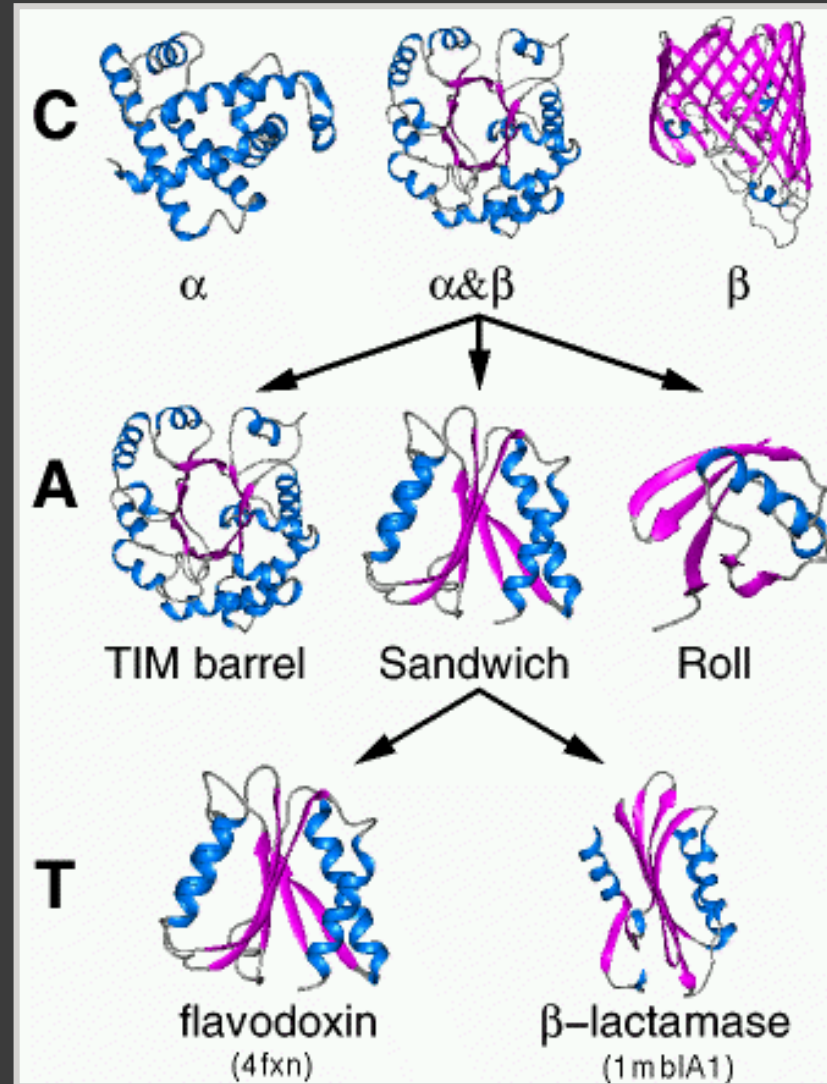
- PROSITE patterns
- PROSITE patterns (frequent match producers)
- PROSITE profiles
- Profile (more profiles)
- Na-channel profiles
- HAMAP profiles
- Pfam HMMs (local models)
- Pfam HMMs (global models)

CATH: FAMIGLIA DI SEQUENZA (*S-level*)

- ★ Proteine di uno stesso *H-level*:
 - identità di sequenza $\geq 35\%$
 - similarità strutturale e/o funzionale

Level	Name	Sequence Identity Overlap
S	35%	80%
O	60%	80%
L	95%	80%
I	100%	80%

CATH: Esempio di livelli gerarchici



CATH:

Class Architecture Topology and Homologous superfamily

CATH v3.1.0

Version	3.1.0
Date	19-01-2007
Number of Domains	93885
Number of Chains	63453
Number of PDBs	30028

C	A	T	H	S	O	L	I	D
Mainly Alpha	5	305	652	1850	2329	3001	5587	19729
Mainly Beta	20	191	415	1860	2531	3846	6503	25537
Alpha Beta	14	496	922	3922	5303	6659	12998	47193
Few Secondary Structures	1	92	102	162	200	275	403	1426
Total	40	1084	2091	7794	10363	13781	25491	93885

CATH

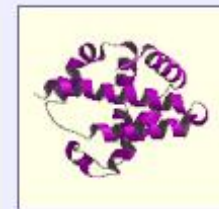
ESEMPIO:

Classificazione dell'emoglobina umana (catena α)

CATH Domain 1spgA0

Classification

C Class	1
Mainly Alpha	
A Architecture	1.10
Orthogonal Bundle	
T Topology	1.10.490
Globin-like	
H Homologous Superfamily	1.10.490.10
Globins	
S Sequence Family (S35)	1.10.490.10.3
Globins	
N Non-identical (S95)	1.10.490.10.3.7
Globins	
I Identical (S100)	1.10.490.10.3.7.1
Globins	



1spgA0



Motif Scan Results

user: anonymous

[log in](#)

Query Protein temporarily stored [here](#).

Database of motifs PROSITE patterns, PROSITE patterns (frequent match producers), PROSITE profiles, Pfam HMMs (local models), Pfam HMMs (global models).

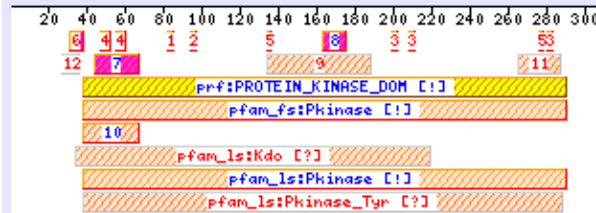
Reference Falquet L, Pagni M, Bucher P, Hulo N, Sigrist CJ, Hofmann K & Bairoch A. (2002) The PROSITE database, its status in 2002. *Nucleic Acids Res.* **30**:235-238

searching PROSITE patterns
searching PROSITE patterns (frequent match producers)
searching PROSITE profiles
searching Pfam HMMs (local models)
searching Pfam HMMs (global models)
postprocessing

Summary

Original output [pat](#), [freq_pat](#), [prf](#), [pfam fs](#), [pfam ls](#).

Matches map
(features from query are above the ruler, matches of the motif scan are below the ruler)



Legends: 1, freq_pat:ASN_GLYCOSYLATION [?]; 2, freq_pat:CAMP_PHOSPHO_SITE [?]; 3, freq_pat:CK2_PHOSPHO_SITE [?]; 4, freq_pat:MYRISTYL [?]; 5, freq_pat:PKC_PHOSPHO_SITE [?]; 6, freq_pat:TYR_PHOSPHO_SITE [?]; 7, pat:PROTEIN_KINASE_ATP [!]; 8, pat:PROTEIN_KINASE_ST [!]; 9, pfam_fs:APH [?]; 10, pfam_fs:Pkinase_Tyr [!]; 11, pfam_fs:Pkinase_Tyr [?]; 12, pfam_ls:Involucrin [?].

FT	MYHIT	82	85	freq_pat:ASN_GLYCOSYLATION [?]
FT	MYHIT	94	97	freq_pat:CAMP_PHOSPHO_SITE [?]
FT	MYHIT	199	202	freq_pat:CK2_PHOSPHO_SITE [?]
FT	MYHIT	208	211	freq_pat:CK2_PHOSPHO_SITE [?]
FT	MYHIT	280	283	freq_pat:CK2_PHOSPHO_SITE [?]
FT	MYHIT	47	52	freq_pat:MYRISTYL [?]
FT	MYHIT	55	60	freq_pat:MYRISTYL [?]
FT	MYHIT	134	136	freq_pat:PKC_PHOSPHO_SITE [?]
FT	MYHIT	256	258	freq_pat:PKC_PHOSPHO_SITE [?]

Match Status Code

! A strong match: it is very unlikely that this match is a false positive.

R Rescued match: despite the low score, it is considered to be a strong match. This concerns primarily domains known to be repeated and that are unlikely to occur as a single copy in a protein.

? Questionable or weak match: determining the true or false negative status of this match requires additional biological evidences.

!! Strong match for a family-specific motif: it is very unlikely that this match is a false positive, in addition it is very likely that this match belongs to the targetted sub-family.

?! Accepted match for a family-specific motif: it is very unlikely that this match is a false positive for the motif, but determining its family assignment requires additional biological evidences.

?? Questionable or weak match for a family-specific motif: determining the true or false negative status of this match requires additional biological evidences.

NA Not Available. The above interpretation rules make no sense (e.g., for a low-complexity region).

Prosite pattern e profile permette di caratterizzare la funzione di una proteina non caratterizzata database di siti e pattern biologicamente significativi.

Pfam è un database in cui vengono raccolte famiglie e motivi di proteine

PfamHMM

InterPro database di famiglie di proteine

HAMAP profile raccolta di famiglie di proteine ortologhe microbiche

TIGRfam raccolta di famiglie di proteine



NiceSite View of PROSITE: PS00128

General information about the entry

Entry name	LACTALBUMIN_LYSOZYME
Accession number	PS00128
Entry type	PATTERN
Date	APR-1990 (CREATED); NOV-1997 (DATA UPDATE); OCT-2006 (INFO UPDATE).
PROSITE documentation	PDOC00119

Name and characterization of the entry

Description	Alpha-lactalbumin / lysozyme C signature.
Pattern	C-x(3)-C-x(2)-[LMF]-x(3)-[DEN]-[LI]-x(5)-C.

Numerical results

- ◆ UniProtKB/Swiss-Prot release number: **51.1**, total number of sequence entries in that release: **241365**.
- ◆ Total number of hits in UniProtKB/Swiss-Prot: **129 hits in 129 different sequences**
- ◆ Number of hits on proteins that are known to belong to the set under consideration: **128 hits in 128 different sequences**
- ◆ Number of hits on proteins that could potentially belong to the set under consideration: **0 hits in 0 different sequences**
- ◆ Number of false hits (on unrelated proteins): **1 hits in 1 different sequences**
- ◆ Number of known missed hits: **2**
- ◆ Number of partial sequences which belong to the set under consideration, but which are not hit by the pattern or profile because they are partial (fragment) sequences: **7**
- ◆ Precision (true hits / (true hits + false positives)): **99.22 %**
- ◆ Recall (true hits / (true hits + false negatives)): **98.46 %**

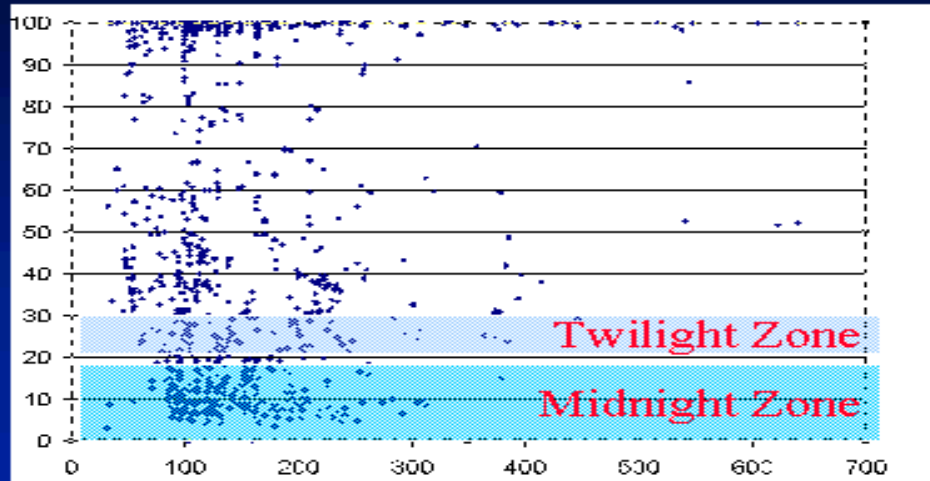
Comments

- ◆ Taxonomic range: **Eukaryotes**



- Segregazione delle catene laterali idrofobiche:

1. Riduzione del numero delle conformazioni possibili per ingombro sterico;
2. I gruppi $-NH$ e $-CO$ dei residui segregati, trovandosi vicino, formano legami ad H \rightarrow Formazione di elementi di struttura secondaria



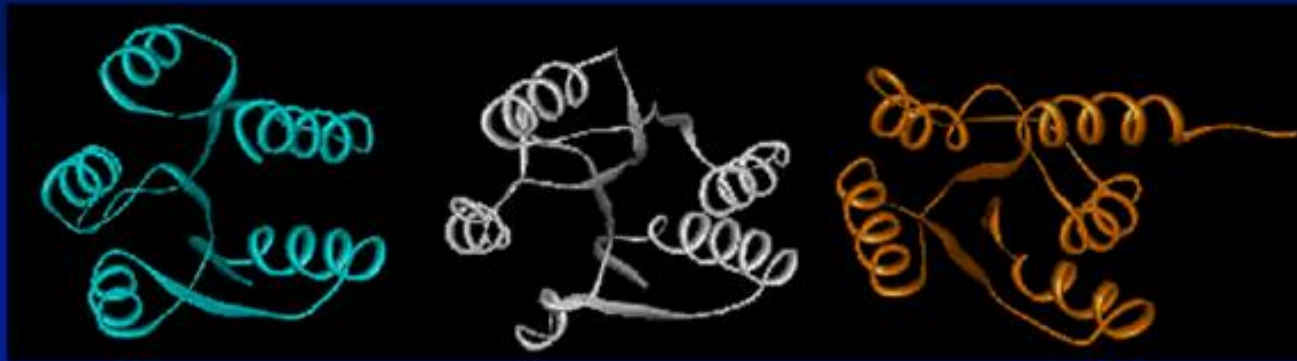
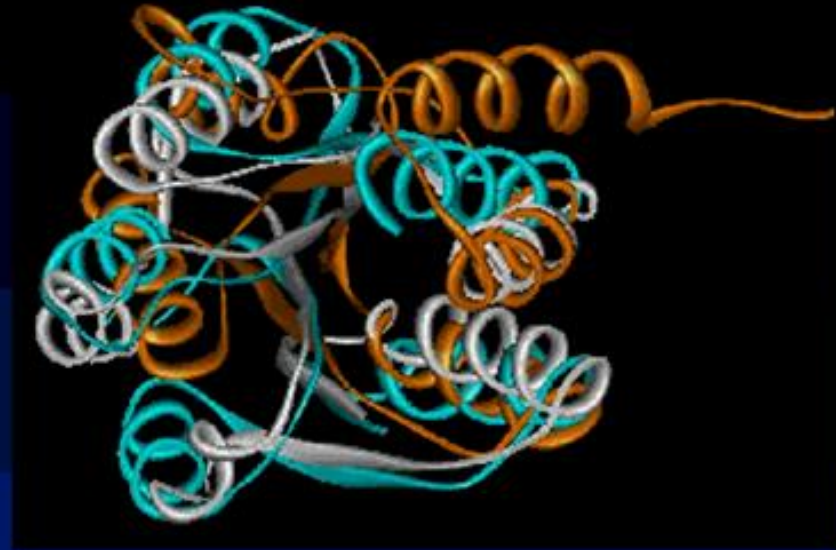
Zona
crepuscolare

Sono state raccolte in modo casuale una certa quantità di sequenze presenti nel PDB e le loro sequenze sono state confrontate tra di loro.

Zona Certa Sequenze con similarità superiore al 40% per il 95% dei casi sono omologhe.

Zona ambigua In cui la similarità è tra il 20% - 30% solo il 10% dei casi le proteine sono omologhe

C'è poi una zona in cui la percentuale di similarità è inferiore al 20% ed in cui non si trovano pochissime sequenze omologhe



1ymv

1fla

1pdo

Alpha/beta proteins characterized as different superfamilies

1PIV:1
Viral Capsid Protein



1HMP:A
Glycosyltransferase



80 Residue Stretch (Yellow) with Over 40% Sequence Identity

Dall'analisi del PDB

Ci sono circa 1000 famiglie di proteine composte da membri di proteine che mostrano una certa similarità di sequenza

Una nuova sequenza mostrerà quasi sicuramente similarità di sequenza con altre e componenti strutturali simili.

C'è un limitato numero di ripiegamenti, nuove sequenze di solito mostrano ripiegamenti simili con strutture già note

Sequenze completamente diverse possono ripiegarsi in strutture simili

Ci sono 3 metodi di approccio per predire la struttura proteica

Metodi Ab-initio

Modeling comparativo

Fold Recognition

La predizione è basata su un calcolo computazionale proveniente dalla posizione di ogni atomo nello spazio e le sue relazioni chimico-fisiche con altri atomi

Teoreticamente possibile

Praticamente poco possibile



SWISS-MODEL

An Automated Comparative Protein Modelling Server

SIB - Biozentrum Basel site provided by:



SWISS-MODEL is a fully automated protein structure homology-modeling server, accessible via the [ExPASy](#) web server, or from the program [DeepView](#) (Swiss Pdb-Viewer). The purpose of this server is to make Protein Modelling accessible to all biochemists and molecular biologists World Wide.

The present version of the server is 3.5 and is under constant improvement and debugging. In order to help us refine the sequence analysis and modelling algorithms, please [report](#) of possible bugs and problems with the modelling procedure.

SWISS-MODEL was initiated in 1993 by Manuel Peitsch, and is now being further developed within the [SIB - Swiss Institute of Bioinformatics](#) in collaboration between Torsten Schwede at the [Structural Bioinformatics Group](#), Biozentrum (University of Basel) and Nicolas Guex at [GlaxoSmithKline](#).

The computational resources for the SWISS-MODEL server are provided in collaboration by the Biozentrum (University Basel) and the [Advanced Biomedical Computing Center](#) (NCI Frederick, USA).

Disclaimer

The result of any modelling procedure is NON-EXPERIMENTAL and MUST be considered with care. This is especially true since there is no human intervention during model building. Carefully read the header section of the files to know what templates and alignments were used during the model building process.



SWISS MODEL

E' un server automatizzato e
comparativo per il modelling proteico
nato nel 1993 da Manuel Peitsch

Reference

- Arnold K., Bordoli L., Kopp J., and Schwede T. (2006). The SWISS-MODEL Workspace: A web-based environment for protein structure homology modelling. *Bioinformatics*, 22, 195-201.
- Kopp J. and Schwede T. (2004) The SWISS-MODEL Repository of annotated three-dimensional protein structure homology models *Nucleic Acids Research* 32, D230-D234.
- Schwede T, Kopp J, Guex N, and Peitsch MC (2003) SWISS-MODEL: an automated protein homology-modeling server. *Nucleic Acids Research* 31: 3381-3385.
- Guex, N. and Peitsch, M. C. (1997) SWISS-MODEL and the Swiss-PdbViewer: An environment for comparative protein modelling. *Electrophoresis* 18: 2714-2723.
- Peitsch, M. C. (1995) Protein modeling by E-mail *Bio/Technology* 13: 658-660.

ATTENZIONE!!

The result of any modelling procedure is **NON-EXPERIMENTAL** and **MUST** be considered with care. This is especially true since there is in human intervention during model building. Carefully read the header section of the files to know what templates and alignments were used during the model building process

Step 1

Search for suitable
through BLASTp2.

Database : ExNRL-3D

template

Step 2

Check sequence identity with
target

Programma : SIM

Step 3

Generate models with ProModII

Programma: ProModII

Database: ExPDB

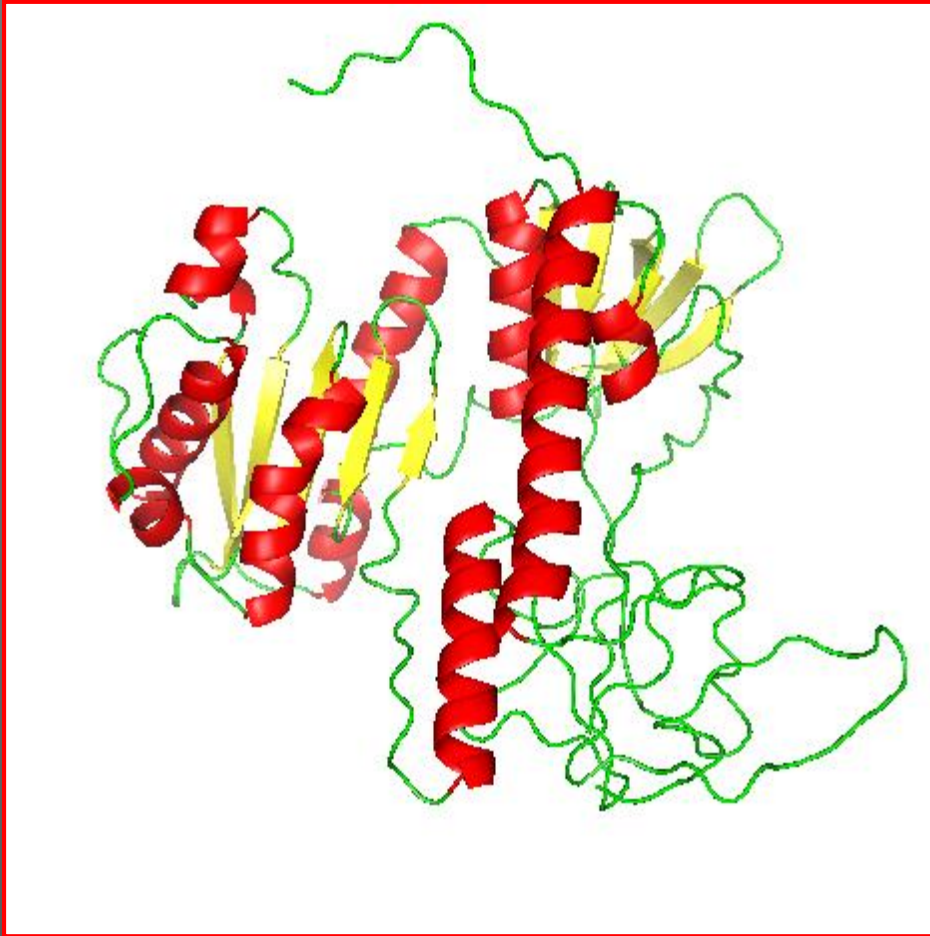
ProModII

**At least one known 3D-structure of a related protein.
Good quality sequence alignment; the reliability of the model is determined by the degree of sequence identity**

- 1. Superposition of related 3D-structure**
- 2. Generation of a multiple alignment with the sequence to be modelled.**
- 3. Generation of a framework for the new sequence.**
- 4. Rebuild lacking loops.**
- 5. Complete and correct backbone.**
- 6. Correct and rebuild side chains.**
- 7. Verify model structure quality and check packing.**
- 8. Refine structure by energy minimisation and Molecular Dynamics.**

Step 4

Energy minimisation with
Gromos 96



Model info:

modelled residue range:

18 to 404

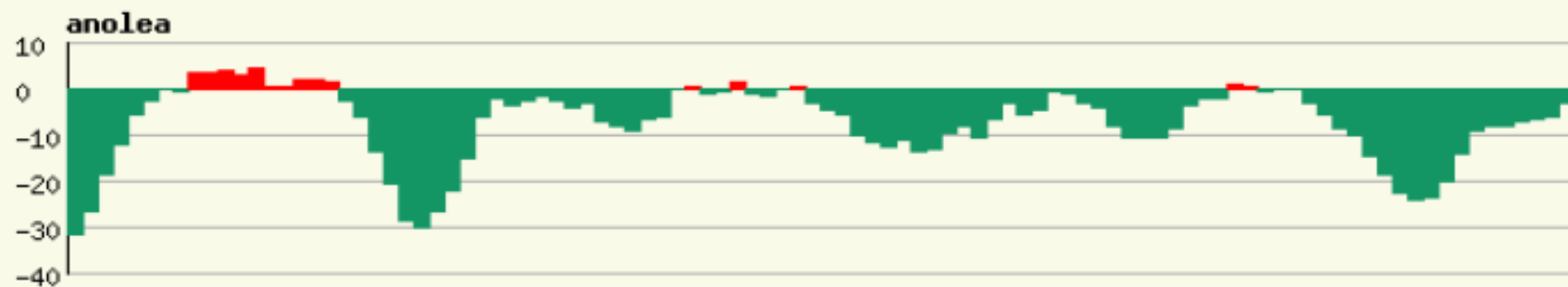
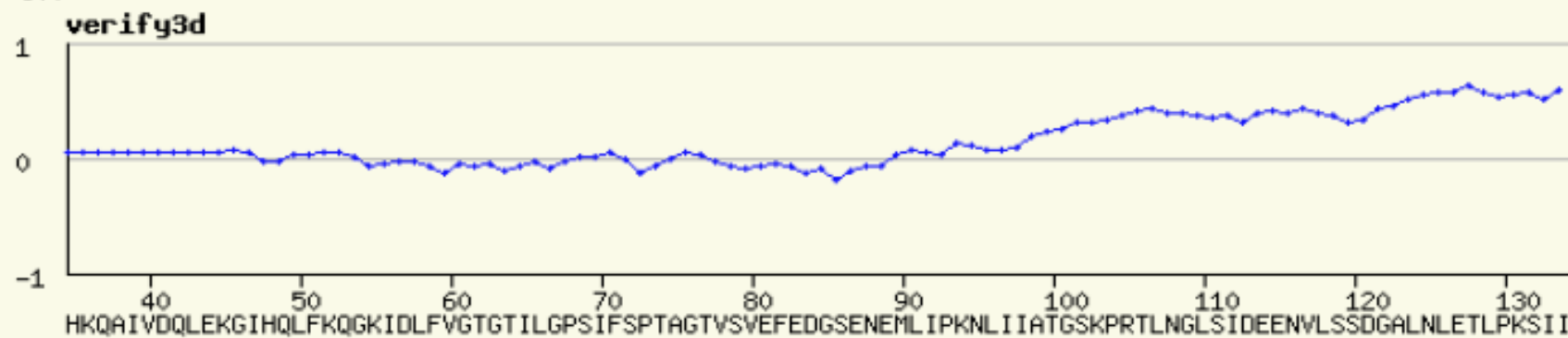
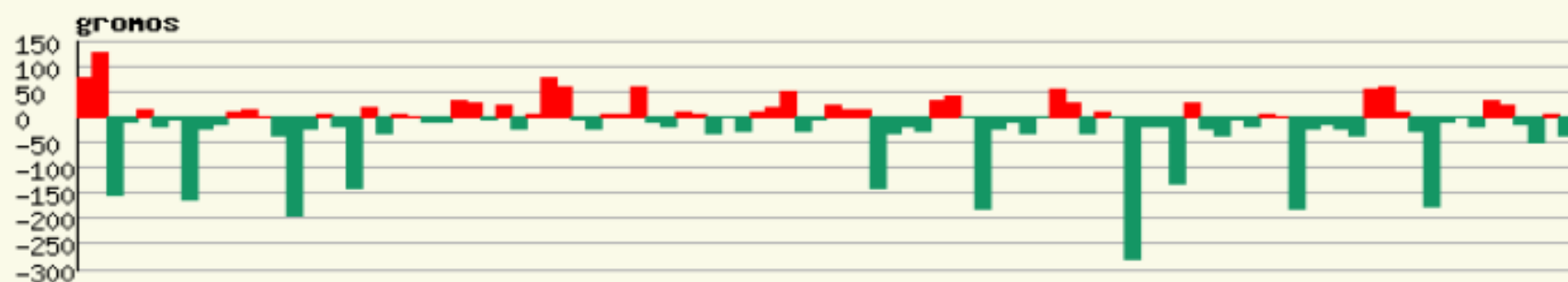
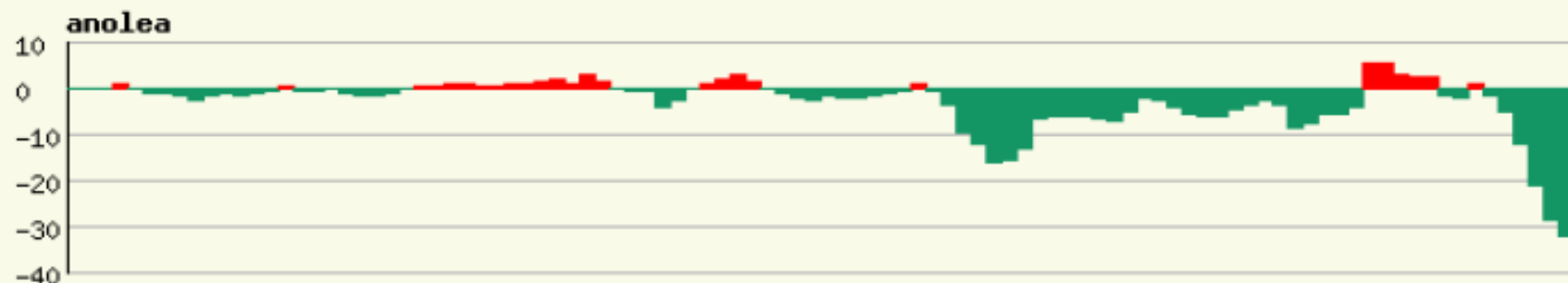
based on template:

2yquB (1.70 Å)

Sequence Identity [%]: 41

Evalue:

5.02e-65



Pôle BioInformatique Lyonnais



Geno3D

- Server utilizzabile per generare modelli proteici 3D non sperimentali, tramite modelling per omologia;
 - E' possibile utilizzare modelli con basso grado di identità (fino al 20%);
 - Genera modelli fino a 500 amminoacidi;
 - E' possibile fare allineamenti multipli (fino a 3 contemporaneamente).
-
- L'operatore fornisce la query che viene confrontata e allineata, tramite PSI-BLAST, con tutte le voci presenti nella banca dati di PDB;
 - L'utente seleziona le voci che vuole siano utilizzate nella modellizzazione molecolare.
 - Il server calcola la predizione di struttura della query basandosi su ciascun template e ne calcola la percentuale di accordo.
 - In caso di selezione di più template è visualizzato lo scarto quadratico medio tra il carbonio α e i vari template.
 - Le restrizioni degli angoli e le distanze tra gli atomi vengono calcolati a partire dall'allineamento con i template 3D, per poi essere applicati sulla sequenza query. Per i gap vengono utilizzati calcoli statistici.
 - Il software utilizzato per calcolare le restrizioni è CNS.
-
- Alla fine del processo di modellizzazione molecolare si riceve una e-mail in cui è fornito un indirizzo internet dove sono presenti i risultati (disponibili per 7 gg). E' possibile utilizzare questi risultati on line o scaricarli sul proprio pc in formato .archive.tar.gz
 - In output vengono rilasciati i modelli 3D che soddisfano le restrizioni nel migliore dei modi.
 - Vi possono essere porzioni che non vengono risolte basandosi sui template; queste possono essere particolari domini che vengono ricostruiti usando template di proteine omologhe.



GENO3D Release 2 : AUTOMATIC MODELING OF PROTEINS THREE-DIMENSIONAL STRUCTURE

[Abstract] [\[GENO3D help\]](#) [Original server]

Database :

Sequence 1

Paste protein 1 sequence below : [help](#)

```
MTMDKSELVQKAKLAEQAERYDDMAAAMKAVTEQGHELSNEERNLLSVAYKNVVGARRSS
WRVISSIEQKTERNEKKQMGKEYREKIEAELQDICNDVLELLDKYLIPNATQPESKVFY
LKMKGDYFRYLSEVASGDNKQTTVSNSSQAYQEAFEISKKEMQPTHPIRLGLALNFSVFY
YEILNSPEKACSLAKTAFDEAIAELDTLNEESYKDSSTLIMQLLRDNLTLTWSENQGDEGD
AGEGEN
```

Filter query sequence (-F) : (DUST with BLASTN, SEG with others)

Expectation value (-e, real) :

Number of on-line description (-v, int) :

Number of alignments to show (-b, int) :

Matrix (-M) :

Expectation value threshold for inclusion in multipass model(-h, real) :

Maximum number of passes to use in multipass version (-j, int) : (limited to 10)



Pôle BioInformatique Lyonnais

Geno3D

Geno3D is the [IBCP](#) contribution to [PBIL](#) in Lyon, France

[\[HOME\]](#) [\[GENO3D\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#) [\[NPS@\]](#) [\[SuMo\]](#) [\[PBIL\]](#)

Tuesday, April 8th 2008: more options for PSI-BLAST search ([see news](#))



Job GENO3D2 (ID: 17095) is running on GENO3D server (started on 20080529-122822).
Results will be shown below. Please wait and don't go back.

Run GENO3D2

FIRST STEP :

Select template(s) to use for each chain in one or more pdb target :

PSI-BLAST run 3 for UNK_170950_0

Templato selezionato

segmento

modelling

Link dell'allineamento

Link NPSA

% identità

TEMPLATE	E	FIRST	LAST	ID	ALIGNEMENT	COMMENT	NPSA link
<input checked="" type="checkbox"/> pdb2br9A-0	2.000000e-91	4	231	67.000000	see alignment	CELL REGULATOR PROTEIN 03-MAY-05 2BR9 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2bq0B-0	2.000000e-91	3	233	100.000000	see alignment	CELL REGULATOR PROTEIN 25-APR-05 2BQ0 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2bq0A-0	2.000000e-91	3	232	100.000000	see alignment	CELL REGULATOR PROTEIN 25-APR-05 2BQ0 CROSS_PDB	NPSA
<input type="checkbox"/> pdb1qjbB-0	6.000000e-91	3	233	88.000000	see alignment	COMPLEX (SIGNAL TRANSDUCTION/PEPTIDE) 23-JUN-99 1QJB CROSS_PDB	NPSA
<input type="checkbox"/> pdb2btpA-0	1.000000e-90	2	232	83.000000	see alignment	COMPLEX (SIGNAL TRANSDUCTION/PEPTIDE) 05-JUN-05 2BTP CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c1jA-0	2.000000e-90	3	232	88.000000	see alignment	SIGNALING PROTEIN/COMPLEX 13-SEP-05 2C1J CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c1jB-0	2.000000e-90	3	232	88.000000	see alignment	SIGNALING PROTEIN/COMPLEX 13-SEP-05 2C1J CROSS_PDB	NPSA
<input type="checkbox"/> pdb2o02B-0	2.000000e-90	3	232	88.000000	see alignment	PROTEIN BINDING/TOXIN 27-NOV-06 2O02 CROSS_PDB	NPSA
<input checked="" type="checkbox"/> pdb2c23A-0	5.000000e-90	3	232	97.000000	see alignment	SIGNALING PROTEIN/COMPLEX 24-AUG-05 2C23 CROSS_PDB	NPSA

Geno3D ID	E-value	Length	Score	Alignment	Description	Database	
<input type="checkbox"/> pdb1g92A-0	0.000000e-84	4	233	76.000000	see alignment	PROTEIN BINDING 12-DEC-02 1G92 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c74B-0	6.000000e-84	4	232	76.000000	see alignment	SIGNALING PROTEIN/PEPTIDE COMPLEX 17-NOV-05 2C74 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2b05F-0	7.000000e-84	4	231	76.000000	see alignment	CELL CYCLE 13-SEP-05 2B05 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2b05E-0	2.000000e-83	5	231	76.000000	see alignment	CELL CYCLE 13-SEP-05 2B05 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c63A-0	2.000000e-83	4	232	76.000000	see alignment	SIGNALING PROTEIN/PEPTIDE COMPLEX 07-NOV-05 2C63 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c63B-0	2.000000e-83	4	232	76.000000	see alignment	SIGNALING PROTEIN/PEPTIDE COMPLEX 07-NOV-05 2C63 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c63C-0	2.000000e-83	4	232	76.000000	see alignment	SIGNALING PROTEIN/PEPTIDE COMPLEX 07-NOV-05 2C63 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2c63D-0	2.000000e-83	4	232	76.000000	see alignment	SIGNALING PROTEIN/PEPTIDE COMPLEX 07-NOV-05 2C63 CROSS_PDB	NPSA
<input type="checkbox"/> pdb1yz5B-0	2.000000e-82	4	232	68.000000	see alignment	SIGNALING PROTEIN 28-FEB-05 1YZ5 CROSS_PDB	NPSA
<input type="checkbox"/> pdb1ywtB-0	4.000000e-82	3	231	67.000000	see alignment	SIGNALING PROTEIN/DE NOVO PROTEIN 18-FEB-05 1YWT CROSS_PDB	NPSA
<input type="checkbox"/> pdb1ywtA-0	6.000000e-82	3	231	67.000000	see alignment	SIGNALING PROTEIN/DE NOVO PROTEIN 18-FEB-05 1YWT CROSS_PDB	NPSA
<input type="checkbox"/> pdb1yz5A-0	2.000000e-81	3	233	66.000000	see alignment	SIGNALING PROTEIN 28-FEB-05 1YZ5 CROSS_PDB	NPSA
<input type="checkbox"/> pdb2npmB-0	3.000000e-81	5	232	62.000000	see alignment	PROTEIN BINDING 27-OCT-06 2NPM CROSS_PDB	NPSA
<input type="checkbox"/> pdb2npmA-0	7.000000e-81	5	232	62.000000	see alignment	PROTEIN BINDING 27-OCT-06 2NPM CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a37A-0	7.000000e-73	3	230	75.000000	see alignment	COMPLEX (SIGNAL TRANSDUCTION/PEPTIDE) 28-JAN-98 1A37 CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a37B-0	7.000000e-73	3	230	75.000000	see alignment	COMPLEX (SIGNAL TRANSDUCTION/PEPTIDE) 28-JAN-98 1A37 CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a4oA-0	1.000000e-72	3	230	76.000000	see alignment	SIGNAL TRANSDUCTION 01-FEB-98 1A4O CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a4oB-0	1.000000e-72	3	230	76.000000	see alignment	SIGNAL TRANSDUCTION 01-FEB-98 1A4O CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a4oC-0	1.000000e-72	3	230	76.000000	see alignment	SIGNAL TRANSDUCTION 01-FEB-98 1A4O CROSS_PDB	NPSA
<input type="checkbox"/> pdb1a4oD-0	1.000000e-72	3	230	76.000000	see alignment	SIGNAL TRANSDUCTION 01-FEB-98 1A4O CROSS_PDB	NPSA
<input type="checkbox"/> pdb2ijpB-0	8.000000e-40	6	234	24.000000	see alignment	SIGNALING PROTEIN 30-SEP-06 2IJP CROSS_PDB	NPSA
<input type="checkbox"/> pdb2ijpA-0	4.000000e-35	20	236	23.000000	see alignment	SIGNALING PROTEIN 30-SEP-06 2IJP CROSS_PDB	NPSA
<input type="checkbox"/> pdb2ijpD-0	5.000000e-35	11	229	24.000000	see alignment	SIGNALING PROTEIN 30-SEP-06 2IJP CROSS_PDB	NPSA
<input checked="" type="checkbox"/> pdb2ijpC-0	8.000000e-35	21	236	23.000000	see alignment	SIGNALING PROTEIN 30-SEP-06 2IJP CROSS_PDB	NPSA

PSI-BLAST hasn't converged in 3 runs.

Select template



Pôle BioInformatique Lyonnais

Geno3D

Geno3D is the [IBCP](#) contribution to [PBIL](#) in Lyon, France

[\[HOME\]](#) [\[GENO3D\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#) [\[NPS@\]](#) [\[SuMo\]](#) [\[PBIL\]](#)

Tuesday, April 8th 2008: more options for PSI-BLAST search ([see news](#))



Job TEMPLATES VALIDATION ... (ID: 19250) is running on **GENO3D** server (started on 20080529-125249).
Results will be shown below. **Please wait and don't go back.**

SECOND STEP :

Enter your e-mail address :

THIRD STEP :

Choose number of model to generate :

EXPERT OPTIONS :

Inter/intra restraints ratio :

Distance restraints cut off (Angstrom):

Margin in distance restraints (Angstrom):

Margin in angle restraints (degree) :

Maximal number of distance restraints :

Save full template in superposition : No Yes

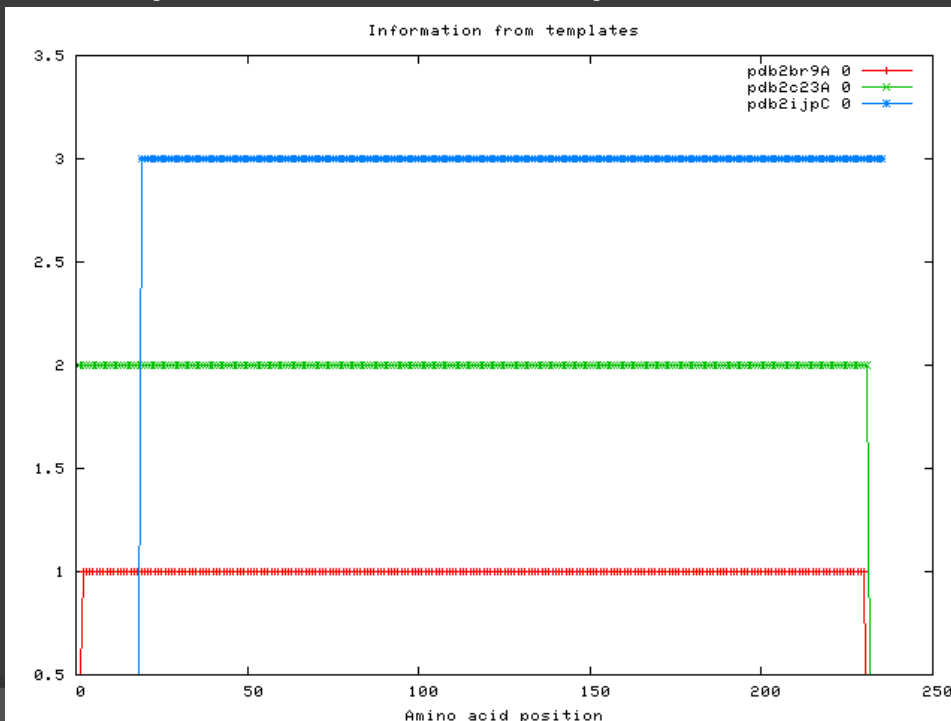
User : public@151.81.11.42. Last modification time : Thu May 29 12:52:57 2008. Current time : Thu May 29 12:52:57 2008 This service is supported by ['Ministere de la recherche'](#) , ['Programme Bioinformatique inter-EPST'](#) , [CNRS \(IMABIO, COMI, GENOME\)](#) and [Région Rhône-Alpes \(Programme EMERGENCE\)](#) . [Comments](#).

-Sequence of this chain :

MDKSELVQKA KLAEQAERYD DMAAAMKAVT
EQGHLSNEE RNLLSVAYKN VVGARRSSWR
VISSIEQKTE RNEKKQQMGK EYREKIEAEL
QDICNDVLEL LDKYLIPNAT QPESKVFYLK
MKGDFRYLS EVASGDNKQT TVSNSQQAYQ
EAFEISKKEM QPTHPIRLGL ALNFSVFYFE
ILNSPEKACS LAKTAFDEAI AELDTLNEES
YKDSTLIMQL LRDNLTLWTS ENQG

Template information (Sov)	Alignment Identity	Secondary
pdb2br9A_0	68.0%	ali_antheprot ali_clustalw
pdb2c23A_0	100.0%	ali_antheprot ali_clustalw
pdb2ijpC_0	26.0%	ali_antheprot ali_clustalw

- Template at each amino acid position :



Stereochemical quality of models with PROCHECK :

Model1	core	allowed	generously	disallowed
	83.7%	14.0%	1.8%	0.5%

-ramachandran plot : [jpeg](#) [postscript](#)

-ramachandran plots for all residue types : [jpeg](#) [postscript](#)

-main-chain parameters : [jpeg](#) [postscript](#)

-side-chain parameters : [jpeg](#) [postscript](#)

-residue properties : [jpeg](#) [postscript](#)

Model 2	86.4%	11.3%	1.8%	0.5%
----------------	-------	-------	------	------

-ramachandran plot : [jpeg](#) [postscript](#)

-ramachandran plots for all residue types : [jpeg](#) [postscript](#)

-main-chain parameters : [jpeg](#) [postscript](#)

-side-chain parameters : [jpeg](#) [postscript](#)

-residue properties : [jpeg](#) [postscript](#)

Model 3	84.2%	14.0%	0.5%	1.4%
----------------	-------	-------	------	------

ramachandran plot : [jpeg](#) [postscript](#)

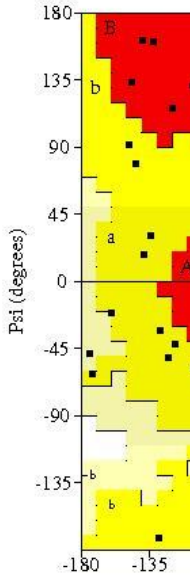
-ramachandran plots for all residue types : [jpeg](#) [postscript](#)

-main-chain parameters : [jpeg](#) [postscript](#)

-side-chain parameters : [jpeg](#) [postscript](#)

-residue properties : [jpeg](#) [postscript](#)

Ramachandran Plot

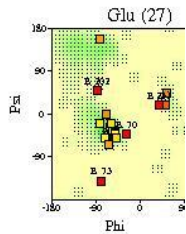
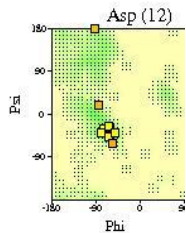
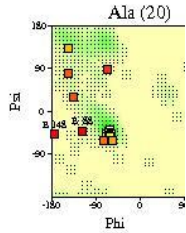


Residues
Residues
Residues
Residues
Number
Number
Number
Number
Total aa

model_1_01.ps

Ramachandran plots for all residue types

model_1

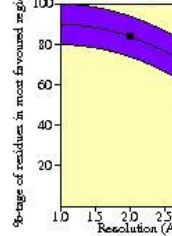


Numbers of residues are shown
Shading shows favourable conformation

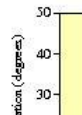
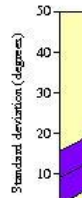
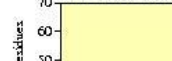
model_1_01.ps

Ma

a. Ramachandran plot

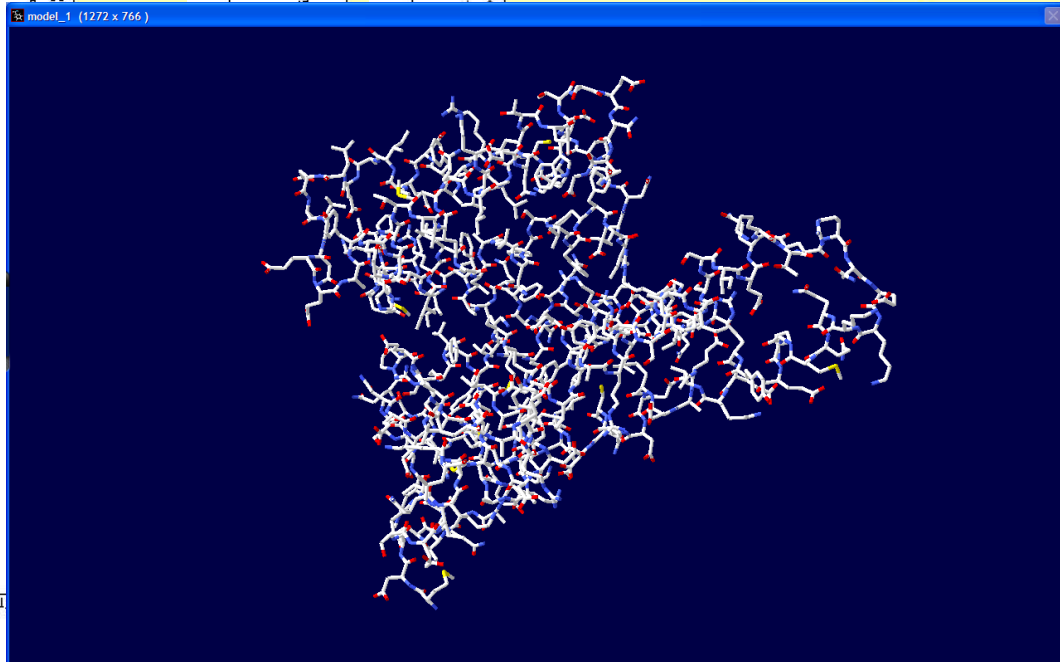
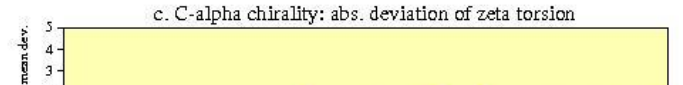
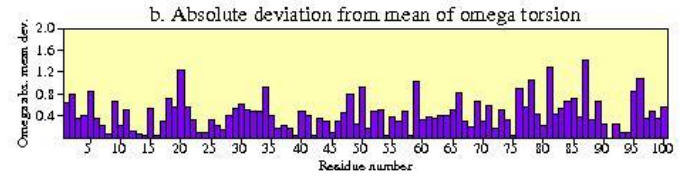
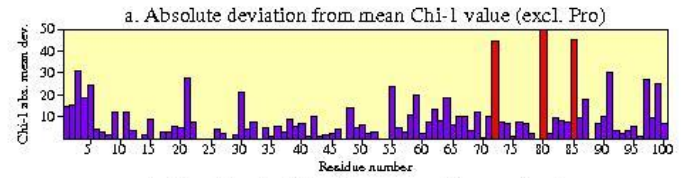


c. Measure of bad non-residues

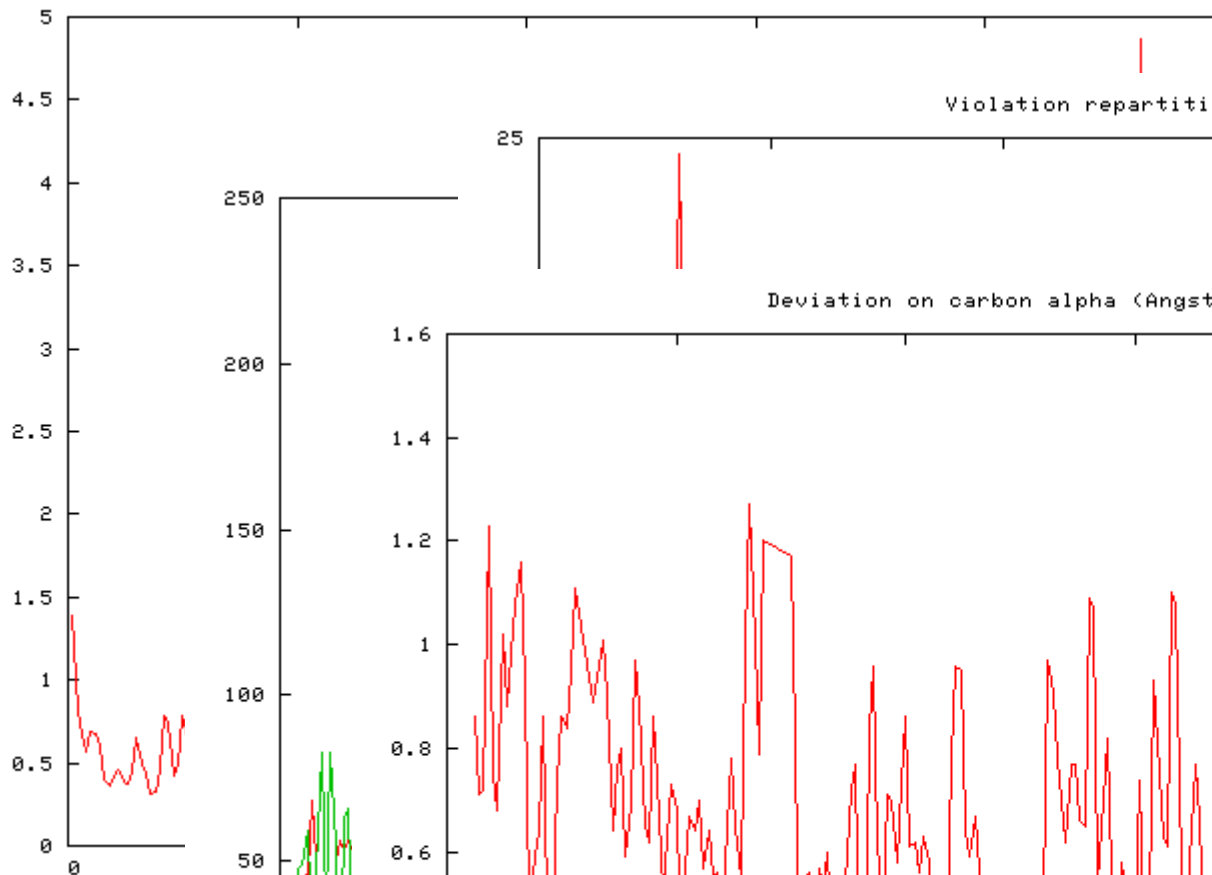


model_1

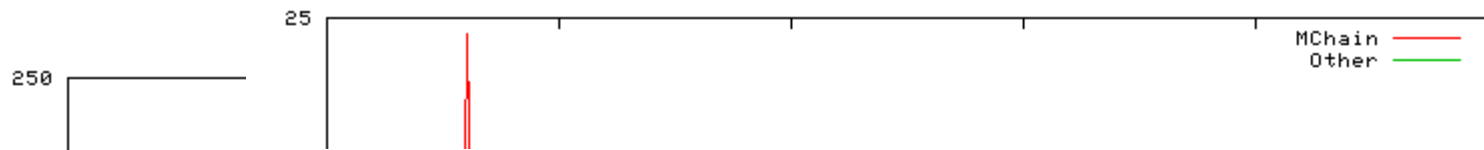
Residue properties model_1



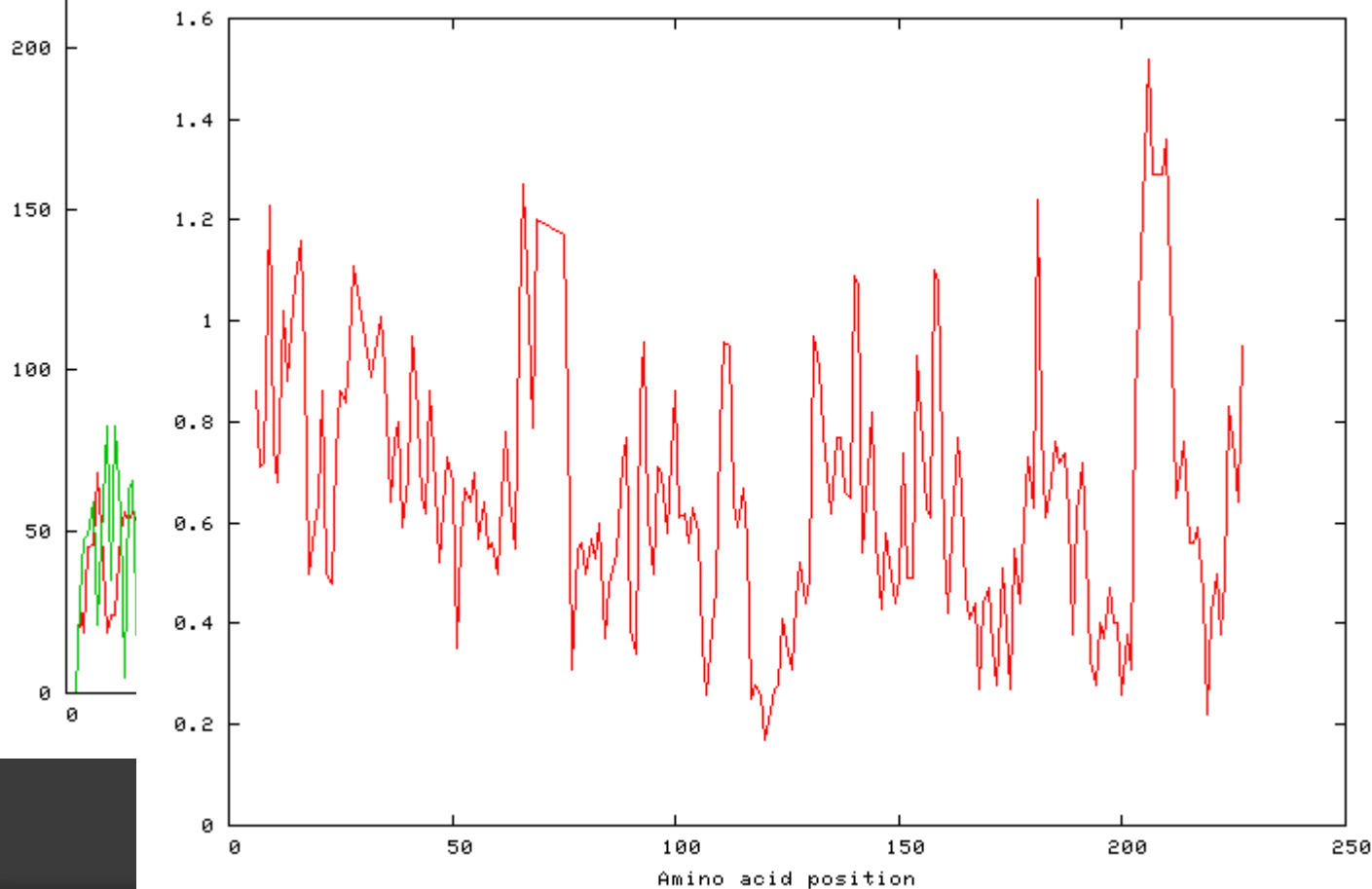
Deviation (Angstrom) on carbon alpha



Violation repartition (%)



Deviation on carbon alpha (Angstrom)



**-Structural agreement between models
(RMSD in angstrom) :**

	Model 1	Model 2	Model 3
Model 1	0.00	0.75	0.78
Model 2	0.75	0.00	0.69
Model 3	0.78	0.69	0.00

Models energy (kcal/mol):

. model 1 : -11081.70
. model 2 : -11577.00
. model 3 : -11479.00

- Number of violation of these intrachain restraints :

	MChain	Other
Model 1	205(2.70%)	111(1.40%)
Model 2	114(1.50%)	65(0.80%)
Model 3	167(2.20%)	95(1.20%)

Sviluppi futuri:

- modellizzazione di dimeri;
- Includere il ligando nel processo di modellizzazione molecolare;
- possibilità di caricare allineamenti personali.

Se ci troviamo nella situazione intermedia dobbiamo fare altre valutazioni

Confronto della predizione di struttura secondaria

Residui idrofobici interni conservati

Ponti a disolfuro conservati

Pattern funzionali conservati

Modellizzazione comparativa (o per similarità di sequenza)

**Permette di costruire la struttura
tridimensionale di una proteina**

**sulla base della SIMILARITÀ DI SEQUENZA con
un'altra proteina**

di struttura NOTA

che viene usata come STAMPO.

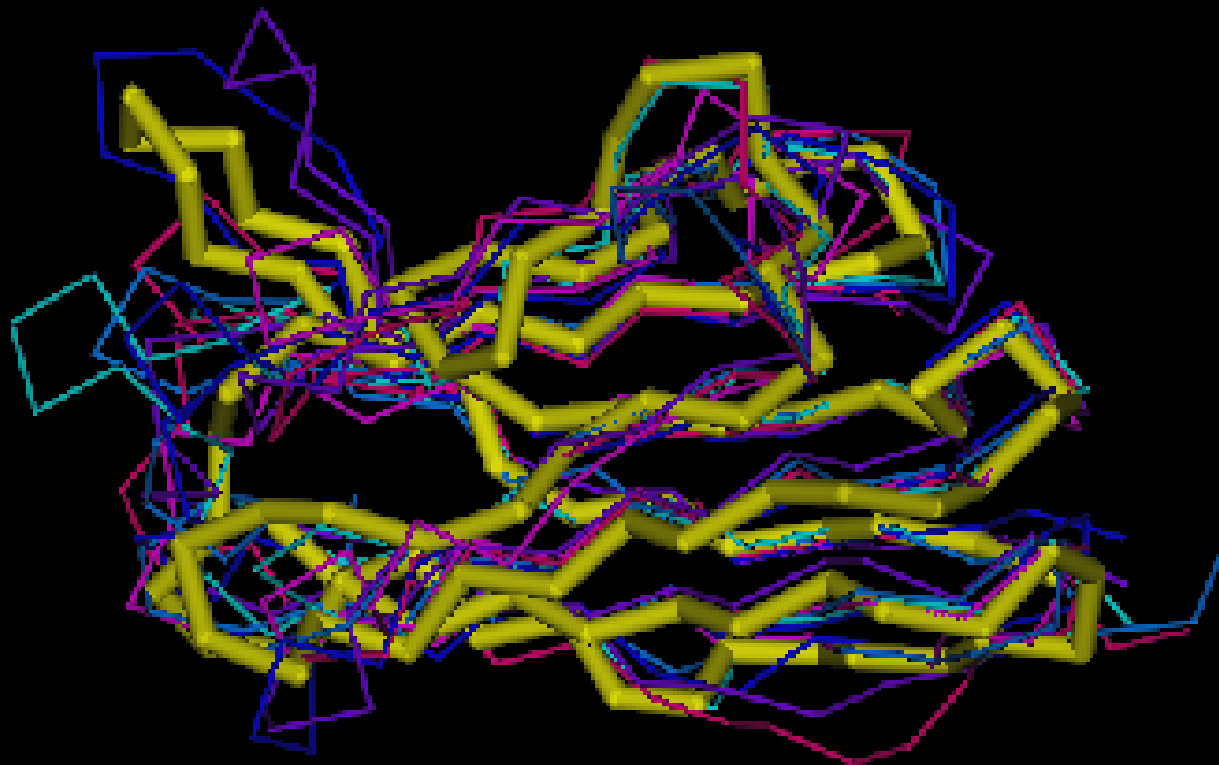
Passi fondamentali

1. Allineamento di sequenza con la/le proteina/e “stampo”

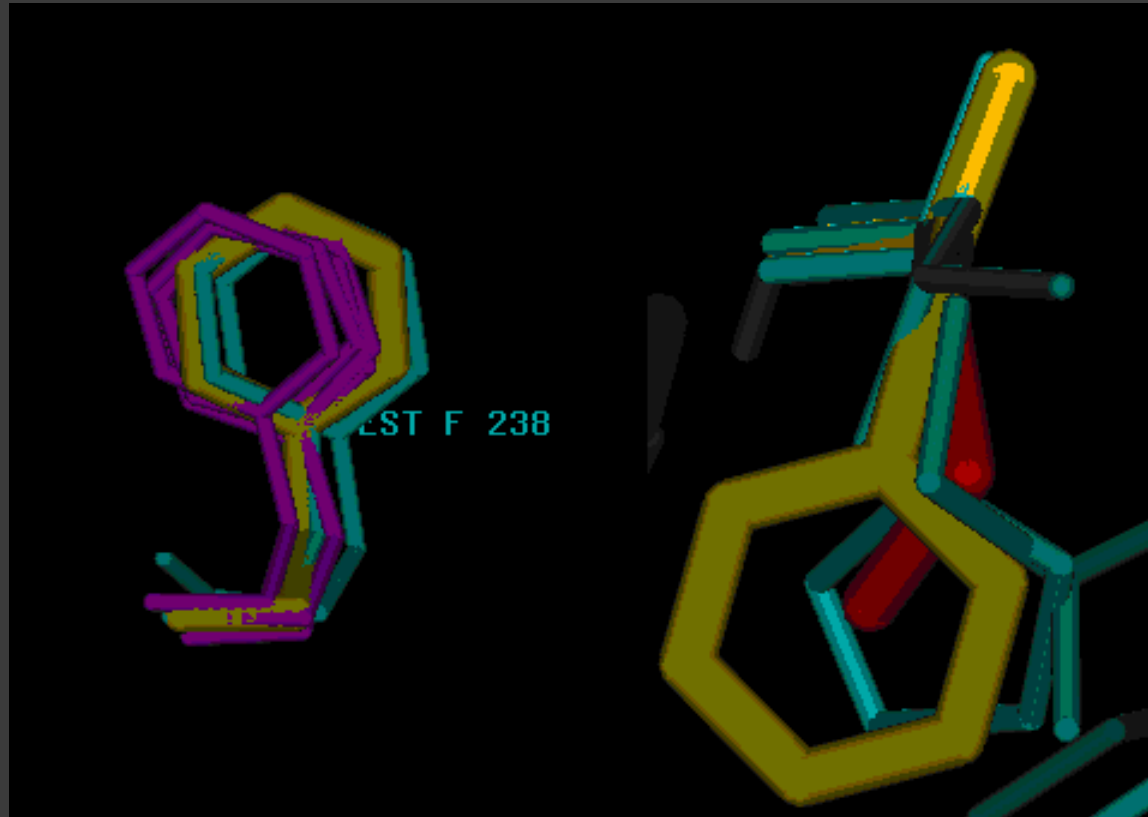
* aa identici
. aa simili

4mdh .aa	1	SEPIRVLVTG	AAGQIAYSL	YSIGNGSVFG	KDQPIILVLL	DITPMMGVLD
11BMD	1	KAPVRVAVTG	AAGQIGYSL	FRIAAGEMLG	KDQPVILQLL	EIPQAMKALE
		* . * * * *	* * * * * . * * * *	. * . * . . *	* * * * . * * * *	. * . * . *
4mdh .aa	51	GVLMEIQDCA	LPLLKDVIA	DKEEIAFKDL	DVAILVGSMP	RRDGMERKDL
11BMD	51	GVMELEDCA	FPLLAGLEAT	DDPDVAFKDA	DYALLVGAAP	RKAGMERRDL
		* * . * * * . * * *	. * * * * . * * *	* . . . * * * *	* * . * * * * . *	* . * * * * . * *
4mdh .aa	101	LKANVKIFKC	QGAALDKYAK	KSVKVIIVVGN	PANTNCLTAS	KSAPSIPKEN
11BMD	101	LQVNGKIFTE	QGRALAEVAK	KDVKVLVVG	PANTNALIAY	KNAPGLNPRN
		* . * * * *	* * * * * * *	* * * * * . * * * *	* * * * * * *	* . * * * . *
4mdh .aa	151	FSCLTRLDHN	RAKAQIALKL	GVTSDDVKNV	IIWGNHSSTQ	YPDVNHAKVK
11BMD	151	FTAMTRLDHN	RAKAQLAKKT	GTGVDRIIRM	TVWGNHSSTM	FPDLFHAEDV
		* . . * * * * * *	* * * * * . * *	* . * * * * * * * *	. * * . * * *
4mdh .aa	201	LQAKEVGVYE	AVKDDSWLKG	EFITTVQQRG	AAVIKARKLS	SAMSAKAIC
11BMD	201	GRP----AIE	LVDME-WYEK	VFIPTVAQRG	AAIIQARGAS	SAASAANAAI
	 * *	* * . * * * *	* * * * * * *	* * . * * * * *	* * * * * . *
4mdh .aa	251	DHVRDIWFGT	PEGEFVSMGI	ISDGNVYGVV	DDLAYSFPVT	IKDKTKWIVE
11BMD	246	EHIRDWALGT	PEGDWVSMVA	PSQGE-YGIP	EGIVYSFPVT	AKDGAYRVVE
		. * * * * . * *	* * * . * * * . .	* * * * * * * * * * * *	* * *
4mdh .aa	301	GLPINDFSRE	KMDLTAKELA	E EKETA FEFL	SSA	
11BMD	295	GLEINEFARK	RMEITAQELL	DEMEQVKALG	LI	
		* * * * * . *	. * * * * * *	. * * *		

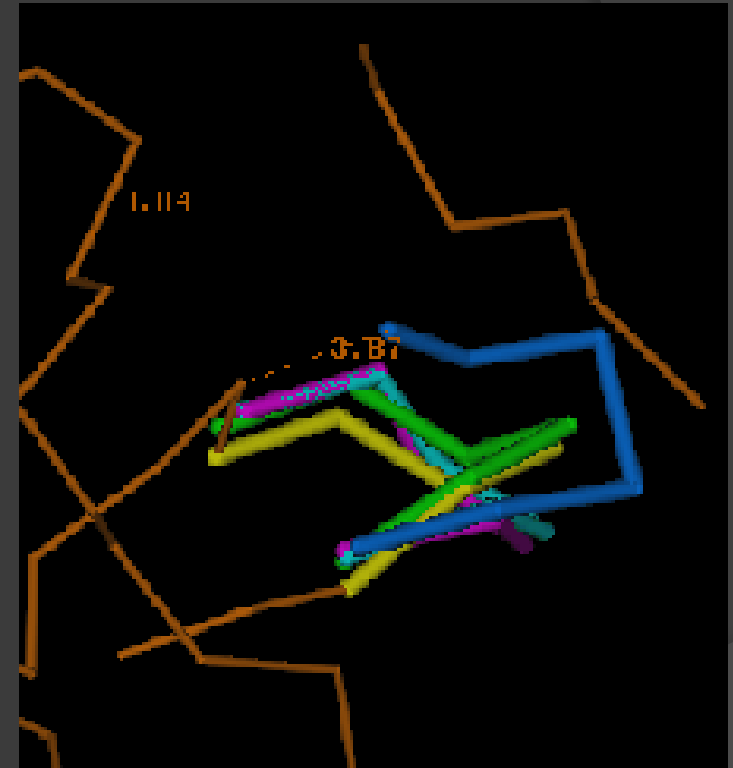
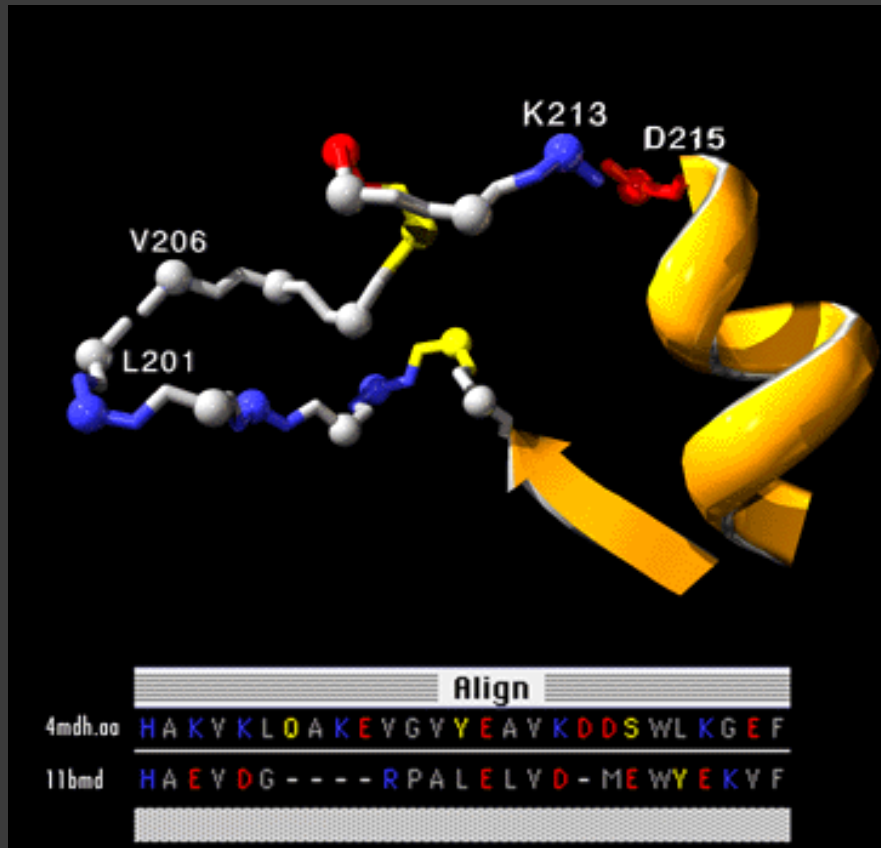
2. Costruzione dello scheletro



3. Visualizzazione delle catene laterali e risoluzione di ingombri sterici



4. Inserimento dei loop corrispondenti a “buchi” nell’allineamento



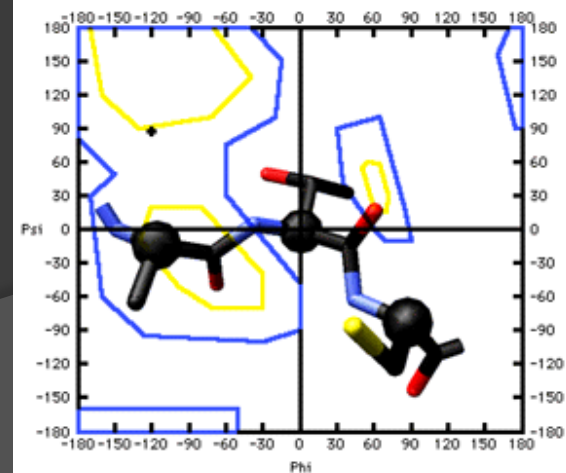
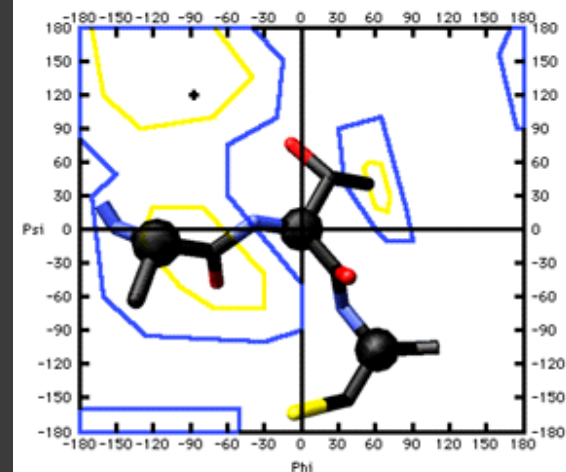
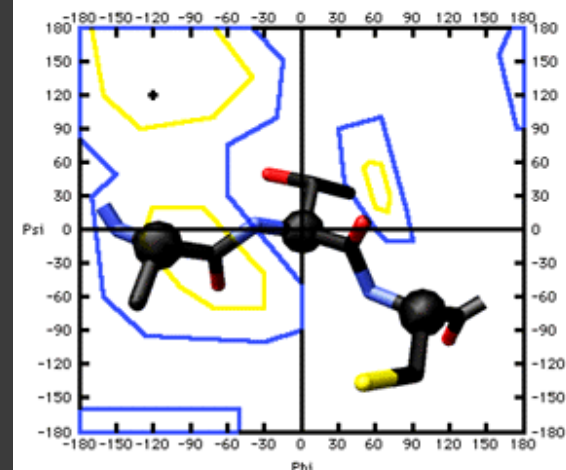
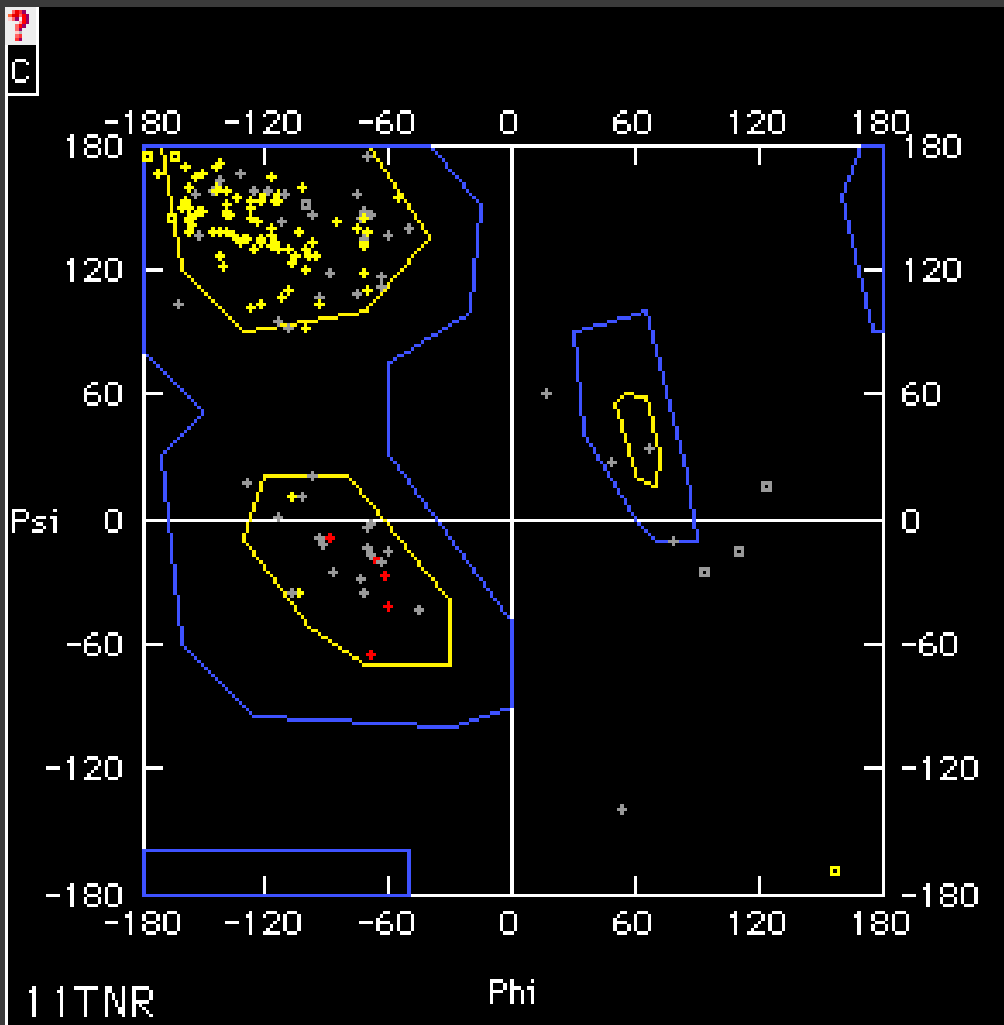
5. Ottimizzazione del modello

Regolarizzazione
di legami, angoli e
torsioni

Eliminazioni di
clash strutturali

Minimizzazione
energetica

6. Controllo della qualità del modello



Le proteine cambiano molto rapidamente la loro sequenza.

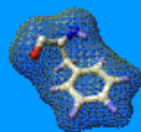
Difficile indicare il grado di similarità necessario per dimostrare in modo non ambiguo che due proteine siano OMOLOGHE.

Doolittle dichiara:

1. Se due sequenze sono più lunghe di 100 aa e l'identità è maggiore del 25% esse sono plausibilmente correlate
2. Se l'identità è tra il 15-25 % potrebbero essere correlate
3. Se l'identità è sotto il 15% probabilmente non sono correlate

[ExPASy Home page](#)[Site Map](#)[Search ExPASy](#)[Contact us](#)

Search for

[Download](#)[User Guide](#)[Tips & Tricks](#)[Tutorial](#)[Feedback](#)[Index](#)[Art Gallery](#)[Mirror sites](#)[References](#)

GlaxoSmithKline R&D
&
the Swiss Institute of Bioinformatics

present

Deep View Swiss-PdbViewer

by

Nicolas Guex , Alexandre Diemand , Manuel C. Peitsch , & Torsten Schwede

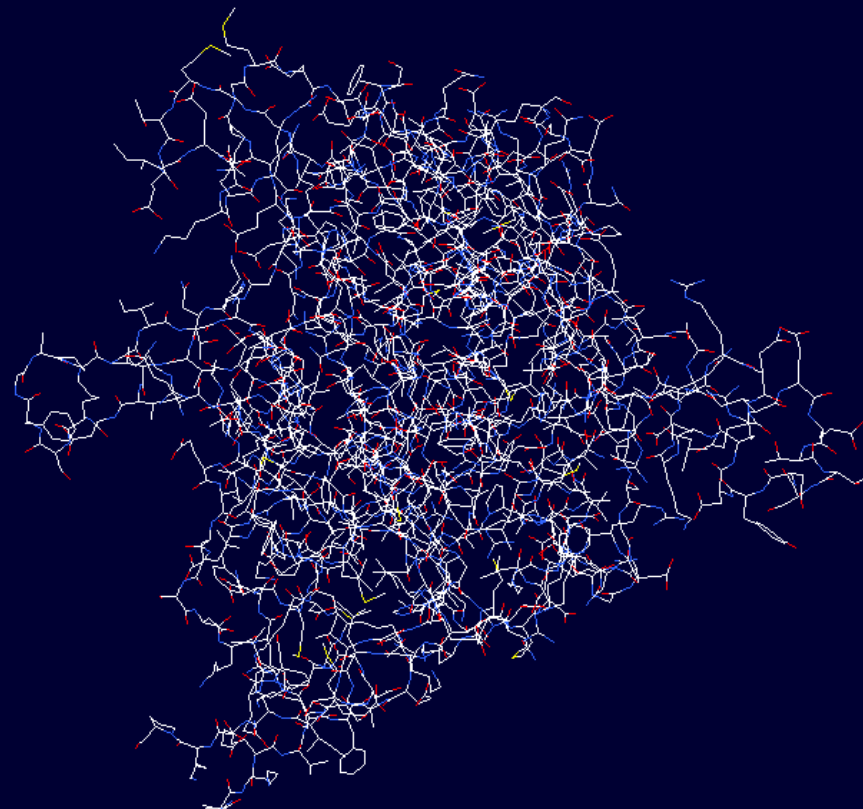
*This mirror site is
hosted by*



News:

- Please update your preferences for the [DeepView Network Service](#)
- Current version: [Deep View - spdbv 3.7](#)

11OK (975 x 768)

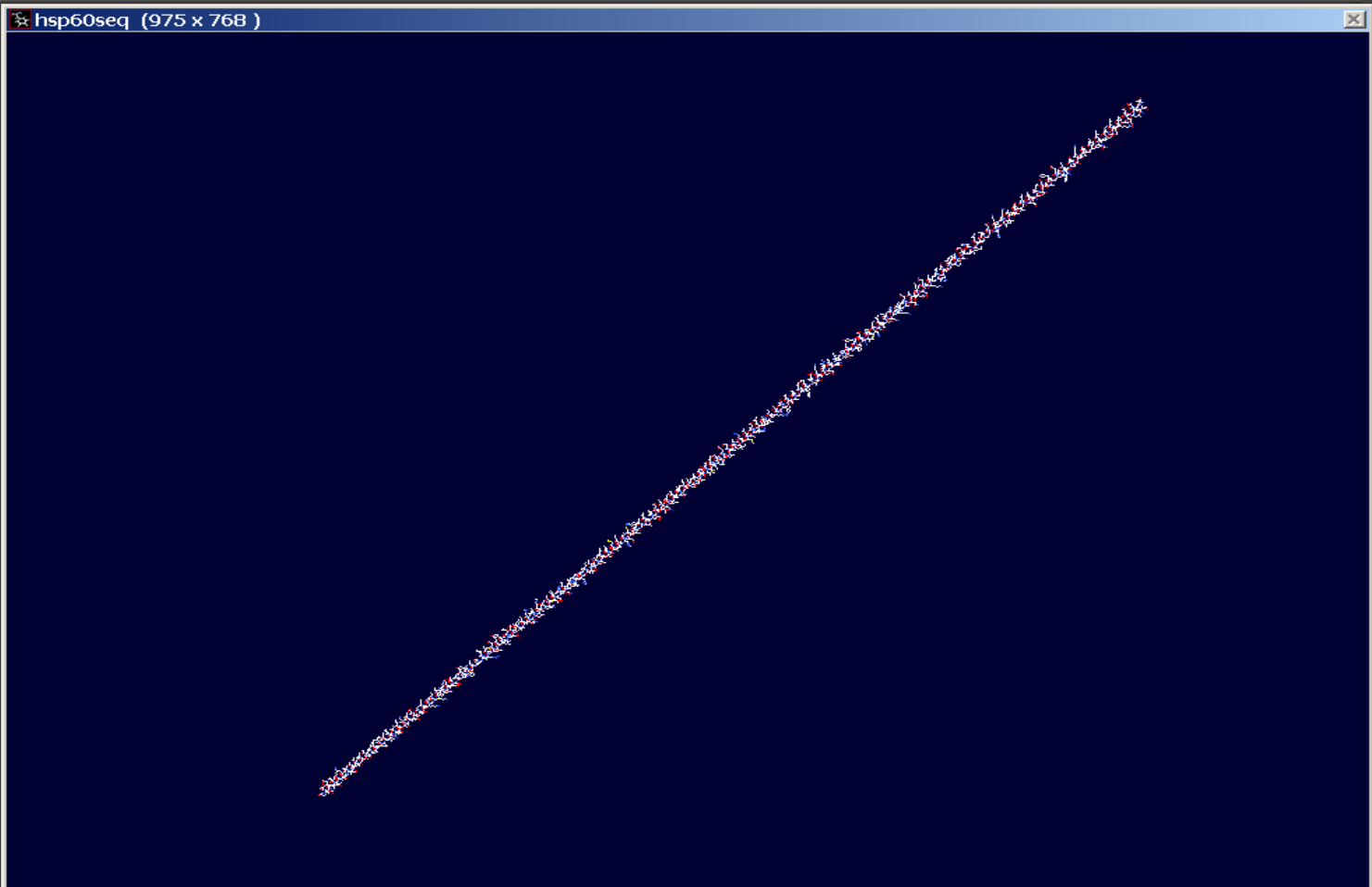


DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help

1.5Å 60.1° COOR W LEU41 1Å ROTATE FORSIOH ?

Move All



Alignment

hsp60seq	MSEQEKLSNYNADKKLFSGI DKLFIQIVKGSYGPKQSLSP T S F F K E R G F Y A I S Q T E L S N S Y E N L G V D F A K A M V N K I H K E H S D G A T
110K	AAKEVKFNSDARDRMLKGVNI LADAVKVT LGPKGRNVVI DKSFGAPRI T K D G V S V A K E I E L S D K F E N M G A Q M V R E V A S R T N D E A

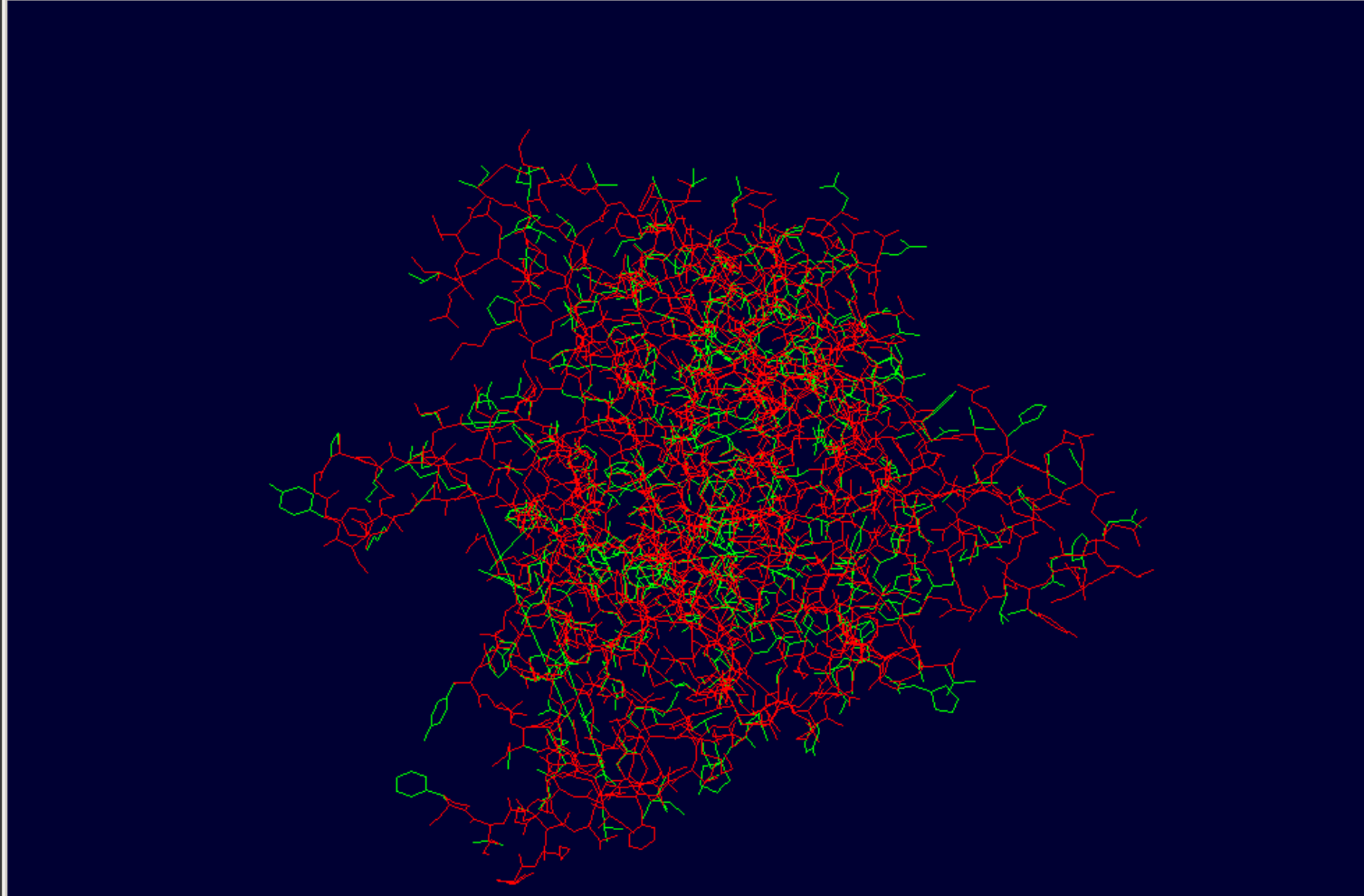
hsp60seq: GLU5 rms: 153.1

DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help



hsp60seq (975 x 768)



Alignment

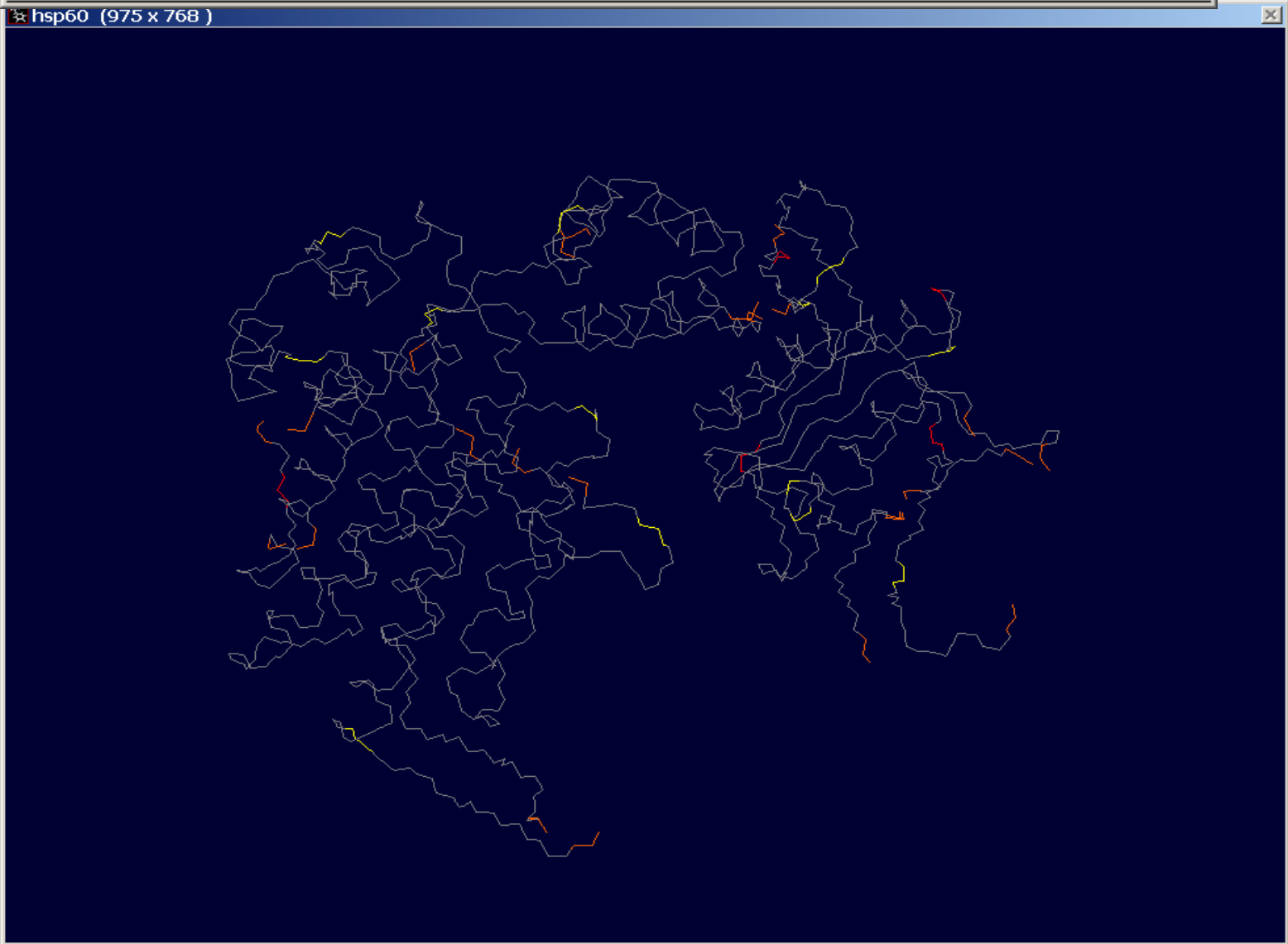
hsp60seq MSEDEKLSNYNADKFLFSGI DKLFQIV KGSYGPKQS SPTSFKERGFIYAI SQTELSNSYENLGVDFAKAM
110K AAKEVKFNS DARDRMLKGVNI LADAVKVTLPKGRNVVI DKSFGAPRI TKIGVSVAKEI ELSDKFENMGAQMVREVASRTNDEA

DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help

1.5 Å 60.1° 60.8° LEU41 MUTATE TORSION ?

Move All



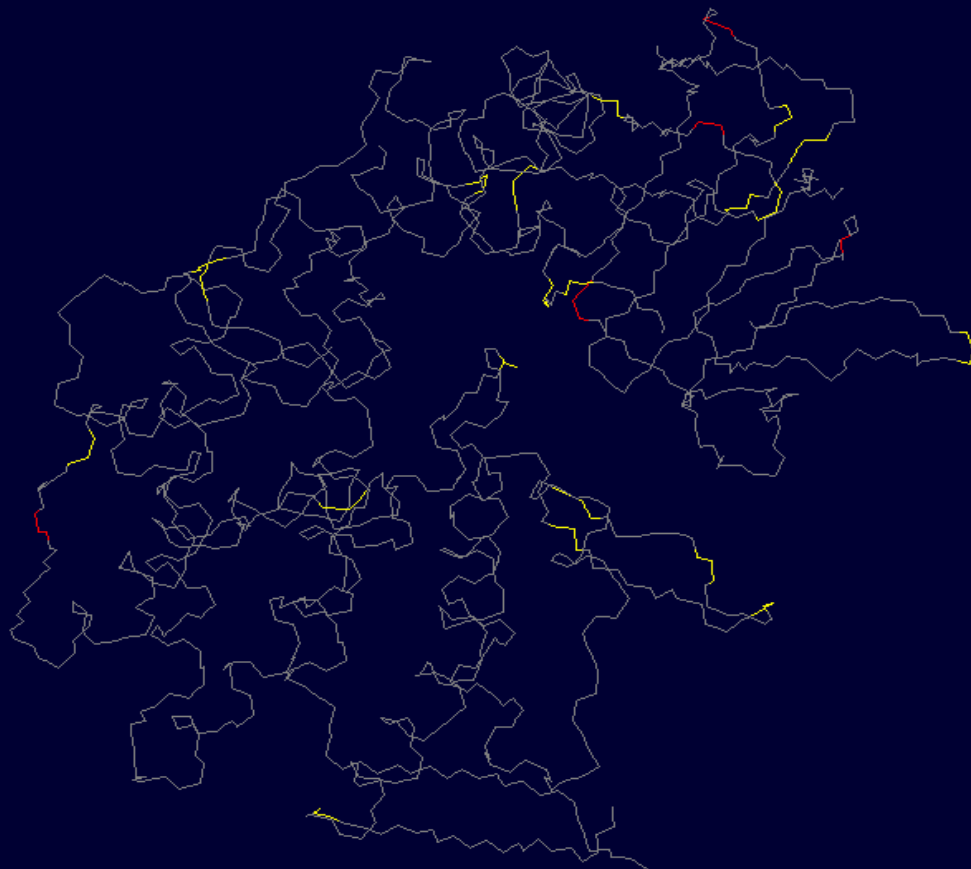
DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help



Move All

hsp60_9 (975 x 768)



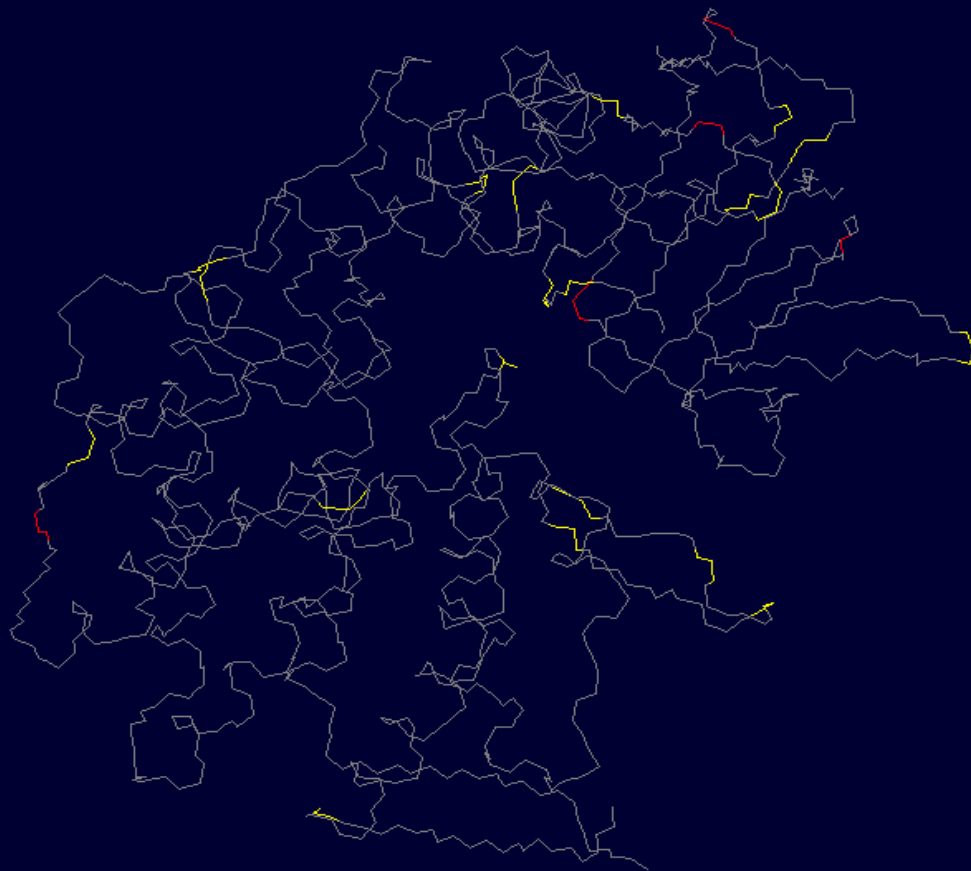
DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help



Move All

hsp60_9 (975 x 768)

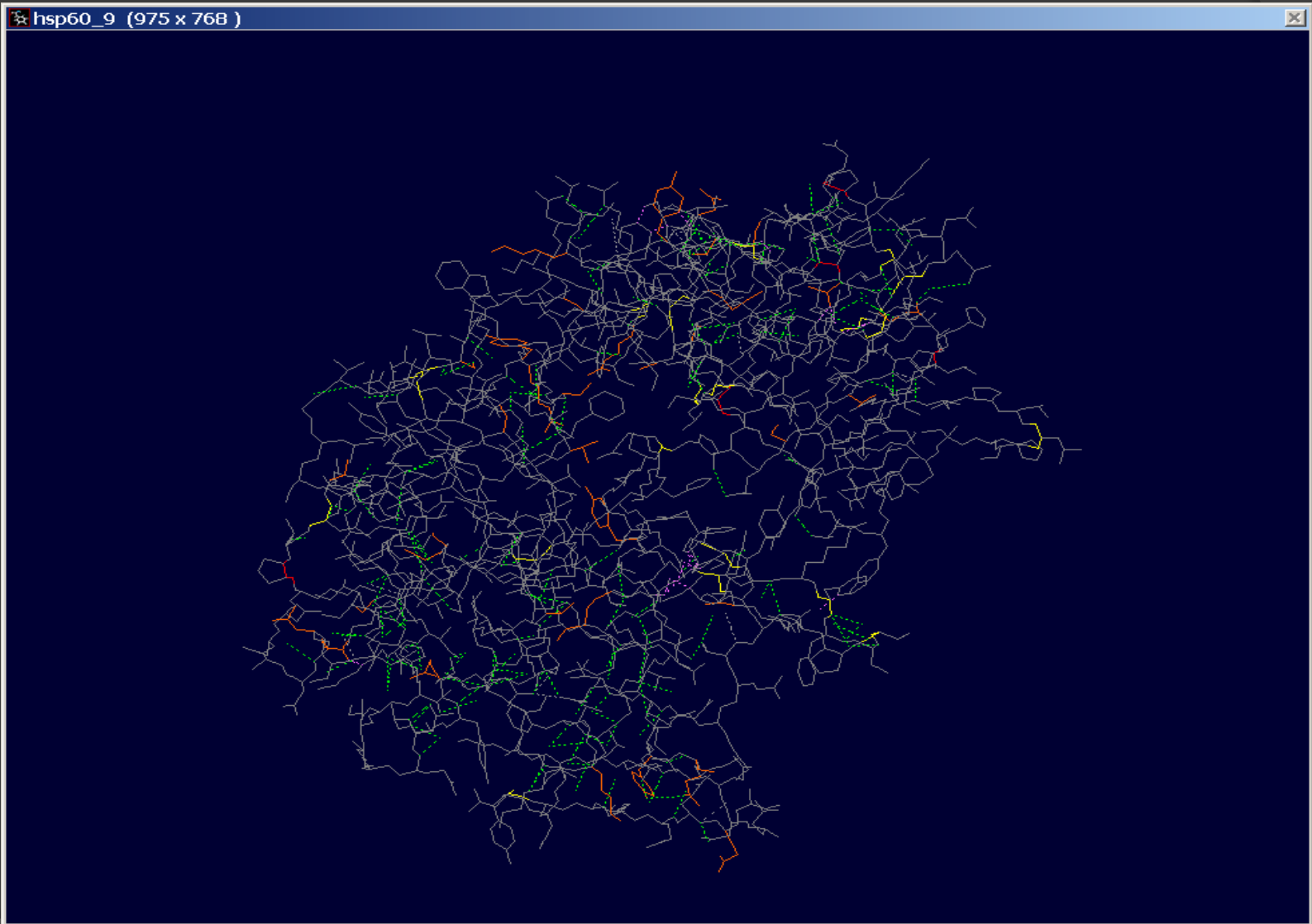


DeepView / Swiss-PdbViewer 3.7 (SP5)

File Edit Select Build Tools Fit Display Color Preferences SwissModel Window Help

15Å 60.1° 60.8° LEU41 1Å MUTATE TORSION ?

Move All



M:\local\wbin\deepview37sp5\temp\hsp60_9.E1										
*LEU	438	0.957	1.778	4.831	2.419	1412.27	-0.27	0.0000	// E=	1421.982
LYSH	439	0.356	2.717	1.912	1.177	-26.96	-6.46	0.0000	// E=	-27.259
LEU	440	1.311	0.734	3.539	0.240	547.80	-6.51	0.0000	// E=	547.117
*LEU	441	1.022	1.911	1.916	0.006	1400.84	-1.57	0.0000	// E=	1404.128
ALA	442	1.459	0.579	2.317	0.621	-26.89	-13.46	0.0000	// E=	-35.374
THR	443	0.360	0.979	2.199	0.596	5.62	-16.20	0.0000	// E=	-6.444
ASN	444	1.969	0.979	3.857	0.500	158.27	-155.41	0.0000	// E=	10.162
ALA	445	0.941	1.571	2.969	0.040	125.52	6.62	0.0000	// E=	137.669
ASP	446	1.253	2.531	5.586	1.638	103.25	-6.67	0.0000	// E=	107.587
LEU	447	0.731	1.262	3.441	0.461	20.67	1.45	0.0000	// E=	28.013
ASP	448	0.428	2.998	1.492	0.019	-4.60	15.61	0.0000	// E=	15.943
GLY	449	0.651	4.519	1.317	0.221	-10.35	24.68	0.0000	// E=	21.037
ASP	450	1.380	0.482	3.657	1.642	-17.89	11.91	0.0000	// E=	1.173
ALA	451	0.962	1.162	3.081	0.073	-8.43	-3.29	0.0000	// E=	-6.441
VAL	452	2.441	3.877	1.873	1.924	5.50	-0.24	0.0000	// E=	15.381
ILE	453	0.161	4.876	3.165	1.495	-17.45	-3.05	0.0000	// E=	-10.801
ALA	454	0.921	1.594	0.539	1.009	-13.00	3.19	0.0000	// E=	-5.752
*LYSH	455	2.422	1.094	19.372	0.417	49888.76	1.07	0.0000	// E=	49913.137
*LEU	456	3.081	1.742	1.621	0.439	1094.58	0.94	0.0000	// E=	1102.408
SER	457	7.411	25.125	4.264	1.514	19.54	-6.72	0.0000	// E=	51.130
SER	458	0.304	2.984	4.219	0.016	-5.72	-11.67	0.0000	// E=	-9.862
*LEU	459	0.844	1.922	2.514	0.011	51025.36	48.02	0.0000	// E=	51078.675
GLY	460	0.942	1.122	0.887	0.043	-5.97	54.92	0.0000	// E=	51.349
THR	461	0.289	2.615	3.382	1.272	-11.29	-1.08	0.0000	// E=	-4.811
THR	462	0.750	5.123	3.009	1.359	-8.95	0.60	0.0000	// E=	1.884
SER	463	0.536	1.606	2.512	0.493	-17.45	-10.70	0.0000	// E=	-23.005
LEU	464	1.823	6.804	3.995	1.526	29.20	34.46	0.0000	// E=	77.806
GLY	465	1.349	1.478	2.197	2.839	33.23	40.19	0.0000	// E=	81.283
ILE	466	0.835	5.533	16.903	1.263	-14.95	-0.43	0.0000	// E=	9.157
SER	467	0.338	0.808	2.261	1.532	49.52	-26.64	0.0000	// E=	27.826
VAL	468	0.415	3.748	2.012	0.042	63.38	10.86	0.0000	// E=	80.458
PHR	469	2.341	20.725	1.876	0.374	-17.44	24.90	0.0000	// E=	32.772
SER	470	0.565	1.405	2.308	0.015	-3.16	16.55	0.0000	// E=	17.684
ARG	471	1.133	5.067	2.547	1.048	23.79	-248.34	0.0000	// E=	-214.759
GLU	472	1.020	2.174	3.459	0.858	-18.37	-3.90	0.0000	// E=	-14.761
*ILE	473	0.459	1.286	3.433	0.517	2139.81	-5.65	0.0000	// E=	2139.862
*GLU	474	2.499	3.981	3.249	2.433	2095.46	-3.45	0.0000	// E=	2104.176
ASP	475	0.842	2.193	5.190	0.003	-10.37	-8.98	0.0000	// E=	-11.127
LEU	476	0.612	2.546	2.252	0.023	111.61	-18.97	0.0000	// E=	98.074
ILE	477	0.732	1.507	3.281	0.564	25.14	4.73	0.0000	// E=	35.948
ALA	478	0.232	0.449	0.031	0.590	-8.04	45.89	0.0000	// E=	39.148
GLY	479	0.318	2.249	0.275	0.003	-5.53	77.83	0.0000	// E=	75.142
GLY	480	0.945	0.235	2.916	1.280	-1.03	44.73	0.0000	// E=	49.071
ILE	481	0.743	2.333	2.782	0.019	63.08	-8.55	0.0000	// E=	60.411
LEU	482	0.998	9.213	4.624	0.872	-12.72	-5.35	0.0000	// E=	-2.366
ASP	483	0.338	1.031	5.407	0.438	-31.30	-10.59	0.0000	// E=	-34.670
SER	484	0.166	1.347	3.335	3.871	-26.02	-18.65	0.0000	// E=	-35.949
LEU	485	2.528	12.989	4.190	3.032	-2.93	-13.66	0.0000	// E=	6.149
ALA	486	0.709	2.149	0.817	0.864	-17.93	3.42	0.0000	// E=	-9.968
THR	487	0.570	1.239	3.199	0.921	4.24	-9.27	0.0000	// E=	0.894
THR	488	0.651	0.864	2.908	1.515	1.41	-9.32	0.0000	// E=	-1.971
SER	489	0.788	1.513	6.150	0.014	-21.80	-7.40	0.0000	// E=	-20.733
THR	490	1.795	1.121	0.913	1.445	-13.68	-27.75	0.0000	// E=	-36.153
ILE	491	0.399	0.960	2.818	0.696	345.12	-9.60	0.0000	// E=	340.397
LEU	492	0.372	4.440	1.394	0.029	31.28	-1.65	0.0000	// E=	35.870
ALA	493	0.826	0.496	2.168	1.521	24.50	-9.82	0.0000	// E=	19.687
GLN	494	0.392	2.690	4.304	0.015	14.39	-167.21	0.0000	// E=	-145.419
*ALA	495	0.844	0.323	2.039	0.851	719.68	-25.20	0.0000	// E=	698.546
LEU	496	0.894	1.218	1.866	0.934	67.96	-3.54	0.0000	// E=	69.333
ASP	497	1.417	0.563	3.898	1.211	27.04	-0.80	0.0000	// E=	33.324
THR	498	3.828	4.124	6.104	2.115	25.05	-17.02	0.0000	// E=	24.209
ALA	499	1.175	0.444	0.254	0.723	29.82	-3.80	0.0000	// E=	28.612
ILE	500	1.250	1.777	3.004	0.973	34.24	5.12	0.0000	// E=	46.360
LEU	501	1.881	3.150	2.349	0.001	-17.20	3.32	0.0000	// E=	-6.492
VAL	502	1.402	0.615	2.950	1.209	-0.26	-1.16	0.0000	// E=	4.763
LEU	503	7.358	18.588	6.675	0.649	38.65	-5.22	0.0000	// E=	66.702
SER	504	0.920	4.101	3.520	0.546	-0.22	-15.38	0.0000	// E=	-6.512
SER	505	0.208	3.488	4.181	0.018	-21.12	-0.40	0.0000	// E=	-13.628
LYSH	506	1.048	3.025	1.790	0.206	-12.00	4.45	0.0000	// E=	-1.484
ILE	507	1.261	1.977	3.291	1.900	19.03	1.87	0.0000	// E=	29.333
LEU	508	0.355	3.534	6.490	1.115	-19.48	-12.51	0.0000	// E=	-20.436
ILE	509	3.071	1.031	3.875	1.524	-2.56	1.85	0.0000	// E=	8.784
LEU	510	6.426	31.127	1.888	0.751	38.24	-13.81	0.0000	// E=	64.616
GLU	511	4.445	3.219	23.199	0.435	-24.44	-6.99	0.0000	// E=	-0.127
ASN	512	0.883	1.199	4.421	0.467	-16.04	-151.52	0.0000	// E=	-160.590
GLN	513	1.206	5.582	4.264	0.442	1.90	-74.12	0.0000	// E=	-60.727
OXT	513	0.000	0.000	0.000	0.000	-3.24	0.99	0.0000	// E=	-2.253

KJ/mol		3338.110	6302.236	2884.081	1743.990	240661.61	-7381.23	0.0000	// E=	247548.797



PROCHECK v.3.5.4

[Roman A Laskowski](#), [Malcolm W MacArthur](#), David K Smith, [David T Jones](#), [E Gail Hutchinson](#), A Louise Morris, [David S Moss](#) & [Janet M Thornton](#)

Checks the **stereochemical quality** of a protein structure, producing a number of PostScript plots analysing its **overall** and **residue-by-residue** geometry.

The plots can be in colour, if required.



[How to run the program](#)



[Operating Manual](#)



[Checks carried out](#)



[Sample outputs](#)



[References](#)

Availability

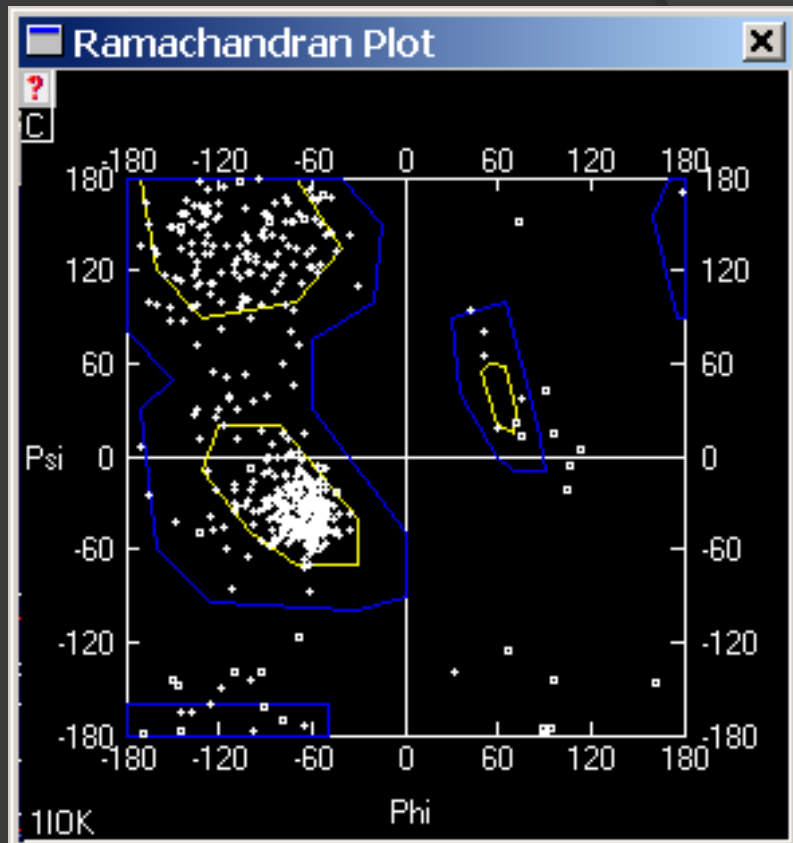
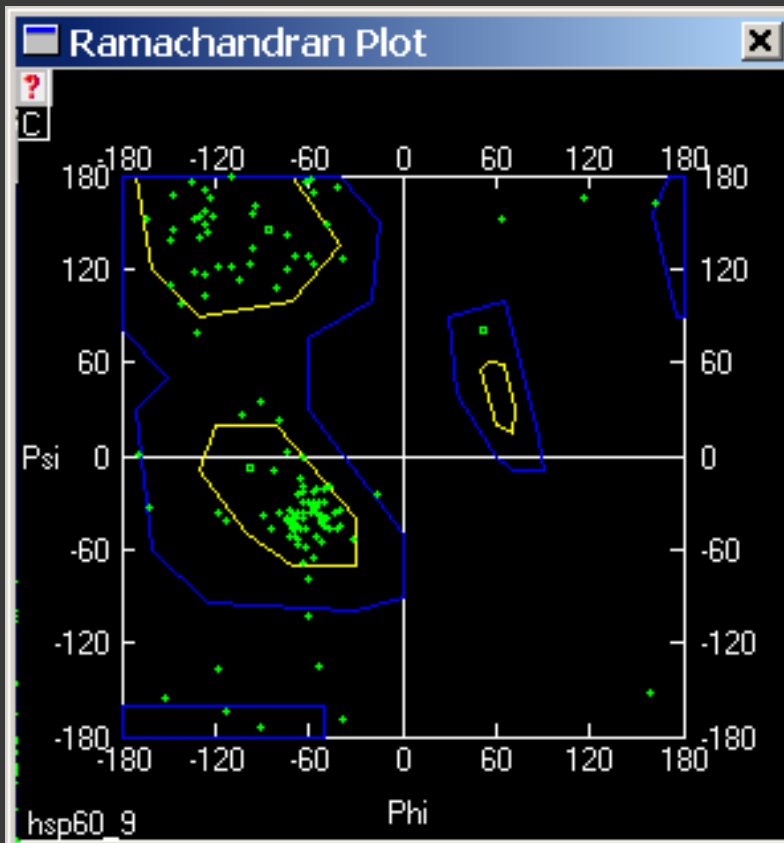
Available by anonymous ftp on: [ftp.biochem.ucl.ac.uk](ftp://ftp.biochem.ucl.ac.uk)

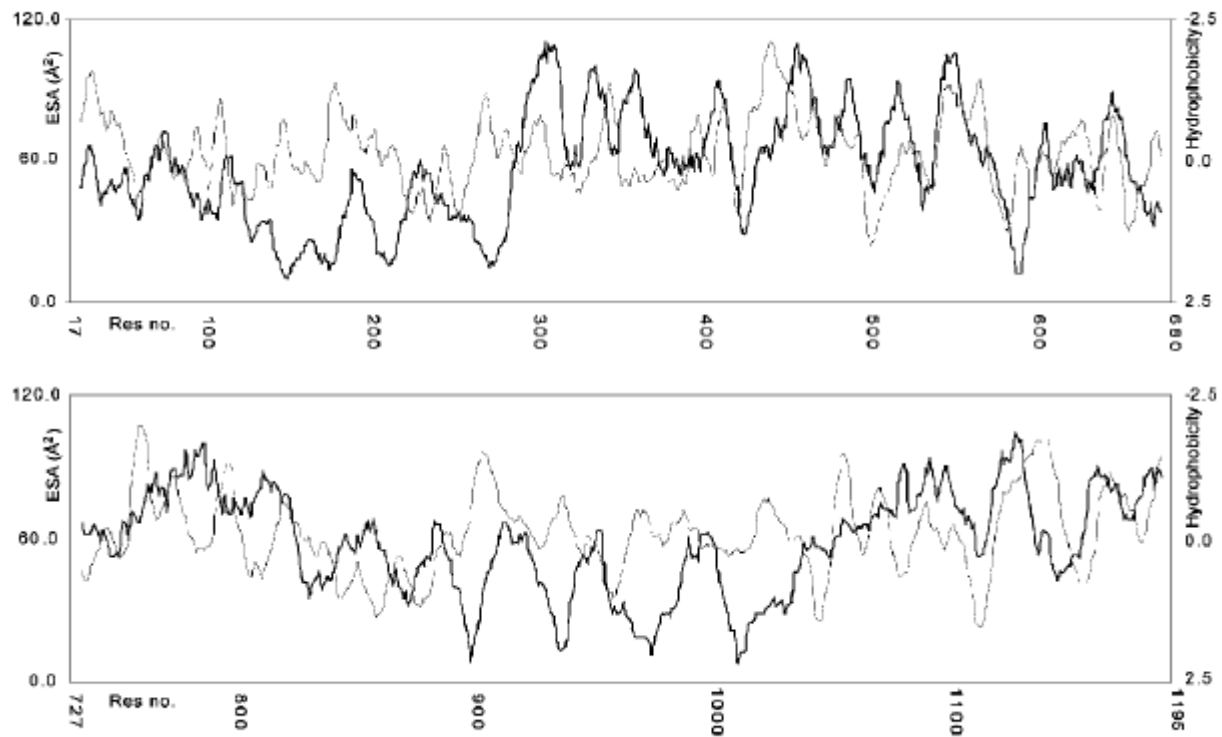
Source code can be picked up from one of two directories:-

pub/procheck/tar3_5

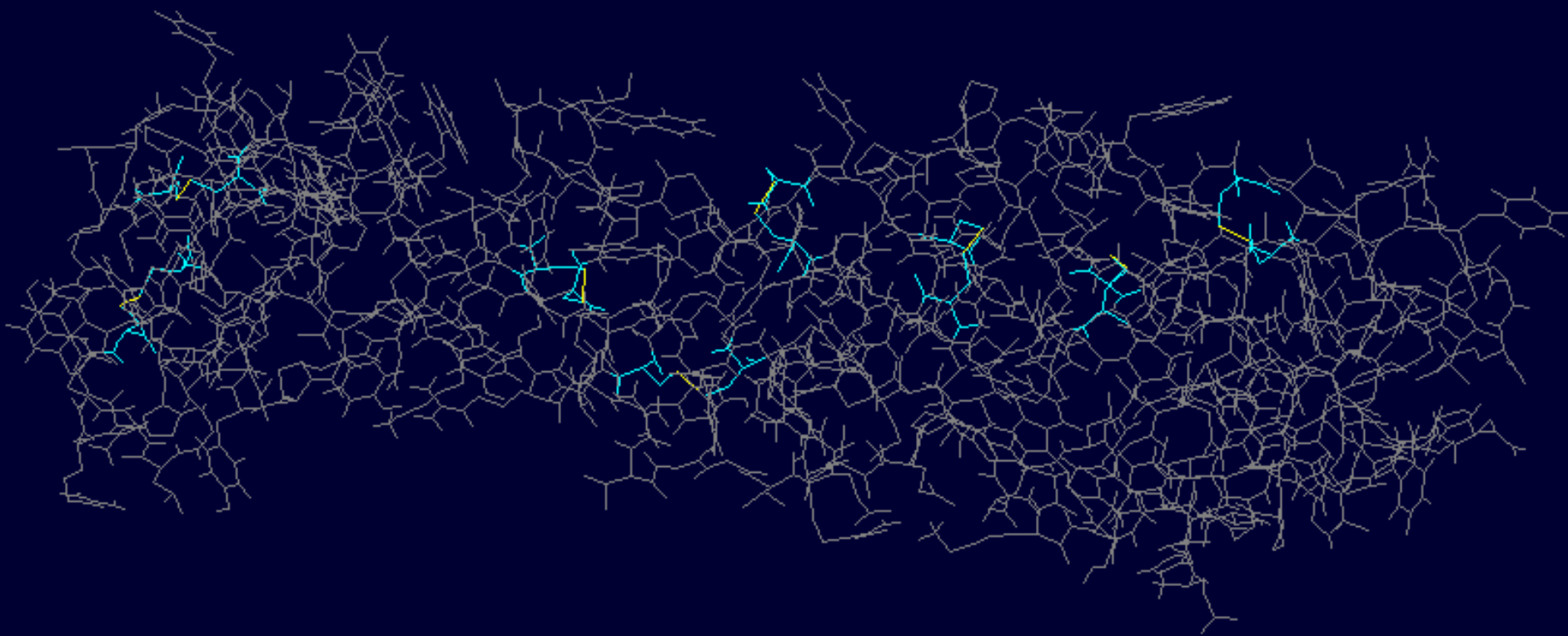
pub/procheck/source3_5

The first contains all the files concatenated into a single tar file (which is the most convenient for unix systems), and the second contains all the individual files

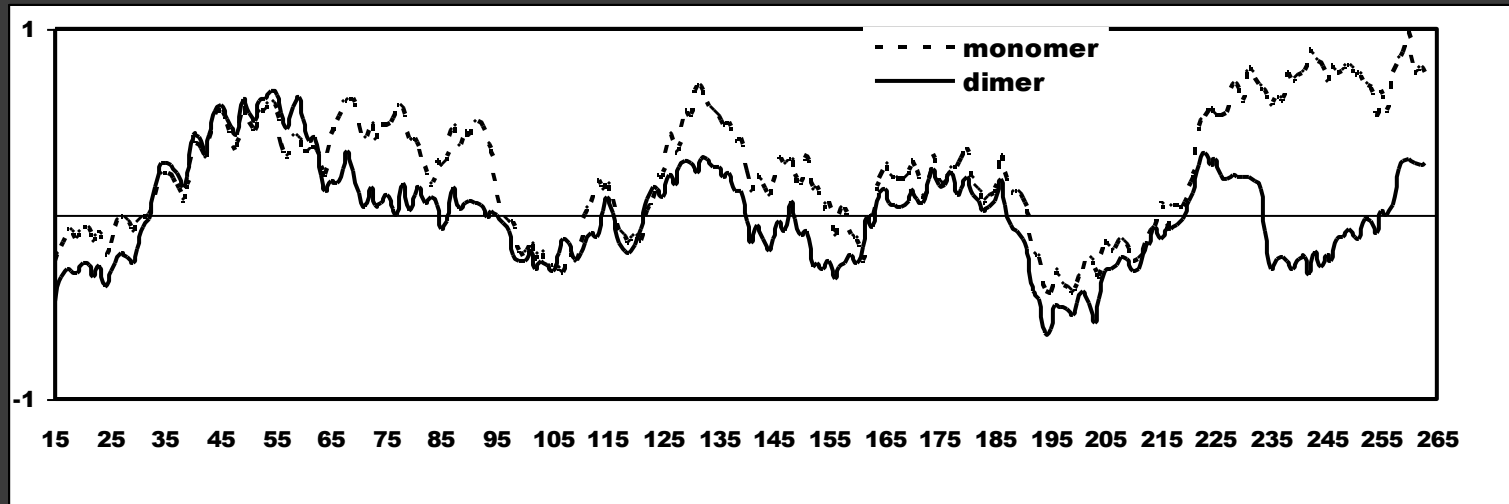




Siti di glicosilazione, siti attivi, siti antigenici....devono essere esposti



PROSA II



Quando l'identità di sequenza fra proteine si trova sotto il 25% un altro metodo di procedere è quello di utilizzare metodi di riconoscimento di fold basati su **profili** e quelli detti di **THREADING**

Metodi basati sui profili:

Ciascun aa ha delle proprietà che possono essere derivate dall'analisi di proteine a struttura nota per ogni aa

Frequenza relativa ad essere presente nelle strutture secondarie note

Frequenza con cui è osservato sulla superficie di una proteina

Frequenza con cui è osservato in un ambiente idrofobico

Ci sono 18 combinazioni 3 per struttura secondaria 3 per l'esposizione 2 per l'ambiente

Per cui l'amminoacido nella sequenza sarà sostituito da una lettera corrispondente ad una delle 18 combinazioni

Se l'operazione viene applicata per tutte le sequenze di struttura nota la banca dati diventerà lineare

Metodi di THREADING=infilare

Questi metodi fanno qualcosa di molto simile a infilare una sequenza nella catena principale delle strutture note. Si ottengono così tanti possibili modelli della proteina usando come template le proteine di struttura nota. I passi successivi sono la valutazione della “bontà” dei modelli approssimati e la selezione di uno o più con punteggio minore da un punto di vista energetico.

ExPASy - Tools - Windows Internet Explorer

http://www.expasy.org/tools/

Cerca web... Traduci la pagina Entra Mail Answers Il Mio Yahoo! Notizie Sport Finanza

ExPASy - Tools

- [Seq2Struct](#) - A web resource for the identification of sequence-structure links
- [STRAP](#) - A structural alignment program for proteins
- [TLSMD](#) - TLS (Translation/Libration/Screw) Motion Determination

Tertiary structure prediction

Comparative modeling

- [SWISS-MODEL](#) - An automated knowledge-based protein modelling server
- [3Djigsaw](#) - Three-dimensional models for proteins based on homologues of known structure
- [CPHmodels](#) - Automated neural-network based protein modelling server
- [ESyPred3D](#) - Automated homology modeling program using neural networks
- [Geno3d](#) - Automatic modelling of protein three-dimensional structure
- [SDSC1](#) - Protein Structure Homology Modeling Server

Threading

- [3D-PSSM](#) - Protein fold recognition using 1D and 3D sequence profiles coupled with secondary structure information (Foldfit)
- [Fugue](#) - Sequence-structure homology recognition
- [HHpred](#) - Protein homology detection and structure prediction by HMM-HMM comparison
- [Libellula](#) - Neural network approach to evaluate fold recognition results
- [LOOPP](#) - Sequence to sequence, sequence to structure, and structure to structure alignment
- [SAM-T02](#) - HMM-based Protein Structure Prediction
- [Threader](#) - Protein fold recognition

- [ProSup](#) - Protein structure superimposition
- [SWEET](#) - Constructing 3D models of saccharides from their sequences

Ab initio

- [HMMSTR/Rosetta](#) - Prediction of protein structure from sequence

Assessing tertiary structure prediction

- [Anolea](#) - Atomic Non-Local Environment Assessment

Internet | Modalità protetta: attivata 100%

MD-@ PC Manager presentazione spiga... GS2.ppt [modalità c... diapositive ExPASy - Tools - Wi... IT 16:39

Choose your system - Windows Internet Explorer

http://www.sbg.bio.ic.ac.uk/~3dpssm/

acer Y! Cerca Traduci la pagina Entra Mail Answers Il Mio Yahoo! Notizie Sport Finanza

Cerca web... Preferiti Spaces

Choose your system

You are advised to use the NEW Phyre fold recognition system. It is more accurate and more up to date than 3D-PSSM

Features in Phyre include:

- Profile-profile matching algorithm (10-15% improved sensitivity over 3dpssm)
- New cleaner interface
- Fully up-to-date fold library

Any comments/criticisms/praise/problems please contact [Lawrence Kelley](#)

phyre OR 3D-pssm

Internet | Modalità protetta: attivata 100%

MD-@ PC Manager presentazione spiga... GS2.ppt [modalità c... diapositive Choose your system... IT 16:36

PHYRE Protein Fold Recognition Server - Windows Internet Explorer

http://www.sbg.bio.ic.ac.uk/~phyre/index.cgi

acer Y! Cerca Traduci la pagina Entra Mail Answers Il Mio Yahoo! Notizie Sport Finanza

Cerca web... Preferiti Spaces Strumenti

phyre

Version 0.2

Protein Homology/analogY Recognition Engine

The Phyre webserver is for **Academic use only**
For in-house and/or commercial use please click [here](#)

Note: [Other tools available from our lab \(function prediction, docking, etc.\)](#)

E-mail Address

Optional Job description

Amino Acid Sequence

Quick Phyre Search

ESEMPIO TIPICO DI RISPOSTA DI PHYRE

QuickPhyre results for job Globin_Example - Windows Internet Explorer

phyre
Imperial College
London

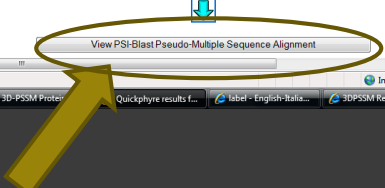
[Phyre Server: 0.1][Fold Library (SCOP): 1.67][Fold Library (PDB): 20050629][NR sequence DB: 20050626][Dynamic FR engine: 1.0]

QuickPhyre Results for Job *Globin_Example*

Email: i.a.kelley@imperial.ac.uk
Job Code: 86cc0c47eba655e6
Description: Globin_Example
Date: Tue Jul 12 12:52:26 BST 2005

[Renew] your results for 6 days
Download a tarred gzipped version of these results

[View PSI-Blast Pseudo-Multiple Sequence Alignment](#)



ALLINEAMENTI POSSIBILI
FORNITI DA PSI-BLAST

Identities computed with respect to: (query) Globin_Example
Colored by: property

HSP processing: ranked
Search cycle: 3
[Click here for FASTA Format Flat File](#)

	Globin_Example	bits	E-value	N	100.0%
1	gi 55635223 ref XP_508244.1 	PREDICTED: similar to HOR5B...	131	3e-30	1 22.4%
2	gi 183885 gb AA050159.1 	gamma-globin	127	5e-29	1 22.4%
3	gi 50731454 ref XP_425673.1 	PREDICTED: similar to Hemog...	125	2e-28	1 23.1%
4	gi 55635481 ref XP_521773.1 	PREDICTED: similar to delta...	124	3e-28	1 22.5%
5	gi 1970 emb CAA25986.1 	epsilon I globin [Capra hir...	124	4e-28	1 24.3%
6	gi 3808 emb CAA37498.1 	gamma-globin-2 (A) [Macaca m...	124	4e-28	1 22.1%
7	gi 6016195 sp Q95239 HBE_PROVE	Hemoglobin epsilon chain >g...	124	6e-28	1 22.1%
8	gi 122722 sp P1495 HBE_TAREV	Hemoglobin epsilon chain >g...	124	6e-28	1 22.9%
9	gi 170812 sp P51438 HBE_ATEBE	Hemoglobin epsilon chain >g...	124	6e-28	1 23.6%
10	gi 122724 sp P11025 HBE_DIDMA	Hemoglobin epsilon-M chain ...	124	6e-28	1 22.1%
11	gi 54037269 sp P68027 HBE_CEBPY	Hemoglobin epsilon chain >g...	123	7e-28	1 22.9%
12	gi 170813 sp P51443 HBE_SAGMI	Hemoglobin epsilon chain >g...	123	8e-28	1 22.9%
13	gi 122731 sp P02103 HBE_RABIT	Hemoglobin epsilon chain (B...	123	8e-28	1 22.9%
14	gi 728697 emb CAA88563.1 	hemoglobin embryonic beta-c...	123	8e-28	1 22.1%
15	gi 52138683 ref NP_001004390.1 	rho-globin [Gallus gallus] ...	123	9e-28	1 22.9%
16	gi 48429238 sp P61948 HBG2_HYLIA	Hemoglobin gamma-2 chain (H...	123	9e-28	1 22.1%
17	gi 47523946 ref NP_999612.1 	hemoglobin, epsilon 1 [Sus ...	123	9e-28	1 25.0%
18	gi 122530 sp P10061 HBB2_SPHPU	Hemoglobin beta-2 chain	123	9e-28	1 23.0%
19	gi 62281955 sp P68079 HBG_PAPCY	Hemoglobin gamma chain >g ...	123	1e-27	1 22.1%
20	gi 49169701 ref NP_990820.1 	beta-globin [Callus gallus] ...	123	1e-27	1 22.9%
21	gi 62440984 ref XP_579276.1 	PREDICTED: hemoglobin, epsi...	123	1e-27	1 25.7%
22	gi 342823 gb AA036930.1 	gamma-1 globin [Pongo pygma...	123	1e-27	1 21.4%
23	gi 54037267 sp P68025 HBE_CHISA	Hemoglobin epsilon chain >g...	123	1e-27	1 22.9%

PROBABILI FUNZIONI DELLA QUERY

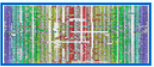
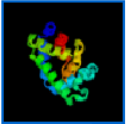
Quickphyre results for job Globin_Example - Windows Internet Explorer

http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/86cc0c47eba655e6/summary.html

Prosite

Functional Keywords	Weight
Erythrocyte	80
Oxygen transport	77
Heme	65
Embryo	62
Iron	61
Metal-binding	57
Transport	57
Polymorphism	42
Acetylation	41

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family
	d3adha (length:145) 100% i.d.	 Download MOL	9.3e-20	100 %	0.90 Biotech	Globin-like	Globin-like	Globins

Internet | Modalità protetta: attivata | 100%

MD-@ PC Manager | 3D-PSSM Protein Fo... | Quickphyre results f... | label - English-Italia... | 3DPSSM Results for ... | Microsoft PowerPoi... | IT | 19:18

GLI ALLINEAMENTI CON MAGGIORE PUNTEGGIO

Quickphyre results for job Globin_Example - Windows Internet Explorer


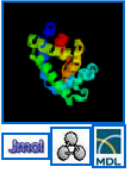

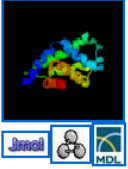

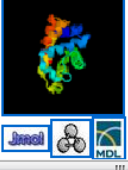
http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/86cc0c47eba655e6/summary.html

acer Y! Traduci la pagina Entra Mail Answers Il Mio Yahoo! Notizie Sport Finanza

Quickphyre results for job Globin_Example

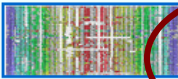
Acetylation 41

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family
	d3sdha (length:145) 100% i.d.		9.3e-20	100 %	0.90 Biotext	Globin-like	Globin-like	Globins
	c2bk9A (length:153) 23% i.d.		7.7e-17	100 %	0.89 Biotext	PDB header:oxygen transport	Chain: A: PDB Molecule:cg9734-pa;	PDBTitle: drosophila melanogaster globin
	d1itha (length:141) 16% i.d.		2.1e-16	100 %	0.88 Biotext	Globin-like	Globin-like	Globins


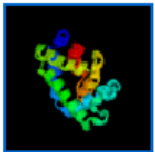





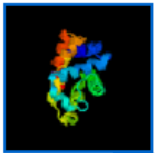

Internet | Modalità protetta: attivata 100%

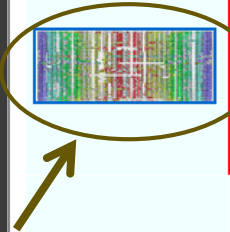
MD-@ PC Manager 3D-PSSM Protein Fo... Quickphyre results f... label - English-Italia... 3DPSSM Results for ... Microsoft PowerPoi... IT 19:19

	d1sc1b (length:150) 55% i.d.		1.8e-15	100 %	0.89 Biotext	Globin-like	Globin-like	Globins
	d1g08a (length:141) 19% i.d.		2.5e-15	100 %	0.94 Biotext	Globin-like	Globin-like	Globins
	d1qpwa (length:141) 18% i.d.		2.8e-15	100 %	0.94 Biotext	Globin-like	Globin-like	Globins
	d2gdm (length:153) 13% i.d.		2.8e-15	100 %	0.84 Biotext	Globin-like	Globin-like	Globins

Acetylation 41

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family
	d3sdha (length:145) 100% i.d.	 	9.3e-20	100 %	0.90 Biotext	Globin-like	Globin-like	Globins
	c2bk9A (length:153) 23% i.d.	 	7.7e-17	100 %	0.89 Biotext	PDB header:oxygen transport	Chain: A: PDB Molecule:cg9734-pa;	PDBTitle: drosophila melanogaster globin
	d1itha (length:141) 16% i.d.	 	2.1e-16	100 %	0.88 Biotext	Globin-like	Globin-like	Globins




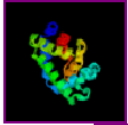


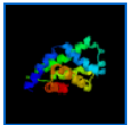


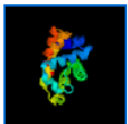

CODICE SCOP = Structural Classification Of Proteins

Quickphyre results for job Globin_Example - Windows Internet Explorer

http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/86cc0c47eba655e6/summary.html

Acetylation 41

Fold Recognition

View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family
	d3adha (length:145) 100% i.d.	 	9.3e-20	100 %	0.90 Biotext	Globin-like	Globin-like	Globins
	c2bk9A (length:153) 23% i.d.	 	7.7e-17	100 %	0.89 Biotext	PDB header: oxygen transport	Chain: A: PDB Molecule: cg9734-pa;	PDBTitle: drosophila melanogaster globin
	d1itha (length:141) 16% i.d.	 	2.1e-16	100 %	0.88 Biotext	Globin-like	Globin-like	Globins

Internet | Modalità protetta: attivata 100%

MD-@ PC Manager Quickphyre results f... Quickphyre alignme... Microsoft PowerPoi...

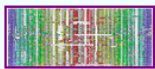

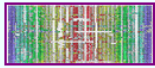
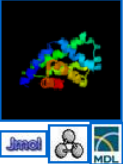

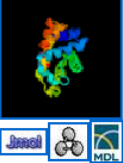
View Model

Quickphyre results for job Globin_Example - Windows Internet Explorer

http://www.sbg.bio.ic.ac.uk/phyre/qphyre_output/86cc0c47eba655e6/summary.html

Acetylation 41

Fold Recognition

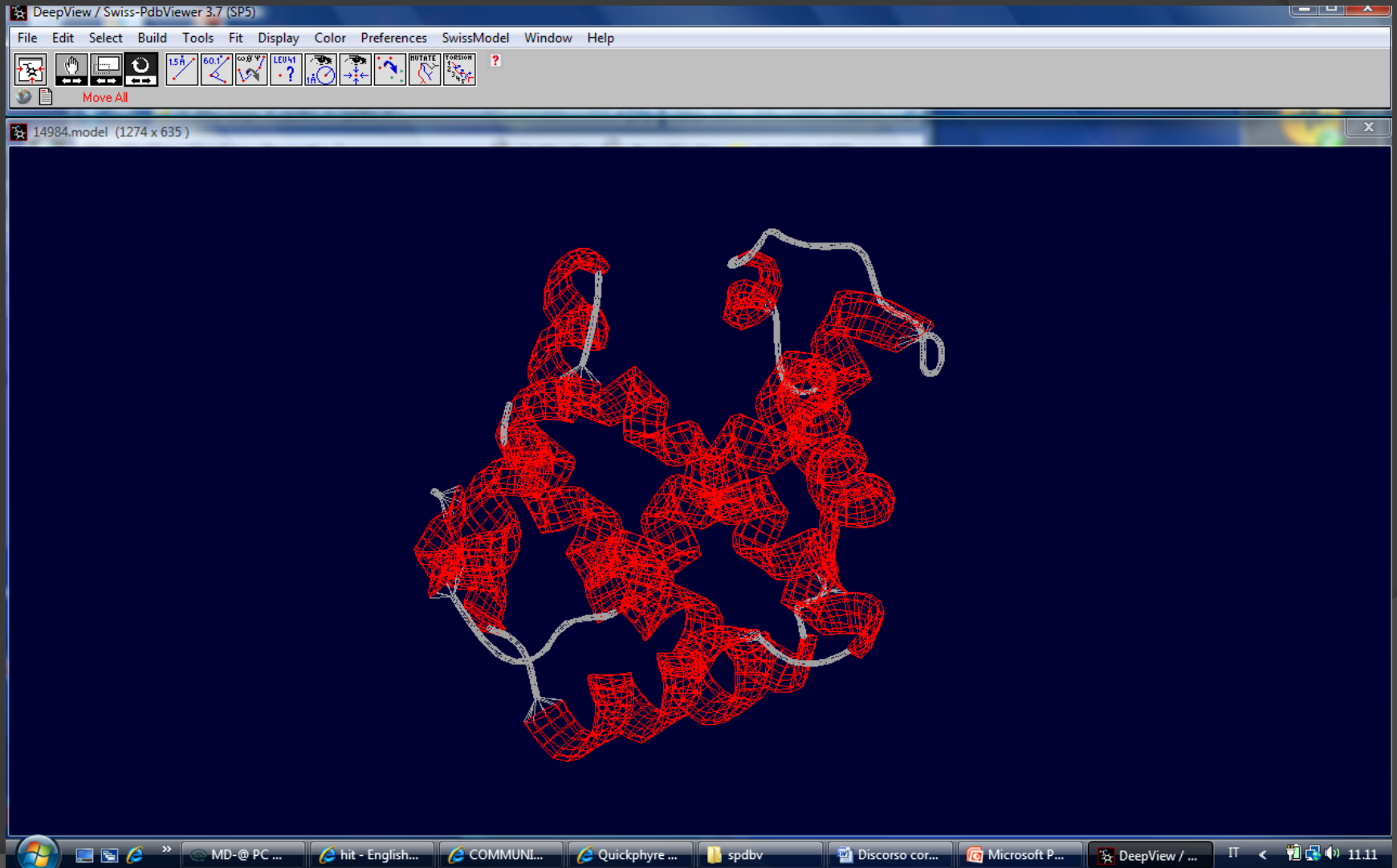
View Alignments	SCOP Code	View Model	E-value	Estimated Precision	BioText	Fold/PDB descriptor	Superfamily	Family
	d3sdha (length:145) 100% i.d.		9.3e-20	100 %	0.90 BioText	Globin-like	Globin-like	Globins
	c2bk9A (length:153) 23% i.d.		7.7e-17	100 %	0.89 BioText	PDB header:oxygen transport	Chain: A: PDB Molecule:cg9734-pa;	PDBTitle: drosophila melanogaster globin
	d1itha (length:141) 16% i.d.		2.1e-16	100 %	0.88 BioText	Globin-like	Globin-like	Globins

Internet | Modalità protetta: attivata

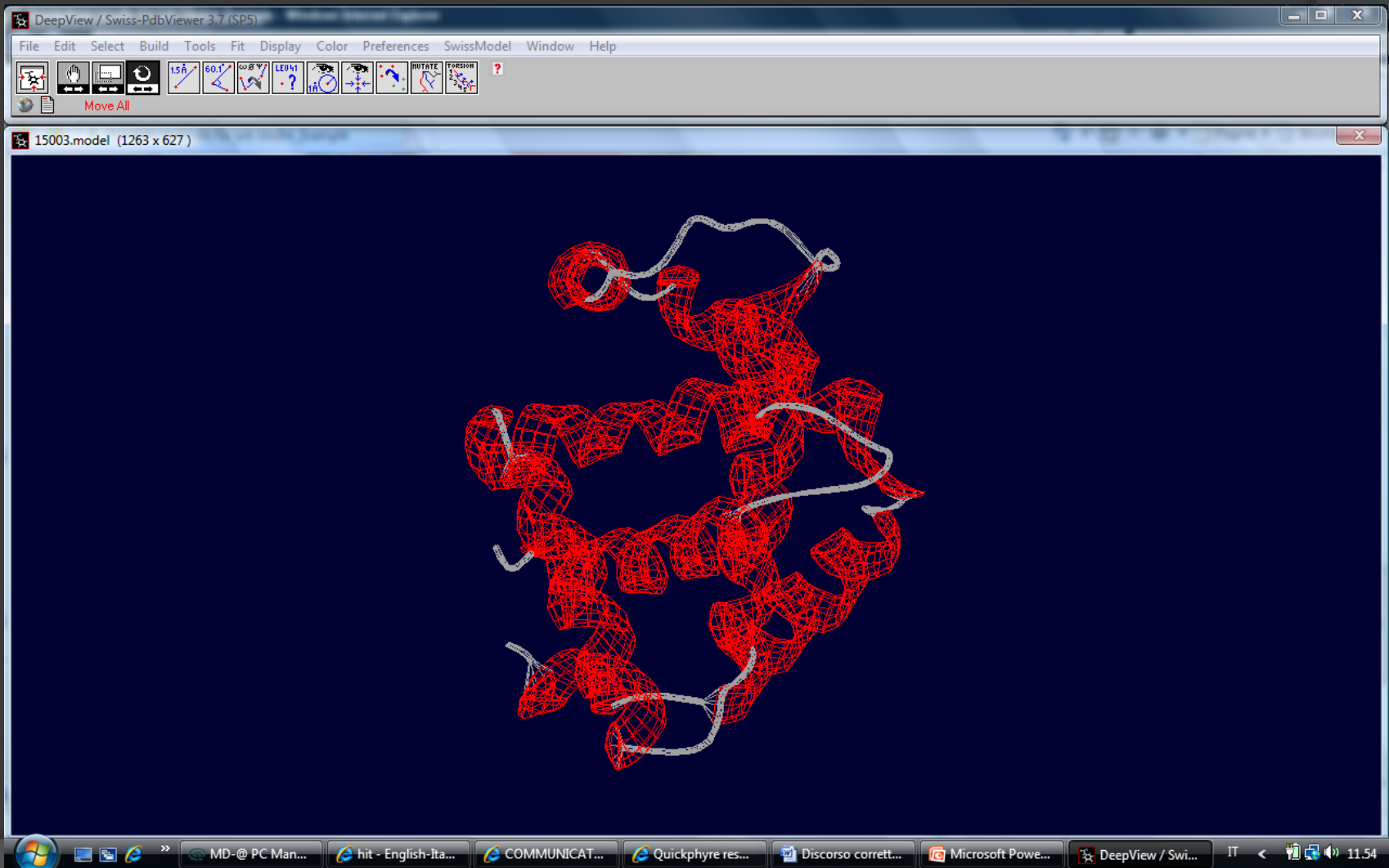
100%

IT < > 0.24

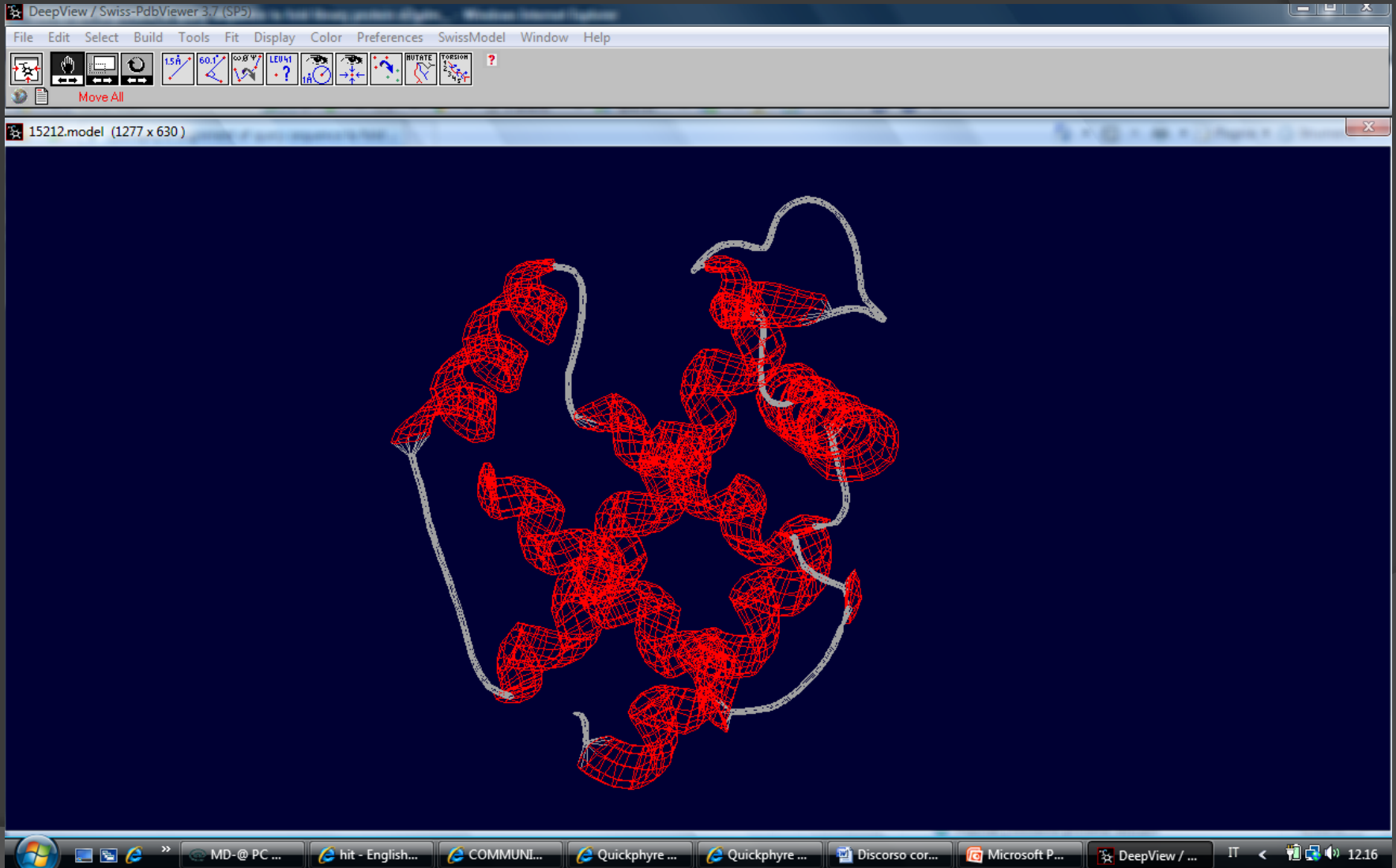
MODELLO DELLA d3sdh_a



MODELLO OTTENUTO USANDO COME TEMPLATE d1sct_b



MODELLO OTTENUTO USANDO COME TEMPLATE d2gdm



Il benchmark di un metodo predittivo

Farlo da soli utilizzando i vari metodi richiederebbe molto tempo e spazio di calcolo.

EVA è un server che ogni giorno scarica dal PDB le nuove strutture ne estrapola la sequenza e prova usando i vari metodi predittivi ad ottenere una struttura usando similitudine di sequenza non superiore al 33% solo automatic servers

CASP (critical assessment techniques for protein structure prediction)

SARS 1

Comparsa nel 2003 di un nuovo virus

?

Sequenziamento del genoma completo

NC_004718

Identificazione della classe di appartenenza

Coronavirus ssRNA

Identificazione delle proteine espresse

14 proteine

Identificazione delle proteine responsabili dell'attacco alla cellula ospite

Spike protein

Identificazione della struttura delle proteine di interesse

S1 + S2

Cosa si può fare con la bioinformatica

Identificazione di una nuova specie

Sars CoVirus

Determinare in tempi brevi la struttura di alcune proteine

Spike protein
Peplomero
Elicasi



Drug design
Diagnostic kits
Vaccine

Sequenza nucleotidica



Sequenza amminoacidica



Metodi bioinformatici di predizione strutturale



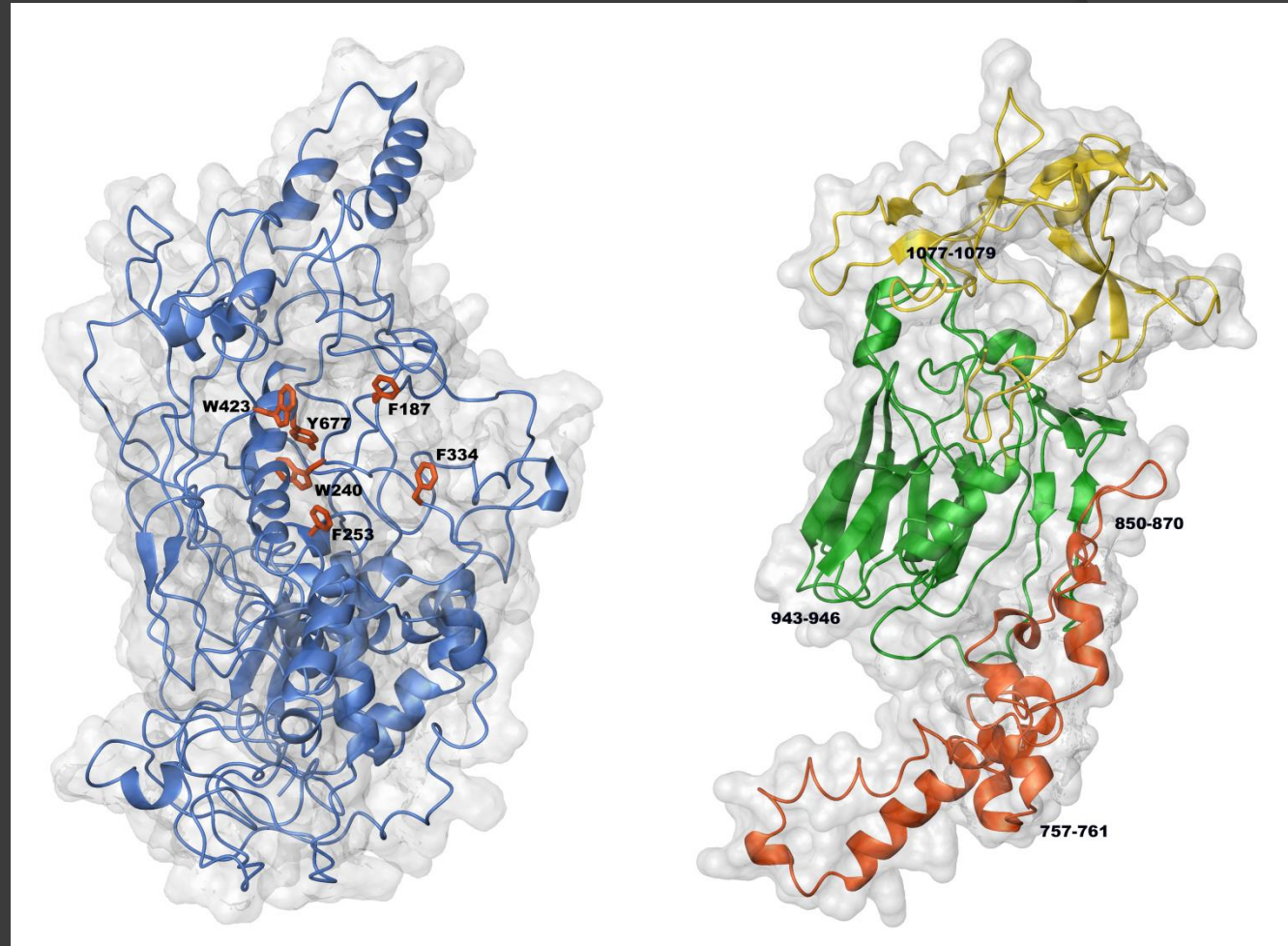
Modelling



Struttura Tridimensionale

1G9D neurotoxin B of *Clostridium botulinum*

15% identity
49% positivity
5% gaps



S1 1Q4Z

S2 1Q4Y

SARS 2

SEM Scanning Electron Microscopy

Struttura
quaternaria

Trimerico della spike protein sulla superficie

Ricostruzione del trimerico con approccio bioinformatico

Struttura di S1 ed S2 precedentemente modellate

Docking molecolare tra le due subunità

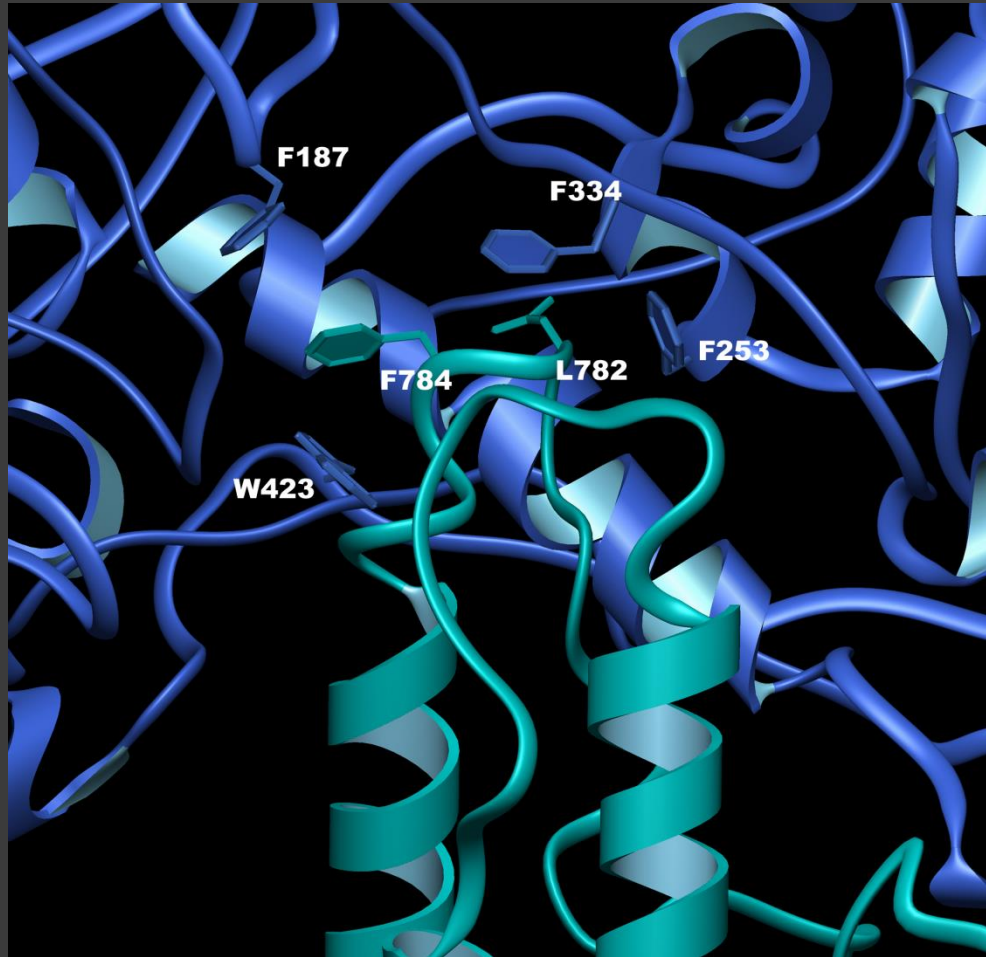
S1 testa globulare del peplomero con attività di legame al recettore cellulare

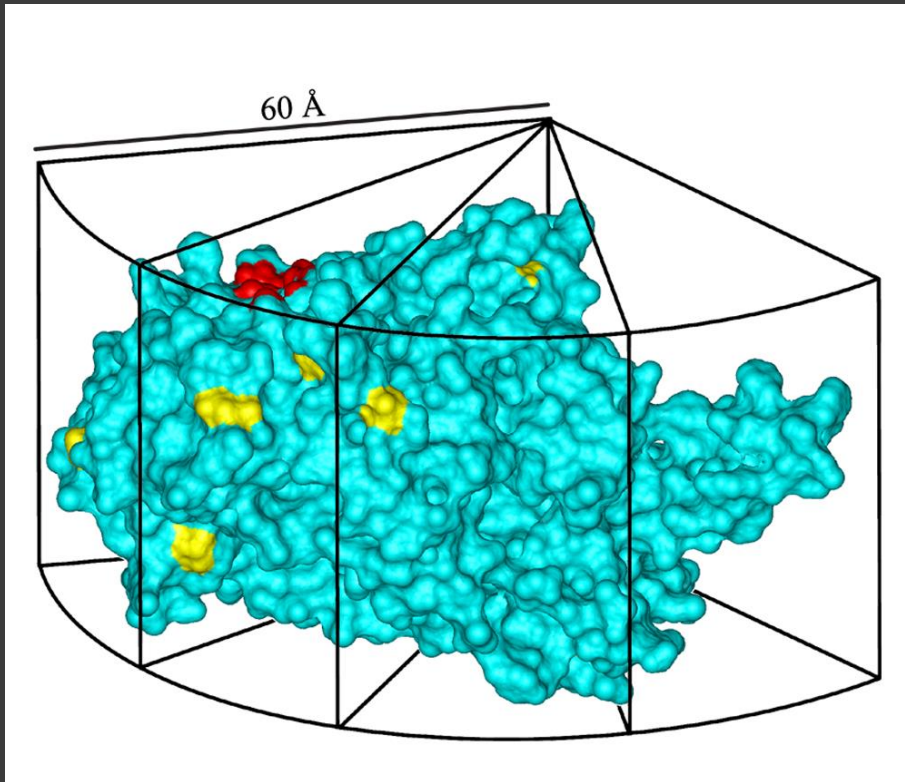
S2 gambo del peplomero

1 step: modelling del gambo di S2 coiled coil e HRregion

2 step: allineamento di 253 sequenze di S1
alcuni residui sono sempre conservati

3 step: tasca idrofobica necessaria per l'assemblaggio del peplomero



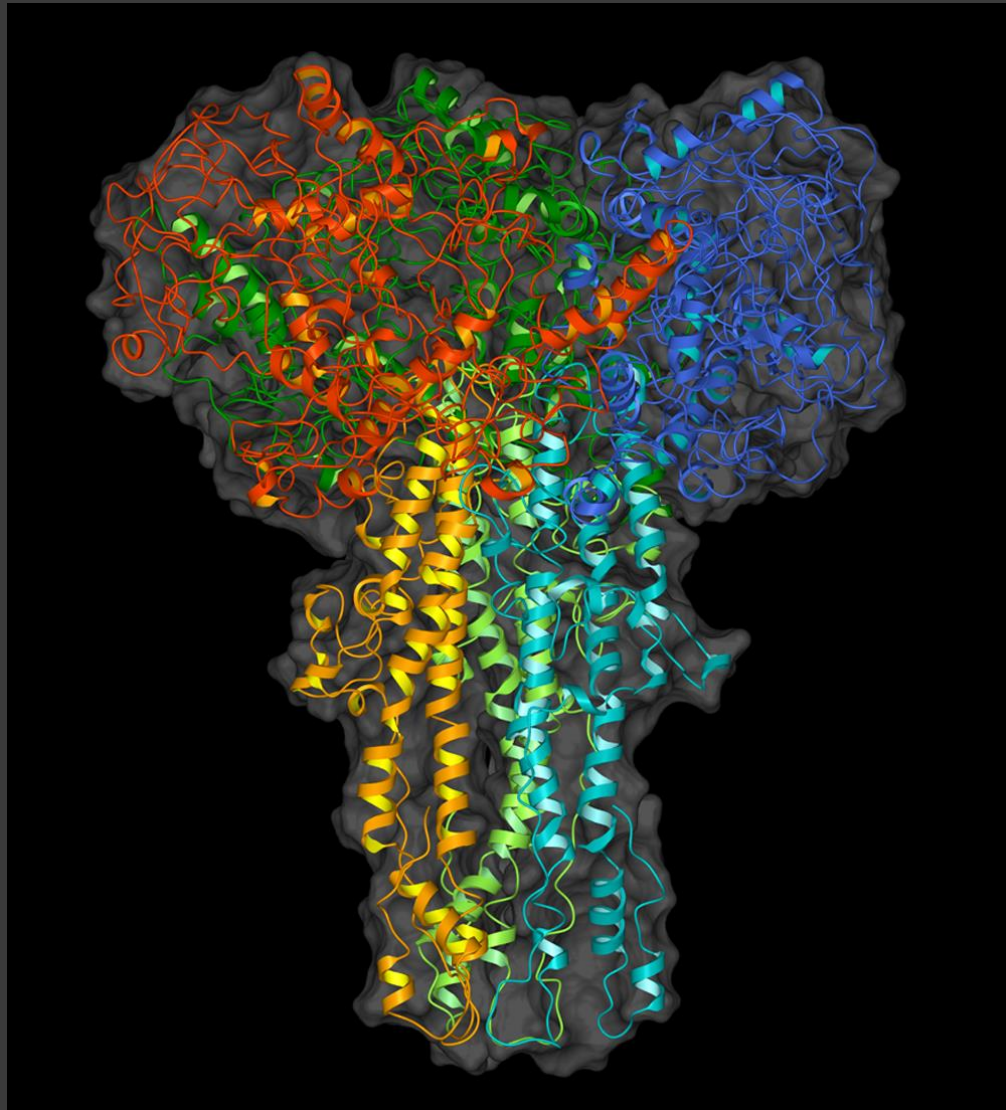


Ricostruzione del
peplomero

Orientazione di
S1+S2

Siti glicosilazione
esterni

Siti idrofobici
interni



Peplomer 1T7G

Brevettato per potenziali farmaci

Confronto con altre elicasi che presentavano un MBD

Arterivirus helicase possiedono un MBD complessante
4 atomi di zinco

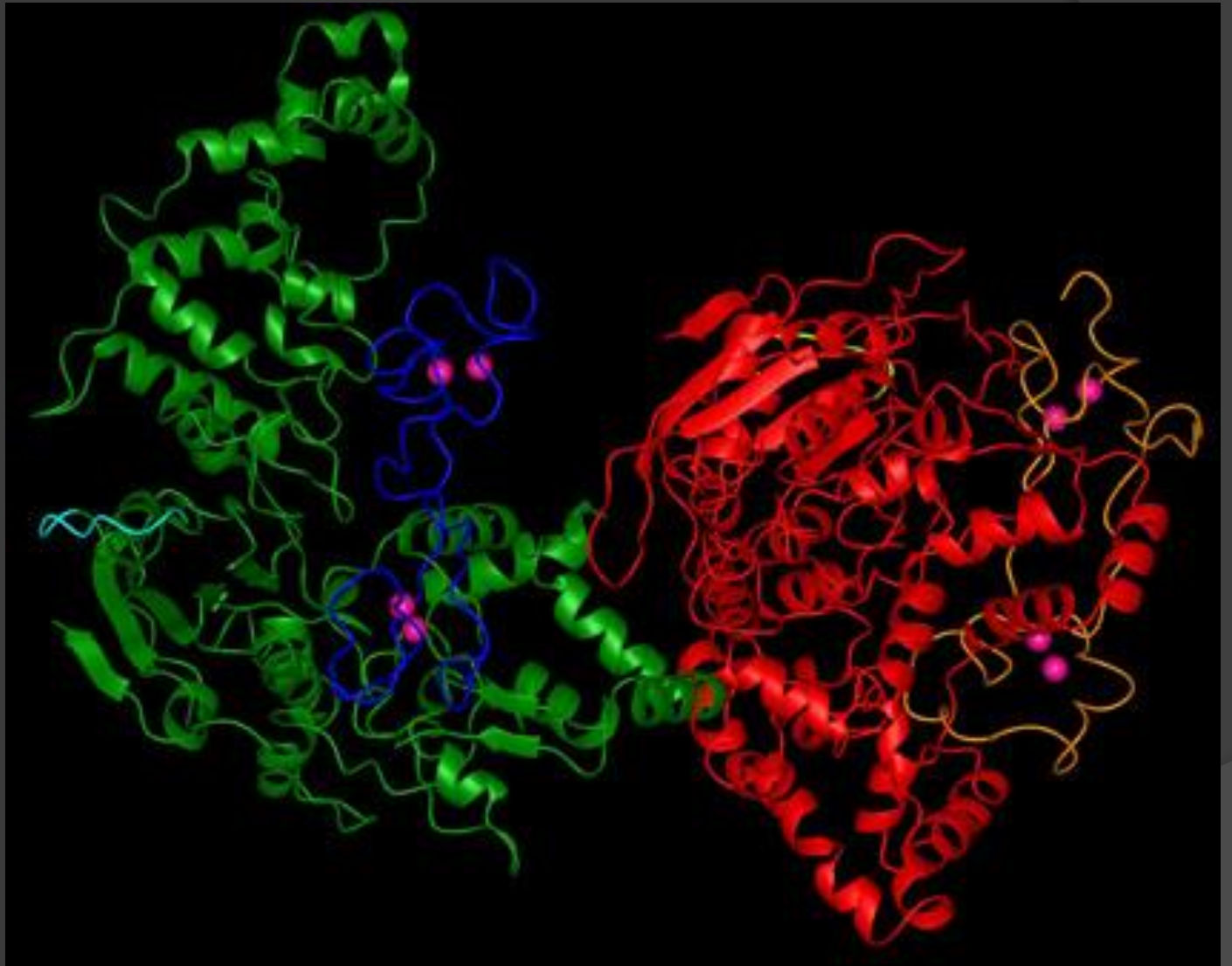
Possibili modelli presenti in altre proteine

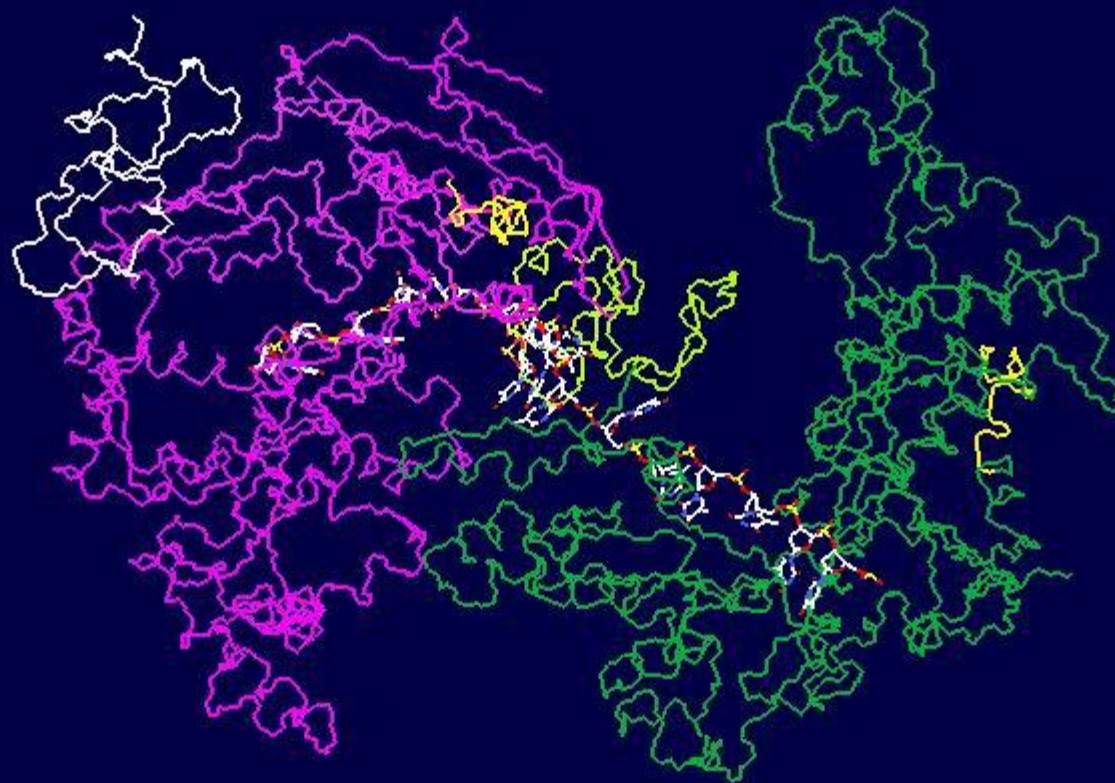
Zn₂Cys₆ GAL4-like protein

Zn₂Cys₄His₂ RAG1 domain

Ricostruzione della struttura geometrica tetraedrica







ALBUMINA

Albumina umana > cristallo

Albumina bovina > modello

Albumina ratto > modello

Albumina suino > modello

Valutazione sperimentale della loro reattività nei confronti di diversi tioli

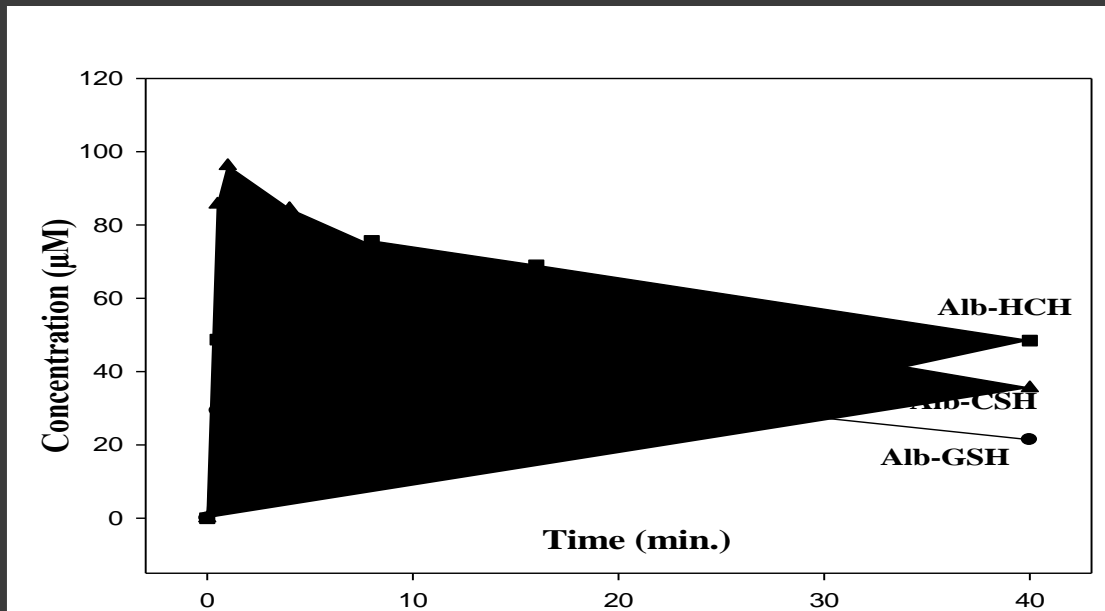
Sequenze molto simili spiegazioni attraverso le strutture

Valutazione sperimentale della loro reattività con diversi tioli:

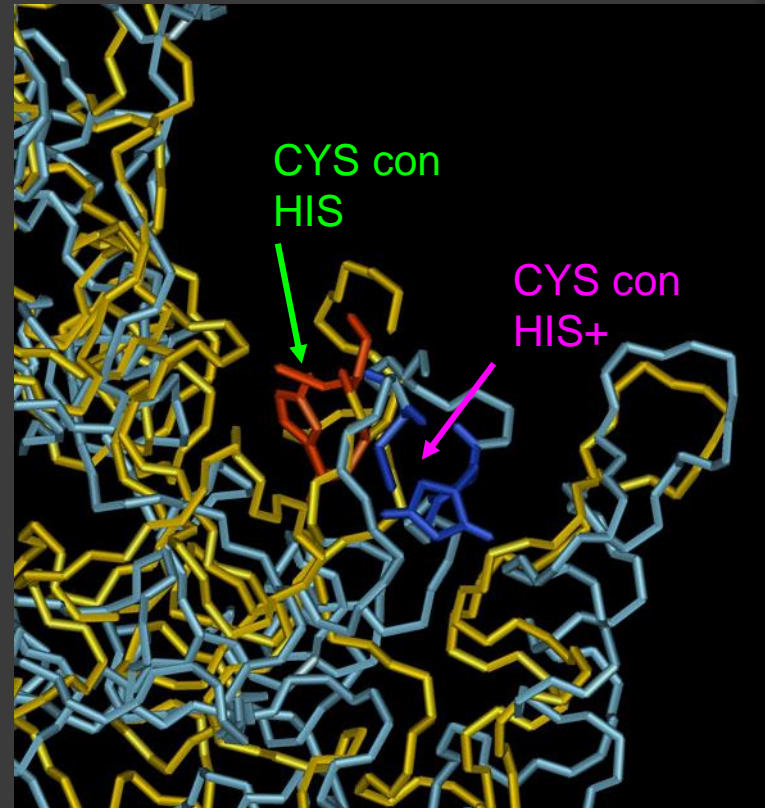
Glutathione

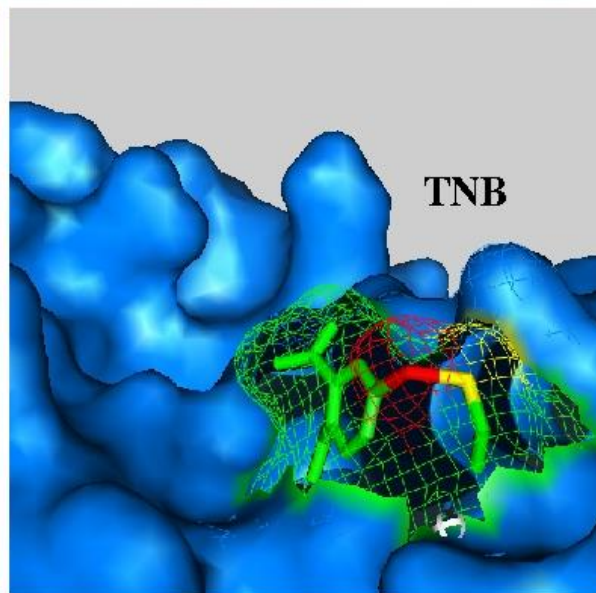
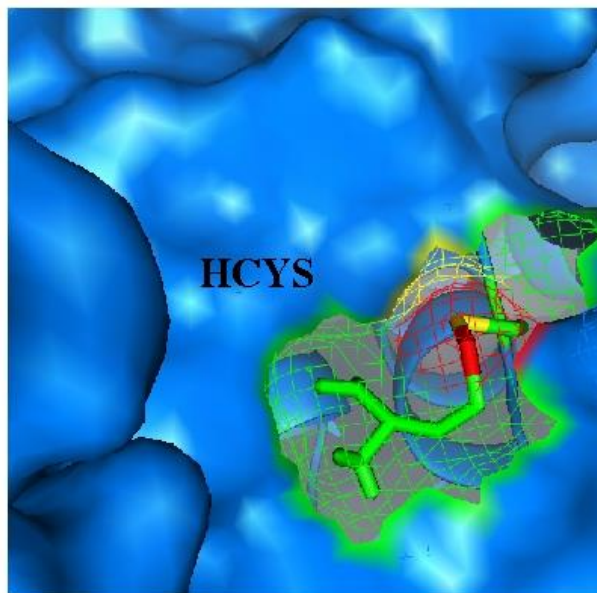
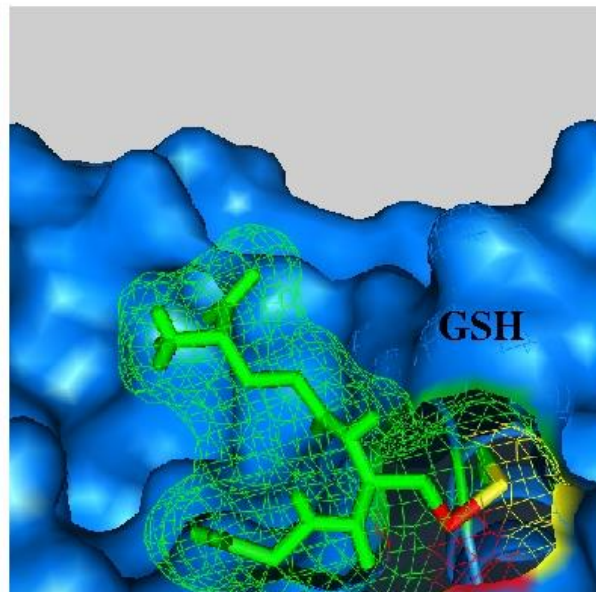
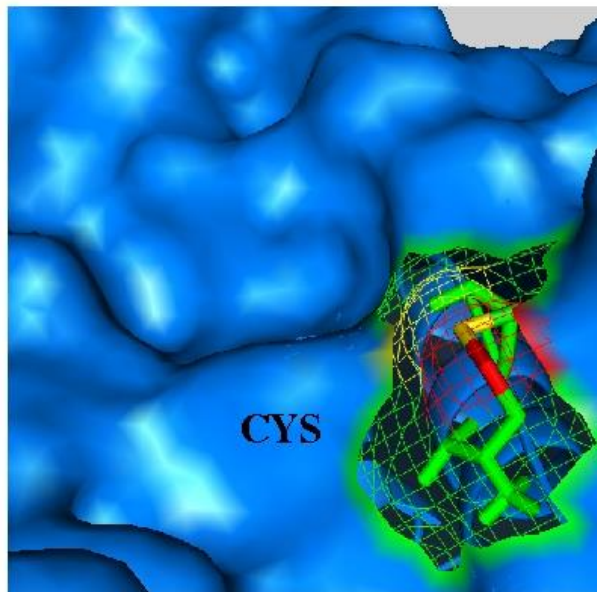
Cistina

Omocisteina



Variazioni del pH influenzano
l'esposizione della cisteina
Flip-Flop





Quali sono le differenze strutturali che spiegano il differente comportamento delle 4 albumine?

```

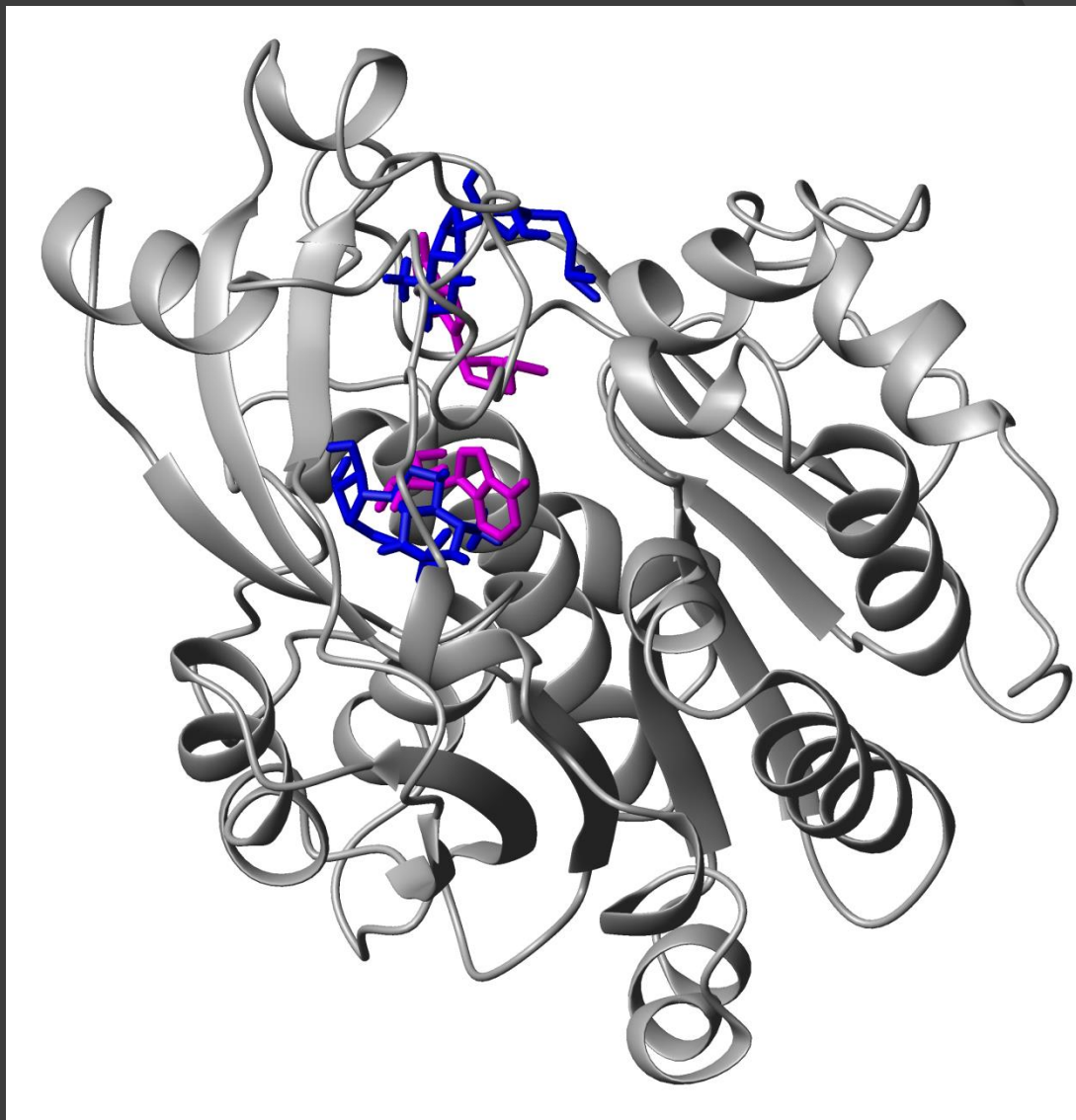
                10         20         30         40         50         60
                |         |         |         |         |         |
ssa_his  WVTFISLLFLFSSAYSRGVFRRDYKSEIAHRFKDLGEQYFKGLVLIAFSQHLQOCPYEE
bsa_his  -----SEIAHRFKDLGEEHFKGLVLIAFSQYLQOCPFDE
hsa_his  -----SEVAHRFKDLGEEHFKGLVLIAFSQYLQOCPFDE
rsa_his  -----SEIAHRFKDLGEQHFKGLVLIAFSQYLQKCPYEE
                *:*****: **:*****:***:***:
Prim.cons.  WVTFISLLFLFSSAYSRGVFRRDYKSEIAHRFKDLGE2HFKGLVLIAFSQYLQOCP2EE

                70         80         90         100        110        120
                |         |         |         |         |         |
ssa_his  HVKLVREVTTEFAKTCVADESAENCDKSIHTLFGDKLCAIPSLREHYGDLADCCKEKEEPER
bsa_his  HVKLVNELTEFAKTCVADESHAGCEKSLHTLFGDELCKVASLRETYGDMADCCKEQEPER
hsa_his  HVKLVNEVTTEFAKTCVADESAENCDKSLHTLFGDKLCTVATLRETYGEMADCCAQKQEPER
rsa_his  HIKLVQEVTDFAKTCVADENAENCDKSIHTLFGDKLCAIPKLRDNYGELADCCAQKQEPER
                *:***.***:*****. .:***:*****:** :..** : **:*:**** *:****
Prim.cons.  HVKLVNEVTTEFAKTCVADESAENCDKS2HTLFGDKLCA22SLRETYG22ADCC2KQEPER

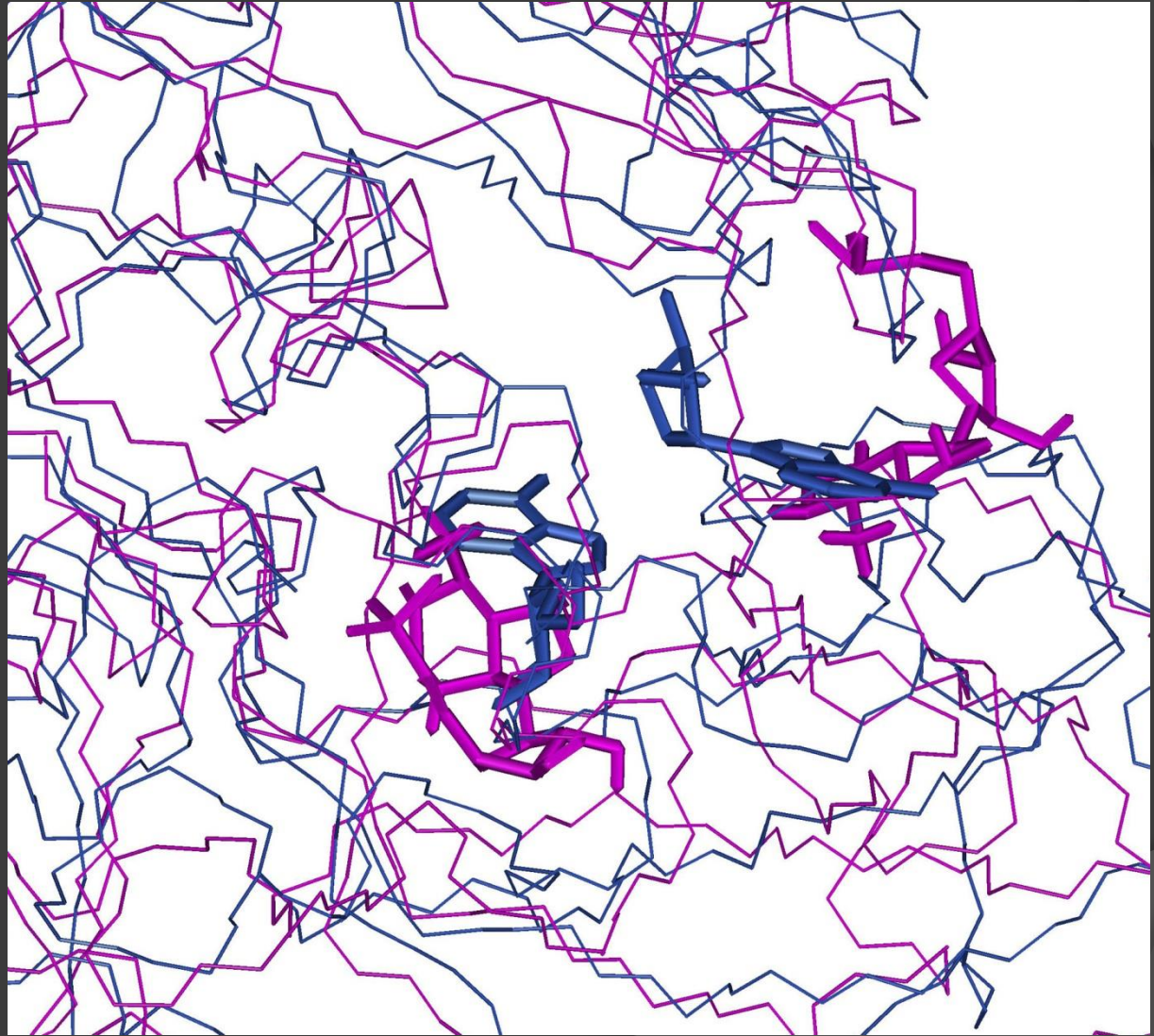
                130        140        150        160        170        180
                |         |         |         |         |         |
ssa_his  NECFLQHKNDNPDIPKLK-PDPVALCADFQEDEQKFWGKYLVEIARRHPYFYAPELLYYA
bsa_his  NECFLSHKDDSPDLPKLK-PDPNTLCDEFKADEKQKFWGKYLVEIARRHPYFYAPELLYYA
hsa_his  NECFLQHKDDNPNLPRLVPRPEVDVMCTAFHDNEETFLKKYLVEIARRHPYFYAPELLFFA
rsa_his  NECFLQHKDDNPNLPPFQRPEAEAMCTSFQENPTSFLGHYLVHEVARRHPYFYAPELLYYA
                *****:***:***:***: : *: .:* *: : .* :**:*:*****:*****:
Prim.cons.  NECFLQHKDDNP2LPKLKRP2P4A2CT4FQE2E4KF2GKYLVEIARRHPYFYAPELLYYA

```

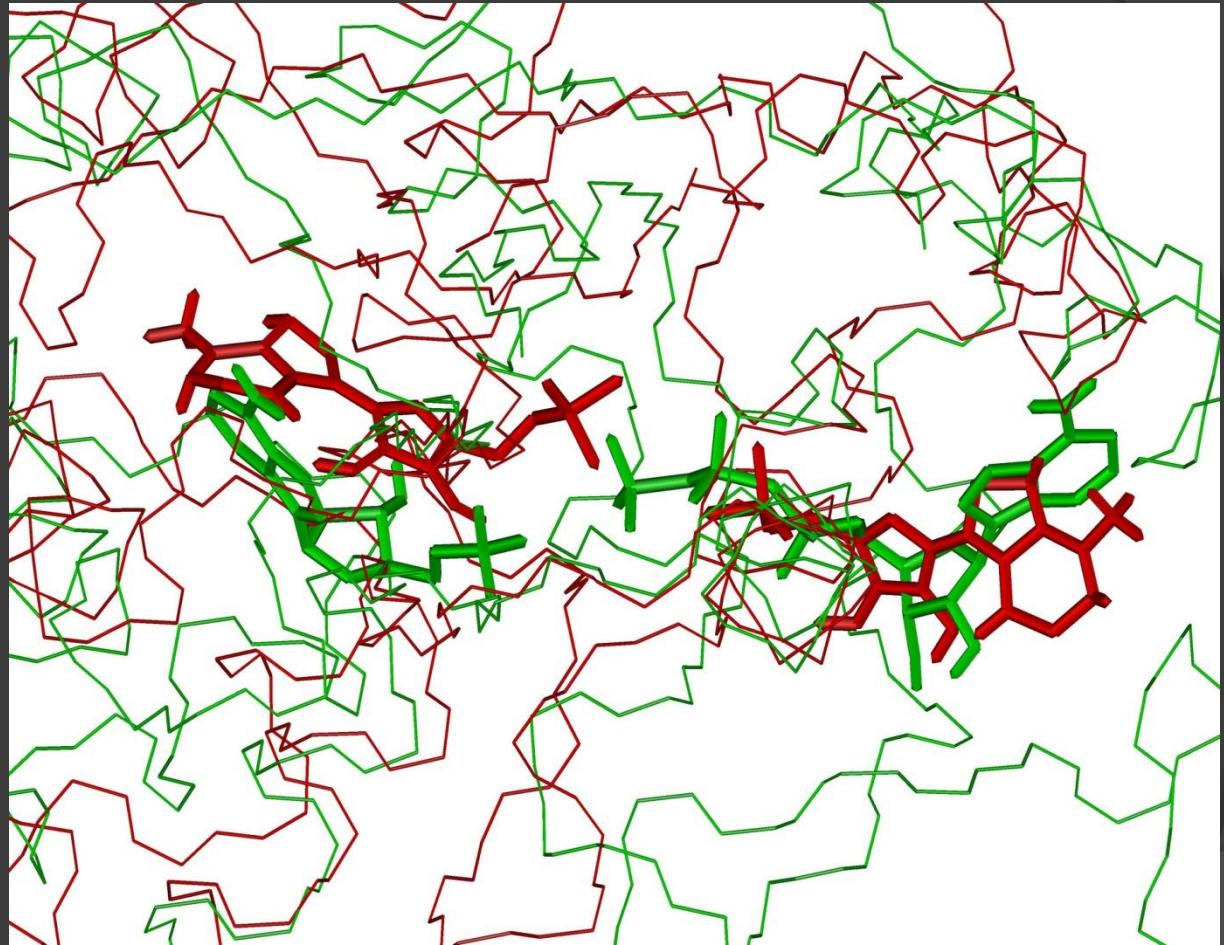

Adenosina Kinasi



Adenosina
Prima e dopo
la dinamica



Adenilato
Prima e dopo
la dinamica

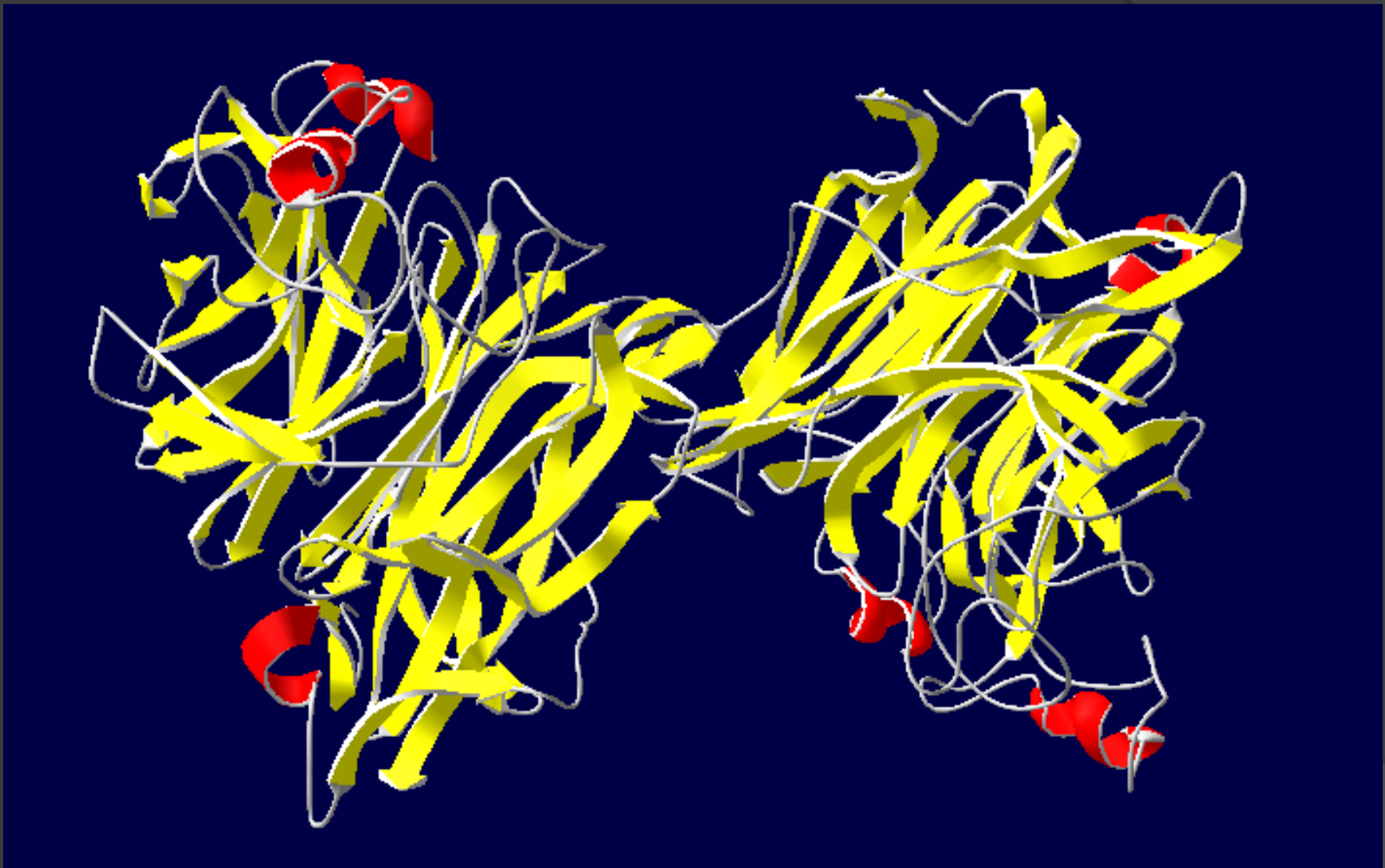




PIVIII

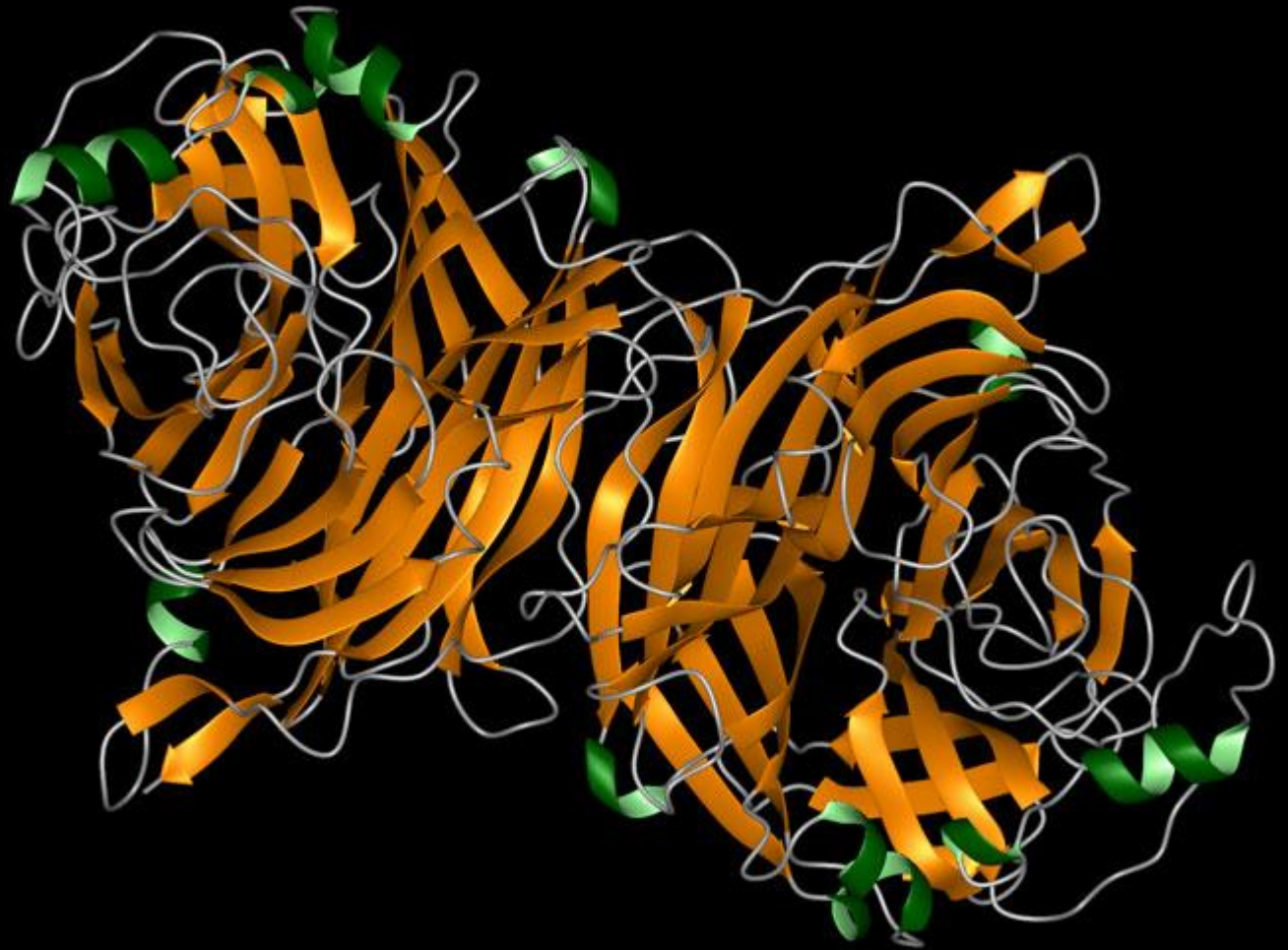
Monomero di HN



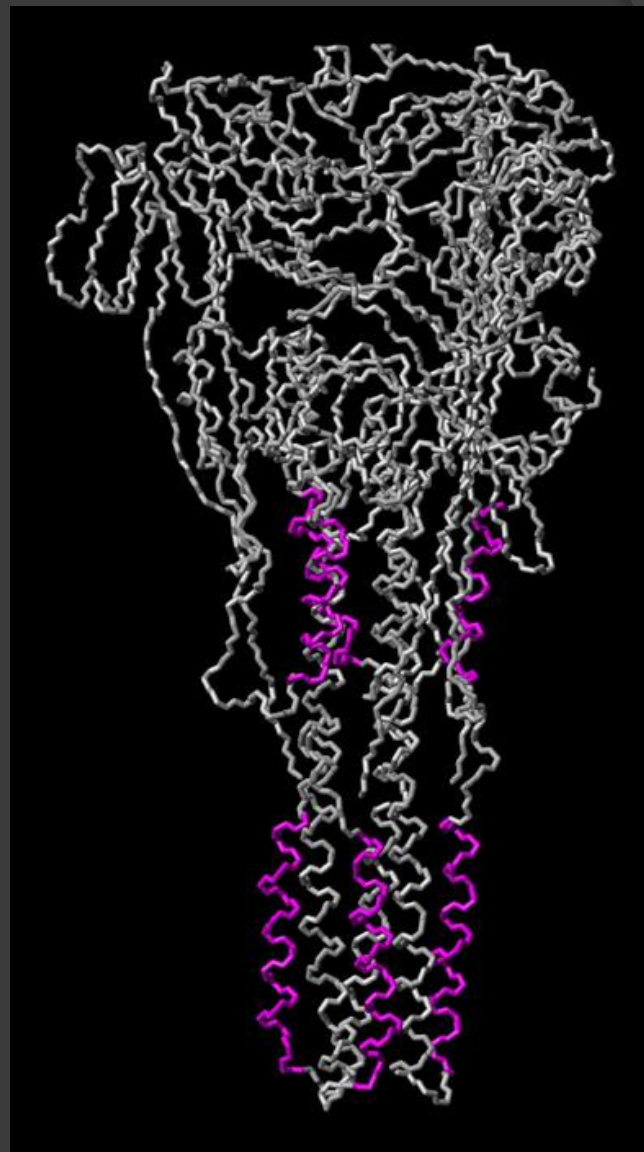


Dimero di HN prima del legame con il recettore cellulare

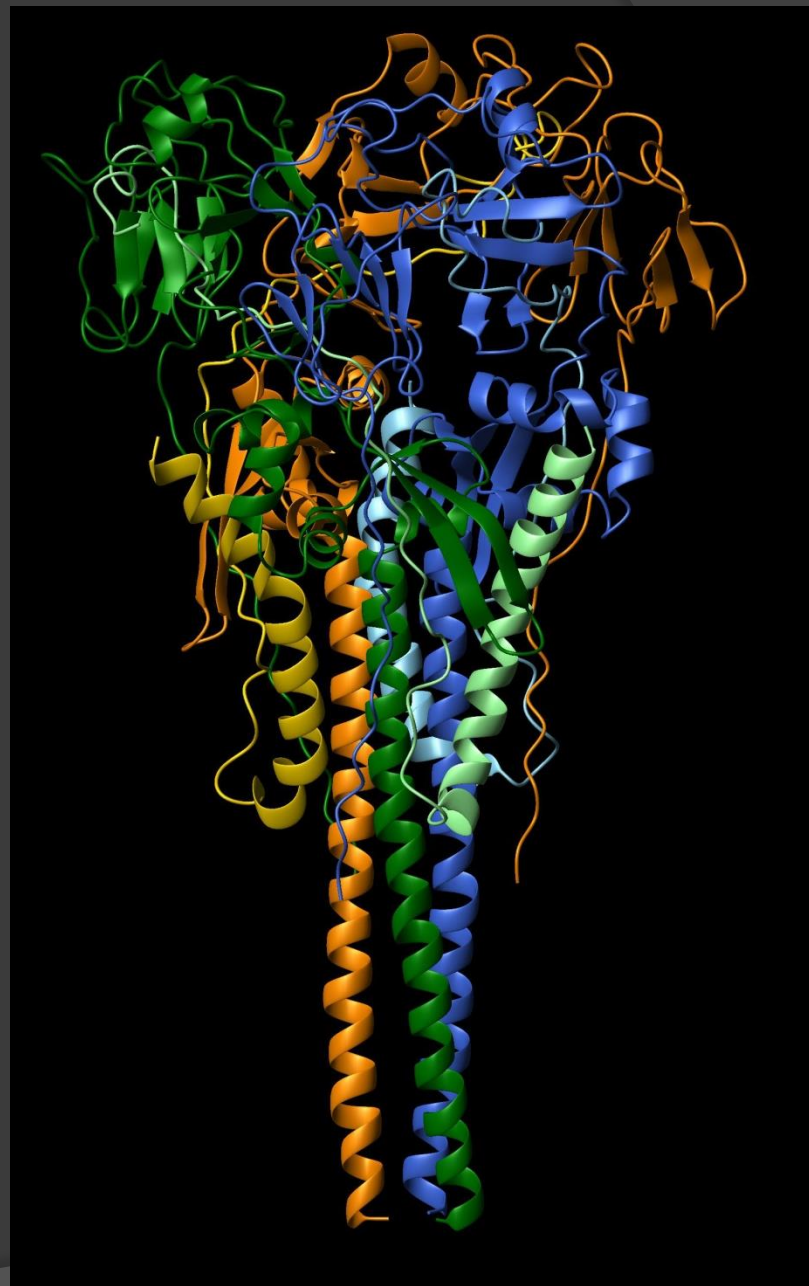
Dimero di
HN dopo
del legame
con il
recettore
cellulare

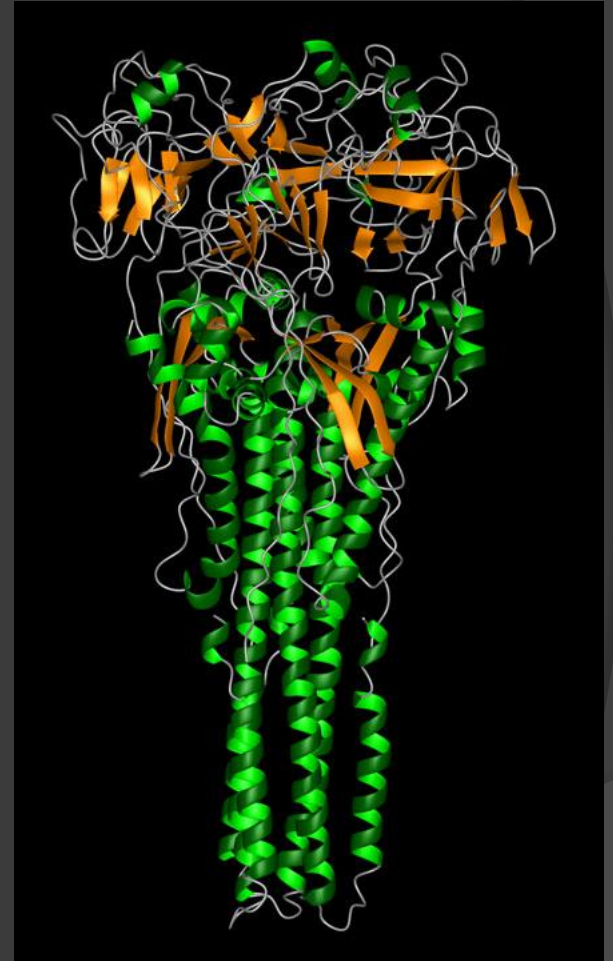
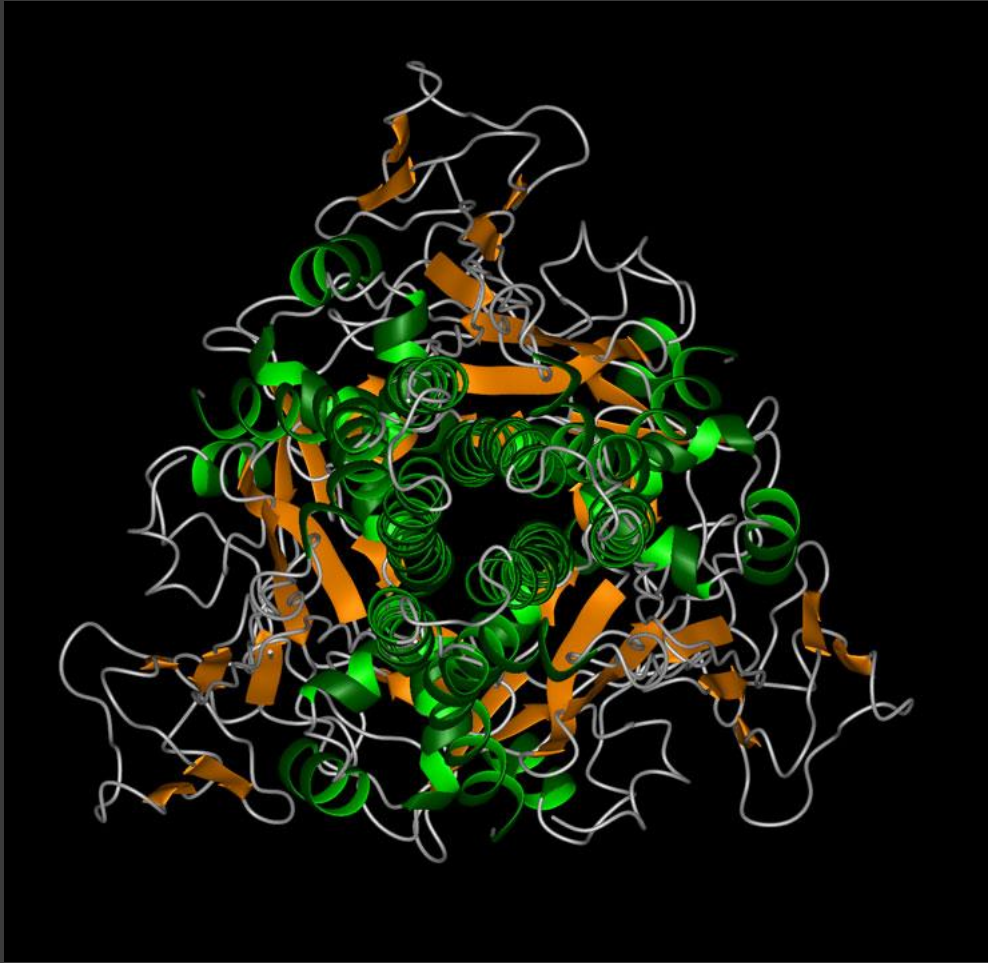


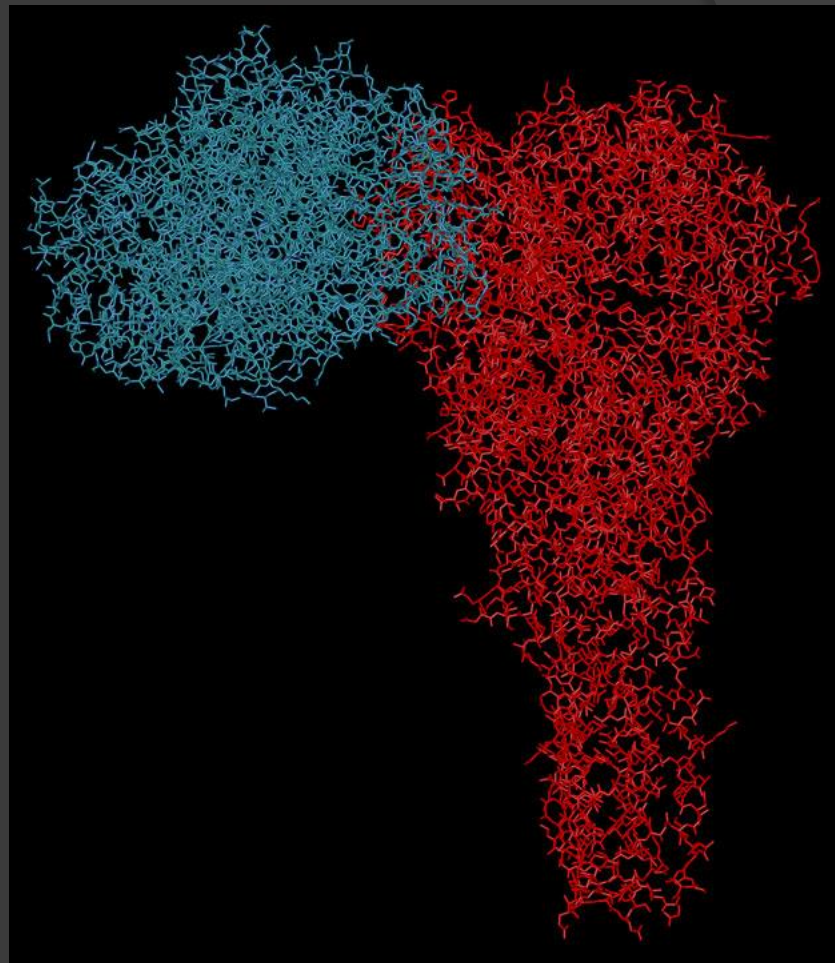
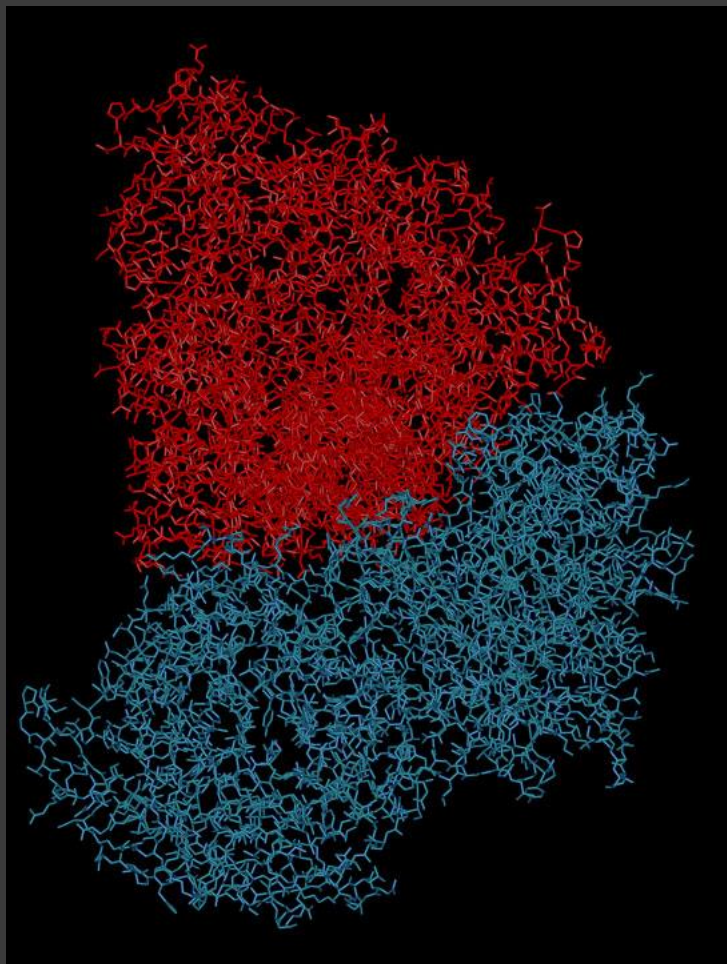
Costruzione del trimero della
proteina F

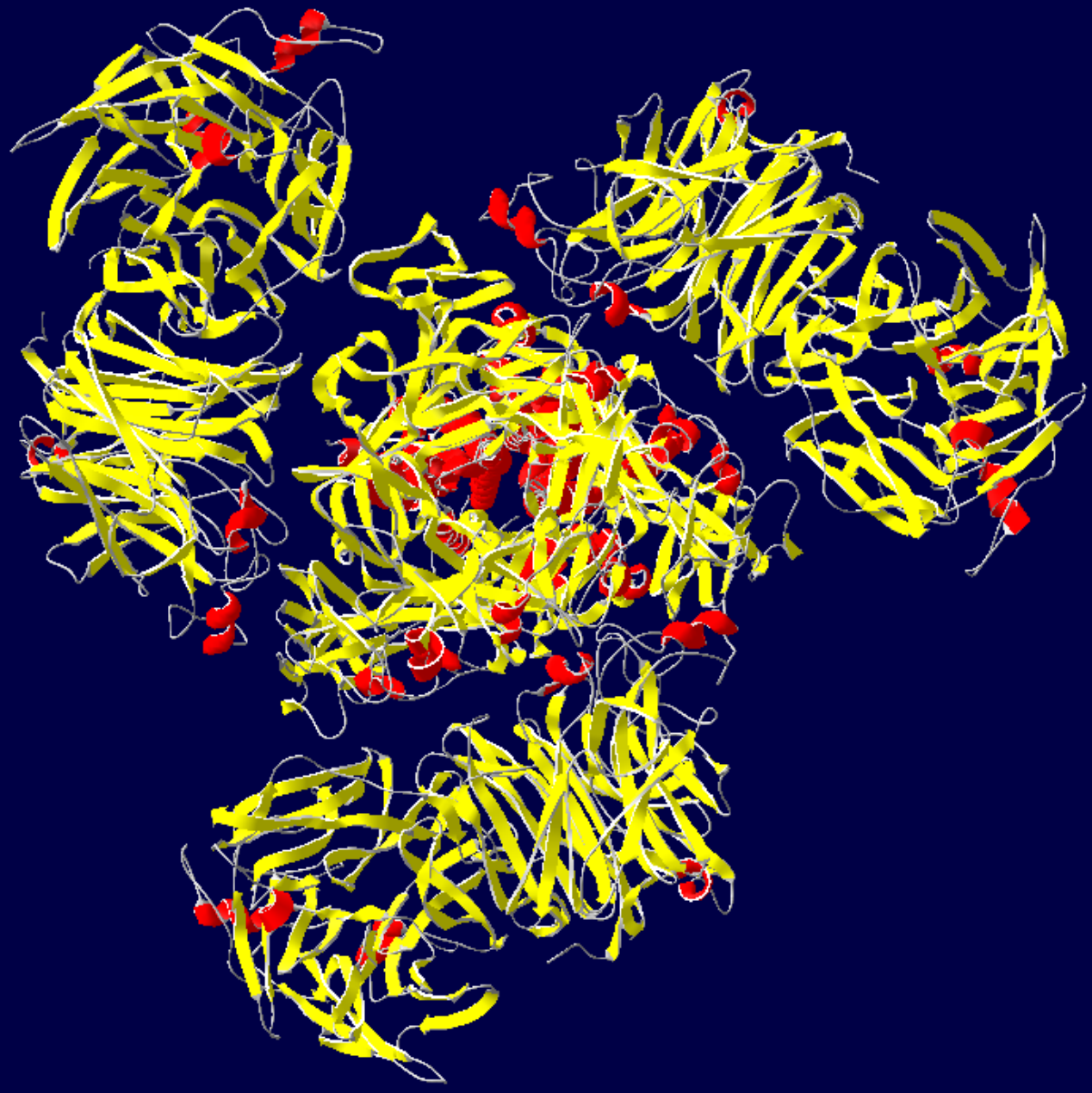


Trimero della proteina F di fusione











INTEGRATED PROJECT MUCOSAL VACCINES FOR POVERTY-RELATED DISEASES

MUVAPRED

Project coordinator: R. Rappuoli, Chiron Srl EC contribution :

15.250.000 Euro

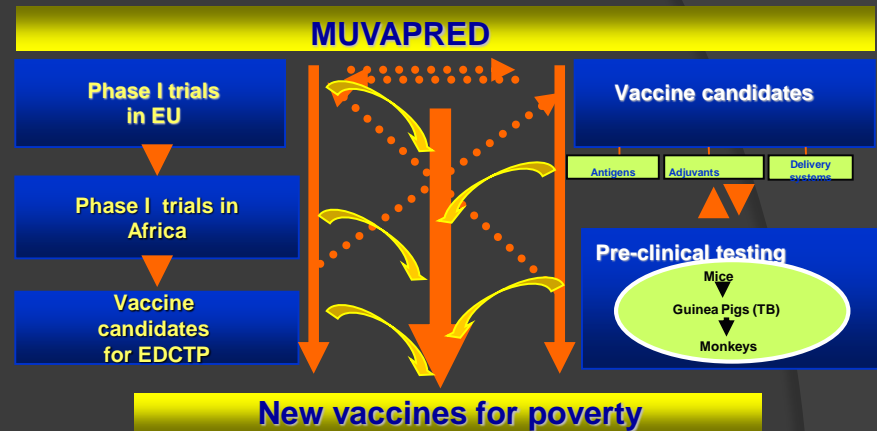
Starting date: 1st December 2003

Duration: 5 years

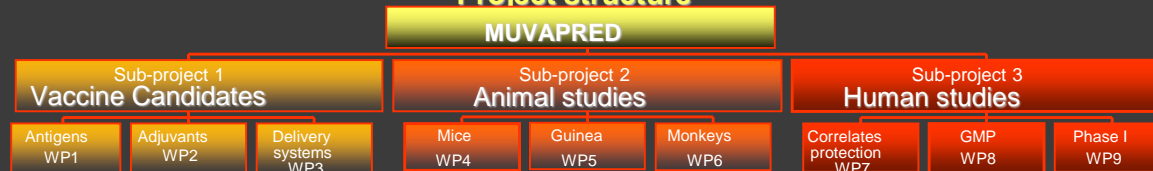
Participants: 24 from 10 countries Project Summary

Human Immunodeficiency Virus and *Mycobacterium tuberculosis* enter the human body at mucosal sites. The aim of the present project is to develop mucosally delivered vaccines against HIV and TB which will induce local immunity able to neutralise the pathogens at their port of entry and systemic immunity able to prevent systemic spread of the infection. The possible development of mucosal vaccines against malaria will be also investigated. The trust of the project derives from the recent proof of concept that mucosal vaccines are feasible in humans.

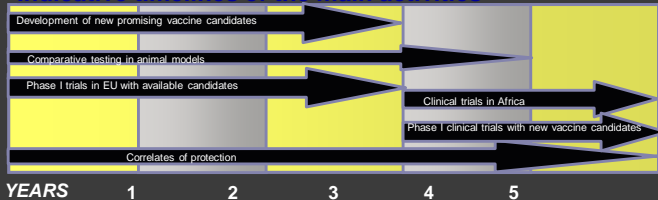
While the first trials are performed, new systems to deliver mucosal vaccines and basic mechanisms of mucosal immune responses and memory in humans will be studied. This will allow better understanding of the clinical results and optimisation of second generation vaccines to be tested in Developing Countries during the second phase of the project.



Project structure



Indicative timelines of the main activities



Objectives

- Phase I clinical trials in Europe with the available mucosal antigens against HIV/AIDS and TB (two to three vaccine candidates).
- Second generation, optimised vaccine candidates against HIV/AIDS, TB and malaria (antigens, adjuvants, delivery systems).
- Selection of the most promising vaccine candidates by comparative testing in animal models
- Phase I clinical trials in Europe with new vaccine candidates developed by the consortium
- Phase I clinical trials in Africa with the vaccine candidates that have proven to be safe and immunogenic in the first clinical trials in Europe

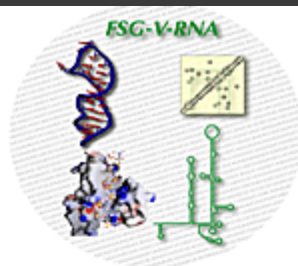
Participants

Chiron Italy, Imperial College of Science UK, University of Goteborg Sweden, University of Siena Italy, Institute Pasteur France, Istituto Superiore di Sanità Italy, Institute for Research in Biomedicine CH, Max-Planck-Institute, Germany University Hamburg Germany, Trinity College Dublin Ireland, Serum State Institute Denmark, Consiglio Nazionale delle Ricerche Italy, Istituto Humanitas Milano Italy, German Arthritis Research Center DRFZ Germany, St George's Vaccine Institute UK, Institute of Microbiology Czech Republic, CAMR UK, University of Florence Italy, Centre-Hospitalier-Universitaire Ignace Deen Guinea, University of Lausanne Switzerland, ALTA Srl Italy, LIONEX Diagnostic & Therapeutics GmbH Germany, INOTECH AG Switzerland

Scientific and Administrative Management



web site: <http://www.mucosalimmunity.org/muvapred/>



Functional and Structural Genomics of Viral RNA



6th Framework Program EU

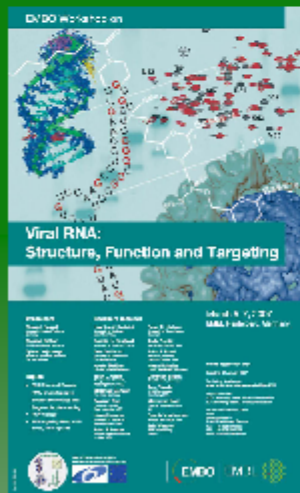


Login:

Pass:

[Home](#) | [Project](#) | [Partners](#) | [Meetings](#) | [Results](#) | [Contact](#) | [Links](#)

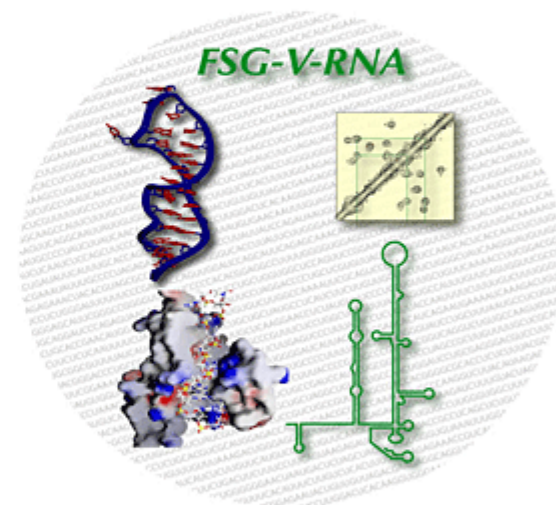
Don't miss the
FSGVRNA WORKSHOP 2007
 (Heidelberg, Germany)
 5-7 March 2007



The current project joins, in an interdisciplinary fashion, leading European labs to integrate their equipment and expertises on:

» The structural, functional and virological analysis of RNA and RNA-protein complexes from viruses.

» The evaluation of these viral RNAs as targets for novel types of drugs, either small RNA binding compounds, or antiviral RNAs.



Acronym: FSG-V-RNA

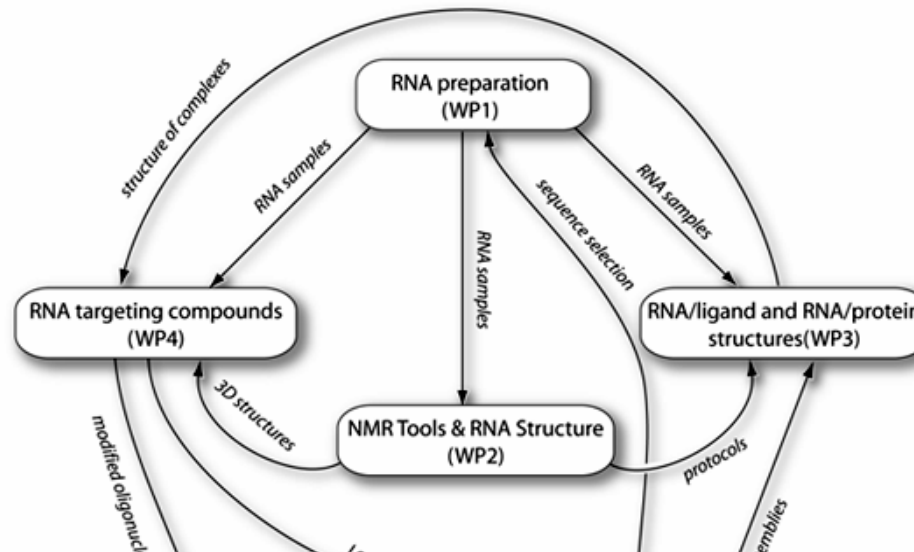
Research Topic: Life sciences, genomics and biotechnology for health, LSH-2002-1.1.0-1: For STREP and CA, research should focus on multidisciplinary functional genomics approaches. Proposals concerned with the development of new tools and approaches, including the standardisation of protocols, to facilitate generation of new knowledge in functional and structural genomics are also envisaged.

The STREP will develop and improve tools and approaches to facilitate the generation of new knowledge in functional and structural genomics of viral RNAs. Specifically, (i) new methods and tools for the rapid and efficient structural analysis of RNA and RNA-protein complexes will be developed, (ii) these optimised tools will be applied to essential RNA elements that are vital for the function of HBV, HCV and HIV viruses, and (iii) complementary screening techniques and structure analysis of RNA-ligand complexes will be implemented to promote the identification of antiviral compounds targeting these RNA structures. The project exploits available RNA sequence data but will also expand our knowledge of viral RNA sequence elements and their variations.

The biomedical importance of RNA as a research target is stressed by the fact that viral infections, such as HBV, HCV and HIV are major global public health problems. The outcomes of the project are expected to initiate the development of novel drugs that target viral RNA molecules and have thus strong implications for public health. The research consortium is strengthened by the involvement of an SME, which contributes a proprietary RNA screening method. This SME will also function as a stakeholder to support the downstream application development towards RNA targeting drugs - even though this is beyond the scope of the current project.

The project involves multidisciplinary research which is only possible by integrating, in an interdisciplinary fashion, the research capacities of a number of leading European labs, i.e. their equipment and complementary expertise on the structural, functional and virological analysis of RNA and RNA-ligand complexes. The innovative tools developed will be made available to other researchers throughout Europe and open the way for efficient analysis of a wide range of RNA-based processes extending far beyond analysis of viral RNAs.

Work packages





Jobs | Contact

Home	Complexes	Press & Publications	Workshops & Meetings	Methods	Structural Genomics	Private
Project	Summary	Objectives	Workpackages	Partners		

Welcome

In the last two decades biological science has made huge progress in many areas. We live in an era when the entire genetic code of many organisms has been established. However, to make full use of this treasure of information, we need to bring it together with the knowledge of what the products of all these genes, the proteins, are doing. At present, even though we appear to have a vast amount of information at our disposal, it would be unreasonable to expect that in the near or medium term we could reach this level of understanding for an entire cell. So, a pre-requisite for this goal is comprehensive knowledge of the biological functions of the complete set of genes and proteins within a genome (post-genomic biology).

Introduction

Proteins rarely act alone: they typically interact with other macromolecules to perform particular cellular tasks. The resulting functional assemblies (complexes) achieve more than the sum of their parts and these complexes have functions that cannot easily be understood by even the most systematic analyses of individual proteins. So, the discovery and analysis of particular cellular protein complexes under physiological conditions provides key insights into their function and takes characterisation of cellular systems well beyond the limits of other experiments. Prominent examples include the ribosome, the chaperonin GroEl/GroEs, the spliceosome, the cytosome, the proteasome, the nuclear pore complex and the synaptosome. Analyses of results from genome-scale interaction experiments in yeast show a clear tendency for many yeast assemblies to mirror their equivalents in animals, including the model organisms and man. Complexes essential for normal cell activity overlap significantly and represent the building blocks of a Eukaryotic core proteome, covering basic cellular function. More importantly, those conserved between yeast and man contribute significantly to the understanding of multifactor diseases, particularly those related to key cellular processes. Elucidation of three-dimensional (3D) structures for protein complexes will open new avenues to unravel the molecular pathology and physiology of human diseases, leading to rational, target-oriented therapeutic approaches. Moreover developments in 3D tomography show that it will soon be possible to fit such structures into a whole cell tomogram. This will be a great leap for systems biology, since it will place complexes in their precise cellular context and provide critical concentration information essential for the quantitative understanding of a living cell. However, without the individual complexes, it will be exceedingly difficult to understand such whole cell images.

3D-Repertoire is an integrated project funded by the European Commission under Framework 6. It brings together scientists from 14 institutions and private companies, in seven countries across Europe, with the aim to resolve structures for all amenable protein complexes from budding yeast (or where necessary equivalents from other species) at the best possible resolution by electron microscopy, X-ray crystallography and in silico approximations. The team of top scientists formed through 3D-Repertoire covers the wide range of expertise needed for such an ambitious undertaking and we have secured access to the yeast constructions used to do the complete pull-down study of complexes in Yeast (given to Euroscarf by Cellzome).



News
[2nd Annual Meeting](#),
22-23/02/07



2nd 3D Repertoire Annual Meeting