

Lecture 1 - Introduction to Structural Bioinformatics

Motivation and Basics of Protein Structure

Objectives of the course

- Understanding protein function.
- Applications to Computer Aided Drug Design.
- Development of efficient algorithms to evaluate the above “*in silico*”.
- Emphasis on the “structure” related problems – Geometric Computing in Molecular Biology.
- Show relevance to other spatial “pattern discovery” tasks.

Most of the Protein Structure slides – courtesy of
Hadar Benyaminy.

Textbook

There is no single, double or triple textbook for this course.

Most of the material is based on journal articles and research done by the Wolfson-Nussinov Structural Bioinformatics group at TAU.

Nevertheless :

Recommended Literature (1):

- Setubal and Meidanis, Introduction to Computational Biology, (1997).
- A. Lesk, *Introduction to Protein Architecture*, 2'nd edition (2001).
- S.L. Salzberg, D.B.Searls, S. Kasif (editors), Computational Methods in Molecular Biology, (1998).

Recommended Literature (2):

- Branden and Tooze, Introduction to Protein Structure (2'nd edition).
- D. Gusfield, Algorithms on Strings, Trees and Sequences, (1997).
- Voet and Voet, Biochemistry (or, any other Biochemistry book in the Library).
- M. Waterman, Introduction to Computational Biology.

Strongly Recommended Literature (currently not in the library):

- Protein Bioinformatics.
- Structural Bioinformatics.

Recommended Web Sites:

- Enormous number of sites.
- Search using “google”.
- PDB site <http://www.rcsb.org/pdb/>
- Birbeck course on protein structure.

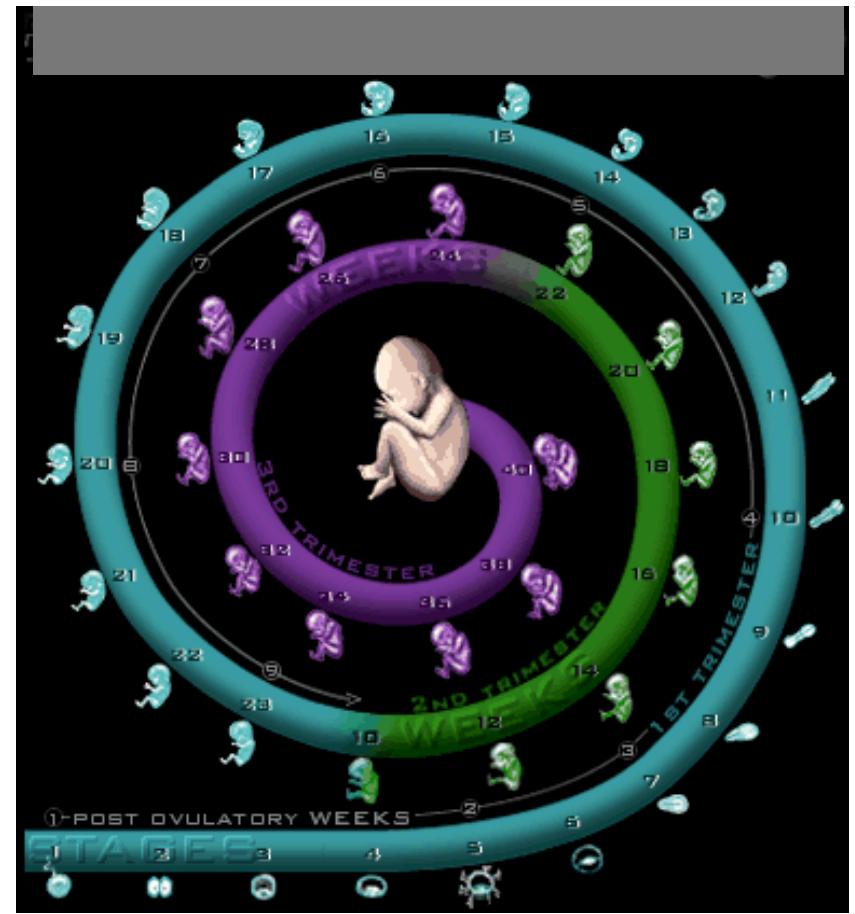
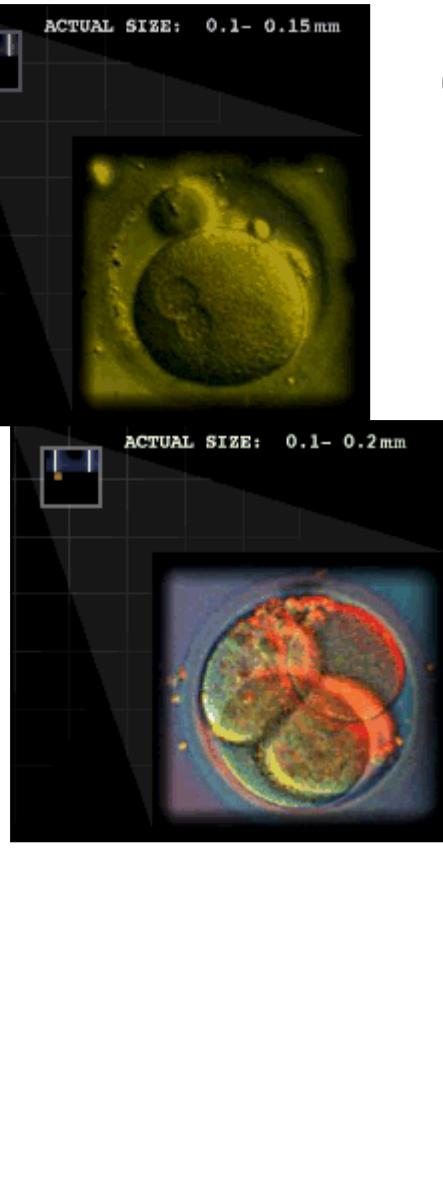
Journals :

- Proteins : Structure, Function, bioinformatics.
- Journal of Computational Biology.
- Bioinformatics (former CABIOS).
- Journal of Molecular Biology.
- Journal of Computer Aided Molecular Design.
- Journal of Molecular Graphics and Modelling.
- Protein Engineering.

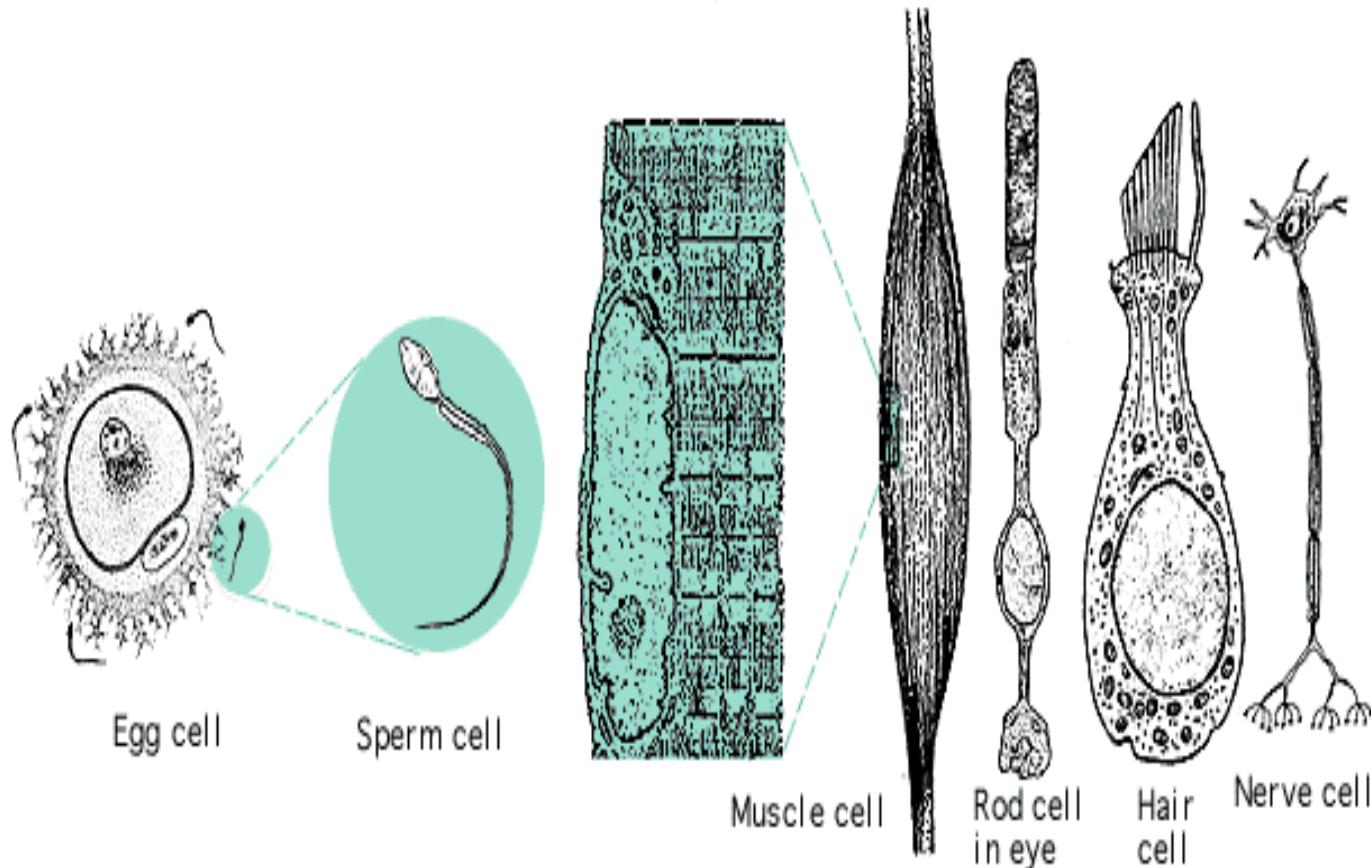
Computational Biology Conferences:

- ISMB - International Conference on Intelligent Systems in Molecular Biology.
- RECOMB - Int. Conference of Computational Molecular Biology.
- ECCB - European Conference on Computational Bio.
- WABI - Workshop of Algorithms in Bioinformatics .

Cell- the basic life unit



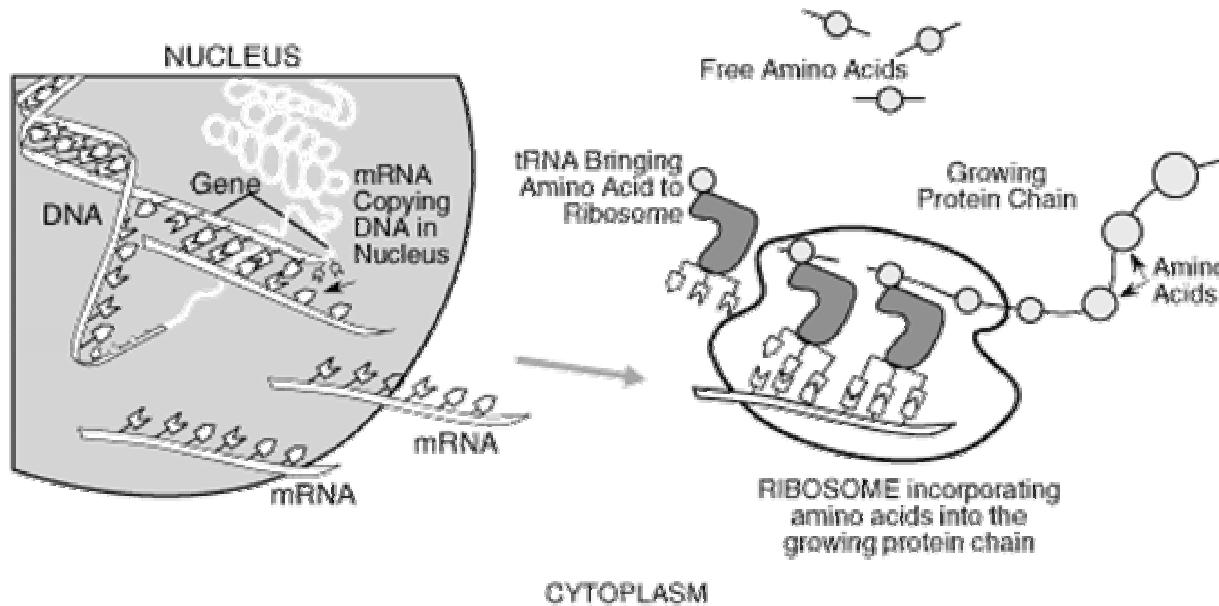
Different cell types



Size of protein molecules (diameter)

- cell $(1 \times 10^{-6} \text{ m})$ μ microns
- ribosome $(1 \times 10^{-9} \text{ m})$ nanometers
- protein $(1 \times 10^{-10} \text{ m})$ angstroms

The central dogma



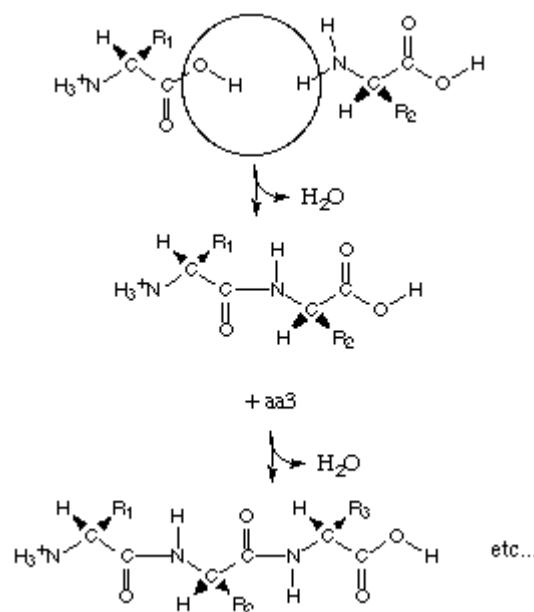
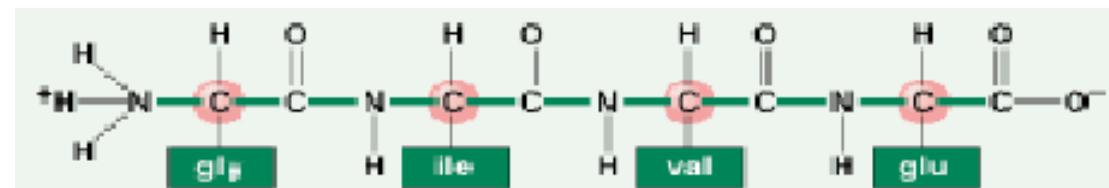
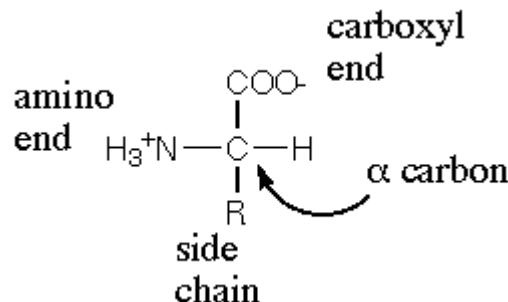
When genes are expressed, the genetic information (base sequence) on DNA is first transcribed (copied) to a molecule of messenger RNA in a process similar to DNA replication. The mRNA molecules then leave the cell nucleus and enter the cytoplasm, where triplets of (codons) forming the genetic code specify the particular amino acids that make up an individual protein.

This process, called translation, is accomplished by ribosomes (cellular components composed of proteins and another class of RNA) that read the genetic code from the mRNA, and transfer RNAs (tRNAs) that transport amino acids to the ribosomes for attachment to the growing protein. (From www.ornl.gov/hgmis/publicat/primer/)

Proteins – our molecular machines (samples of protein tasks)

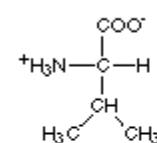
- Catalysis (enzymes).
- Signal propagation.
- Transport.
- Storage.
- Receptors (e.g. antibodies – immune system).
- Structural proteins (hair, skin, nails).

Amino acids and the peptide bond

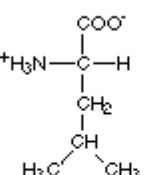


C_β – first side chain carbon (except for glycine).

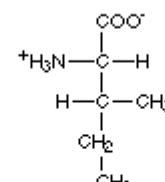
Amino acids with hydrophobic side groups



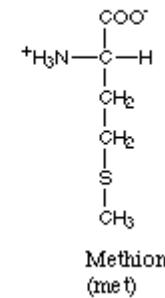
Valine
(val)



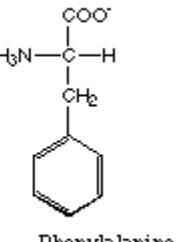
Leucine
(leu)



Isoleucine
(ile)



Methionine
(met)



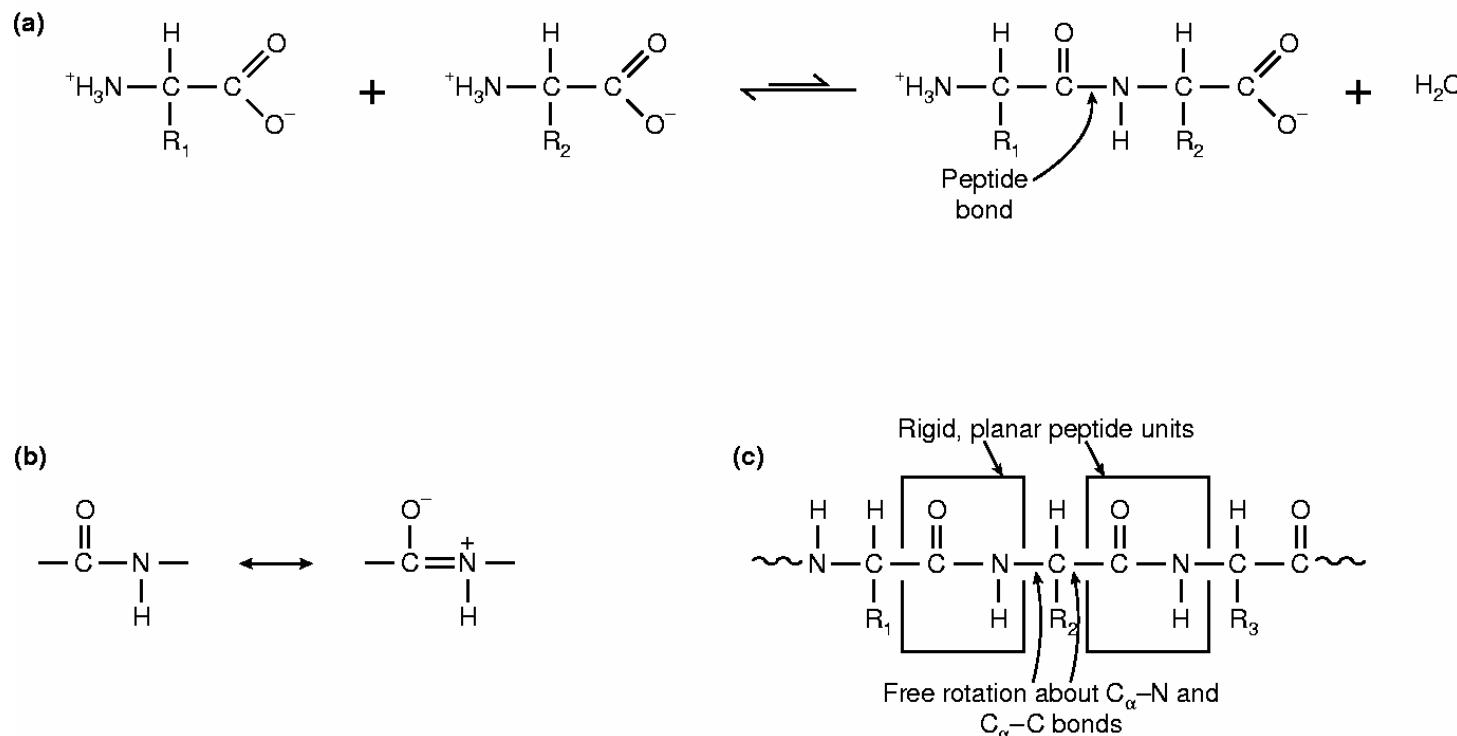
Phenylalanine
(phe)

Primary through Quaternary structure

- Primary structure: The order of the amino acids composing the protein.
- AASGDXSLVEVHXXVFIVPPXIL.....

Folding of the Protein Backbone

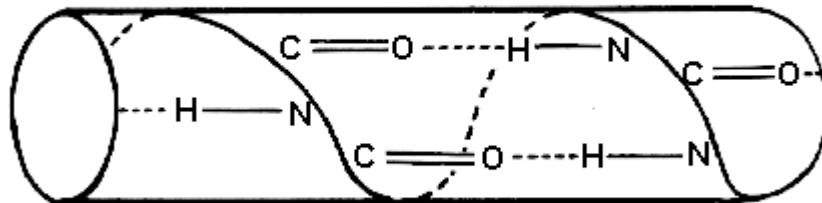
B2.1



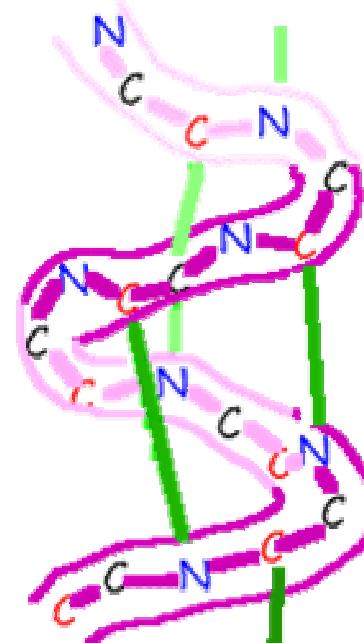
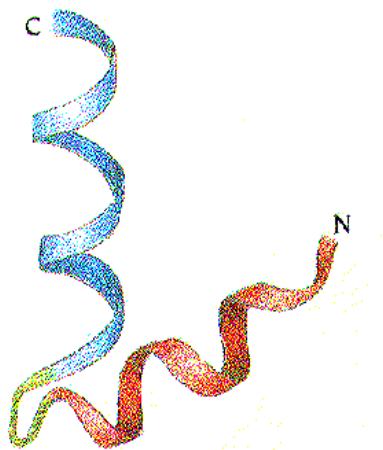
The Holy Grail - Protein Folding

- How does a protein “know” its 3-D structure ?
- How does it compute it so fast ?
- Relatively primitive computational folding models have proved to be NP complete even in the 2-D case.

Secondary structure



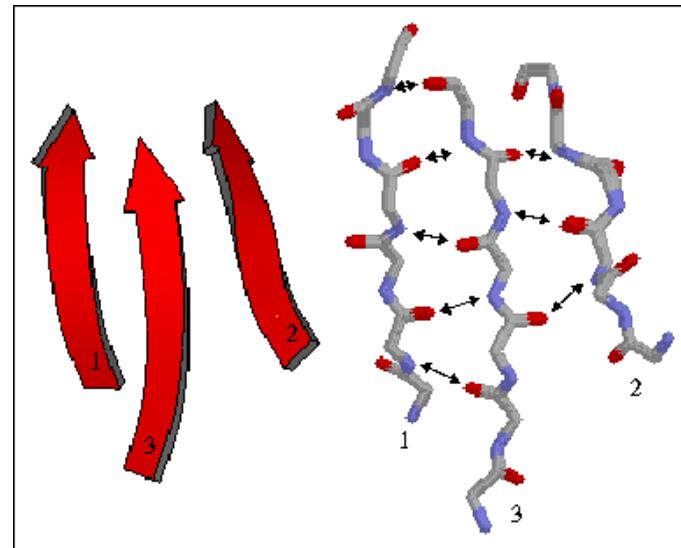
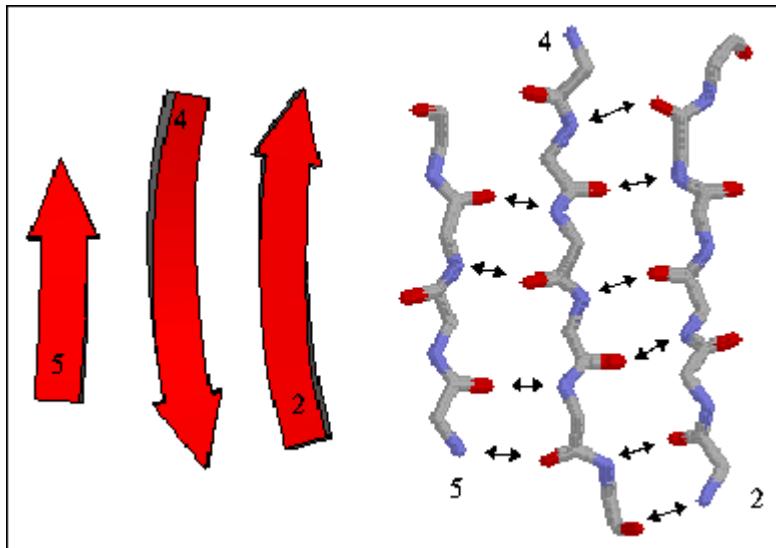
3.6 residues/turn (5.4 Å dist.)



2: Backbone:

- N Nitrogen
- C Alpha Carbon
- C Carboxyl Carbon
- Hydrogen bond

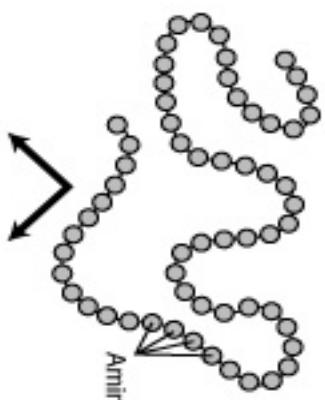
β strands and sheets



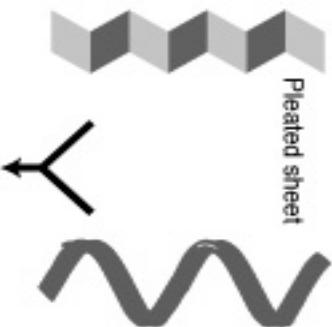
Bond. Hydrogen bond.

Primary protein structure
is sequence of a chain of amino acids

Amino Acids

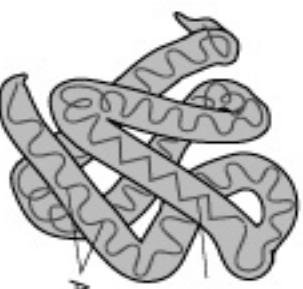


Pleated sheet
Alpha helix



Secondary protein structure
occurs when the sequence of amino acids
are linked by hydrogen bonds

Pleated sheet
Tertiary protein structure
occurs when certain attractions are present
between alpha helices and pleated sheets.



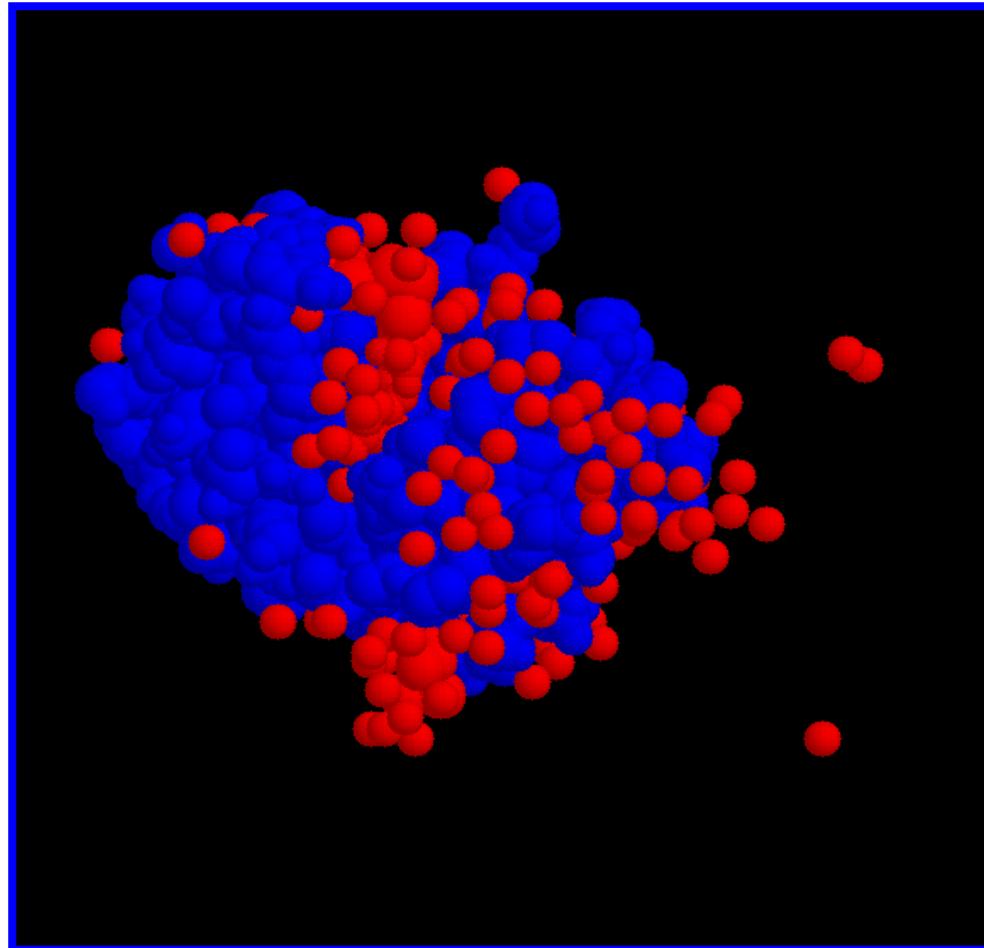
Quaternary protein structure
is a protein consisting of more than one
amino acid chain.



Wire-frame or ribbons display

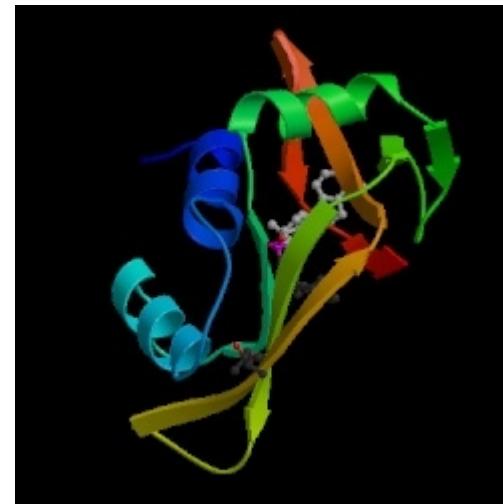


Space-fill display



Tertiary structure: full 3D folded structure of the polypeptide chain

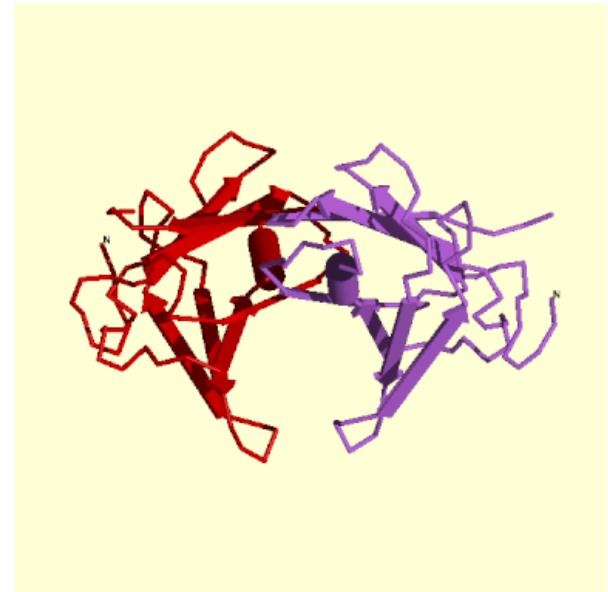
Ribonuclease - PDB code 1rpg



Quaternary structure

The interconnections and organization of more than one polypeptide chain.

Example :Transthyretin dimer (**1tta**)



Determination of protein structures

- X-ray Crystallography
- NMR (Nuclear Magnetic Resonance)
- EM (Electron microscopy)
- Nano – sensors (?)

X-ray Crystallography

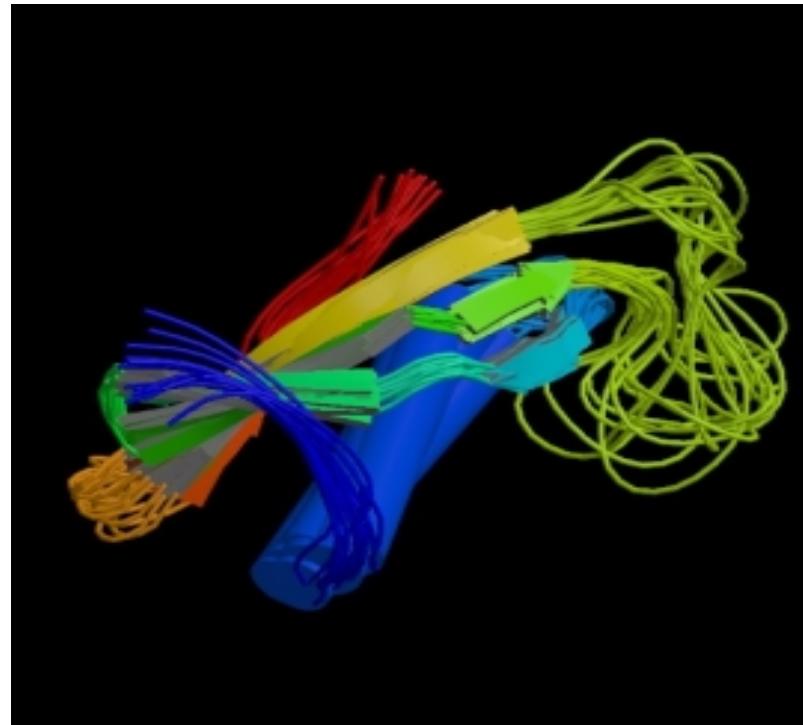
- Crystallization
- Each protein has a unique **X-ray pattern diffraction**.
- The **electron density map** is used to build a model of the protein.

Nuclear Magnetic Resonance

- Performed in an **aqueous solution**.
- NMR analysis gives **a set of estimates of distances between specific pairs of protons (H – atoms)**.
- Solved by Distance Geometry methods.
- The result is an **ensemble of models** rather than a single structure.

An NMR result is an ensemble of models

Cystatin (1a67)



The Protein Data Bank (PDB)

- International repository of 3D molecular data.
- Contains x-y-z coordinates of all atoms of the molecule and additional data.

Welcome to the PDB, the single international repository for the processing and distribution of 3-D macromolecular structure data primarily determined experimentally by [X-ray crystallography](#) and [NMR](#).

[DEPOSIT](#) Contribute structure data

[STATUS](#) Find entries awaiting release

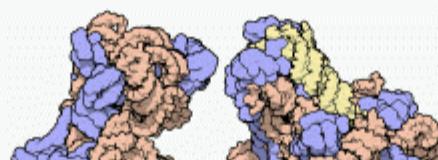
[DOWNLOAD](#) Retrieve structure files (FTP)

[LINKS](#) Browse related information

[PREVIEW](#) Beta-test new features

About the PDB

[General Information](#)
[WWW User Guides](#)
[Get Educated](#)



Current Holdings

[13505 Structures](#)

[Last Update: 24-Oct-2000](#)

[PDB Statistics](#)

Search

Enter a PDB ID: [Explore](#)

[SearchLite](#): simple keyword search

[SearchFields](#): advanced search

News

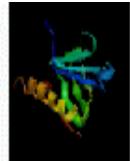
[Complete News](#)
[PDB Newsletter](#)
[Subscribe](#)
[Browse Mailing List](#)

24-Oct-2000

[Issue 7 of the PDB Newsletter Now](#)

Feb. 2003 – about 20,000 structures.
Structural Bioinformatics 2004

Prof. Haim J. Wolfson



[Summary Information](#)

[View Structure](#)

[Download/Display File](#)

[Structural Neighbors](#)

[Geometry](#)

[Other Sources](#)

[Sequence Details](#)

Explore

[SearchLite](#) [SearchFields](#)

Summary Information



Compound: Mol_Id: 1; Molecule: Irs-1; Chain: A; Fragment: PtB Domain; Synonym: Insulin Receptor Substrate 1; Engineered: Yes
 Mol_Id: 2; Molecule: IL-4 Receptor Phosphopeptide; Chain: B; Engineered: Yes

Authors: M.-M. Zhou, B. Huang, E. T. Olejniczak, R. P. Meadows, S. B. Shuker, M. Miyazaki, T. Trub, S. E. Shoelson, S. W. Fesik

Exp. Method: NMR, Minimized Average Structure

Classification: Complex (Signal Transduction/Peptide)

Source: Homo Sapiens

Primary Citation: Zhou, M. M., Huang, B., Olejniczak, E. T., Meadows, R. P., Shuker, S. B., Miyazaki, M., Trub, T., Shoelson, S. E., Fesik, S. W.: Structural basis for IL-4 receptor phosphopeptide recognition by the IRS-1 PTB domain. *Nat Struct Biol* 3 pp. 388 (1996)
[\[Medline \]](#)

Deposition Date: 22-Mar-1996

Release Date: 15-May-1997

Polymer Chains: A, B

Residues: 123

Atoms: 971

HET groups:

ID	Name	Formula
PTR	PHOSPHOTYROSINE	C ₉ H ₁₂ N ₁ O ₆ P ₁

Classification of 3D structures

SCOP

- Provides a description of the structural and evolutionary relationships between all proteins whose structure is known.
- Created largely by manual inspection.
- J. Mol. Biol. 247, 536-540, 1995

SCOP

Structural Classification of Proteins



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

Lineage:

1. Root: [scop](#)
2. Class: [All alpha proteins](#)
3. Fold: [Globin-like](#)
core: 6 helices; folded leaf, partly opened
4. Superfamily: [Globin-like](#)
5. Family: [Globins](#)
Heme-binding protein
6. Protein: Hemoglobin, alpha-chain
7. Species: [Human \(*Homo sapiens*\)](#)

PDB Entry Domains:

1. [1bab](#)
complexed with hem, so4; mutant
 1. [chain a](#)
 2. [chain c](#)
2. [1bz0](#)

CATH - Protein Structure Classification

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



CATH



Protein Structure Classification

Version 1.6 : Released June 1999

Welcome to the **CATH** protein classification home page
[Biomolecular Structure and Modelling Unit](#),
University College London.

*Dr. Frances M.G. Pearl, Mr. James Bray, Ms. Annabel E. Todd,
Dr. David Lee, Dr. Adrian J. Shepherd, Dr. Andrew Harrison, Prof. Janet Thornton
Dr. Christine A. Orengo*

Available options:

- [Browse or search classification](#)
- [Lexicon](#)
- [Glossary](#)

CATH

- **Class**: derived from secondary structure content.
- **Architecture**: gross orientation of secondary structures, independent of connectivities.
- **Topology**: clusters according to topological connections and numbers of secondary structures.

- **Homology**: clusters according to structure and function.



● Class ● Architecture ● Topology ● Homologous superfamily

1.10.490

C ● Mainly Alpha

A ● Non-Bundle

T ● Globin-like

● [H 10] 1hlm (12 S's)

● [S 1] 1hlm
● 1hlm GET DATA... ▾

PDB header: Oxygen Transport
PDB comp: Hemoglobin (Cyano-Met) (Sea Cucumber)
PDB source: Sea Cucumber (Caudina (Molpadia) Arenicola)

● 1hlb GET DATA... ▾

PDB header: Oxygen Transport
PDB comp: Hemoglobin (Sea Cucumber)
PDB source: Sea Cucumber (Caudina (Molpadia) Arenicola)

● [S 2] 1hbg
● 1hbg GET DATA... ▾

PDB header: Oxygen Transport
PDB comp: Hemoglobin (Carbon Monoxo)
PDB source: Marine Bloodworm (Glycera Dibranchiata)

● 2hbg GET DATA... ▾

PDB header: Oxygen Transport
PDB comp: Hemoglobin (Deoxy)

CATH LEXICON

● Mainly Alpha

These proteins consist predominantly of alpha helix secondary structures, although many also contain a small percentage of beta sheet on their peripheries. Mainly alpha proteins have been assigned using a cutoff of >60% alpha and <5% beta secondary structure assignment. In addition these proteins must have >50% alpha-alpha and <5% beta-beta secondary structure contacts (Michie *et al.*, 1996)

● Non-Bundle

The non-bundle architecture is a general

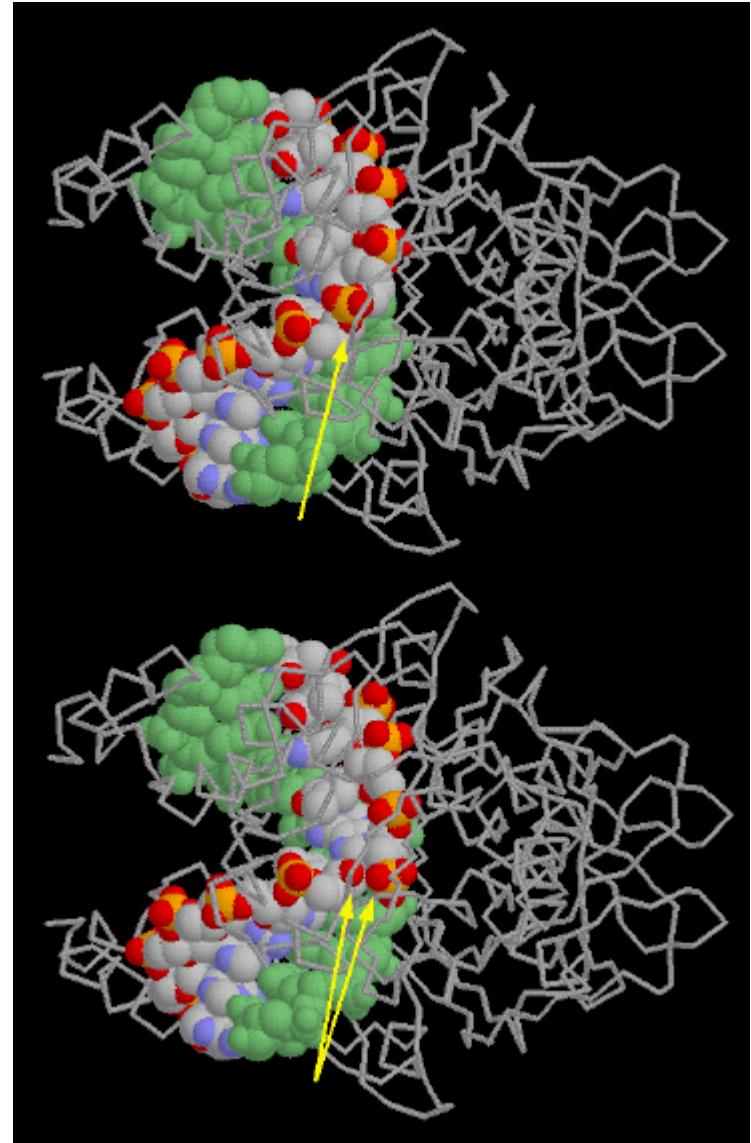
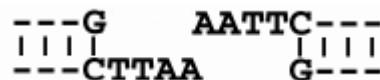
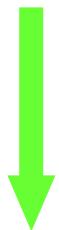


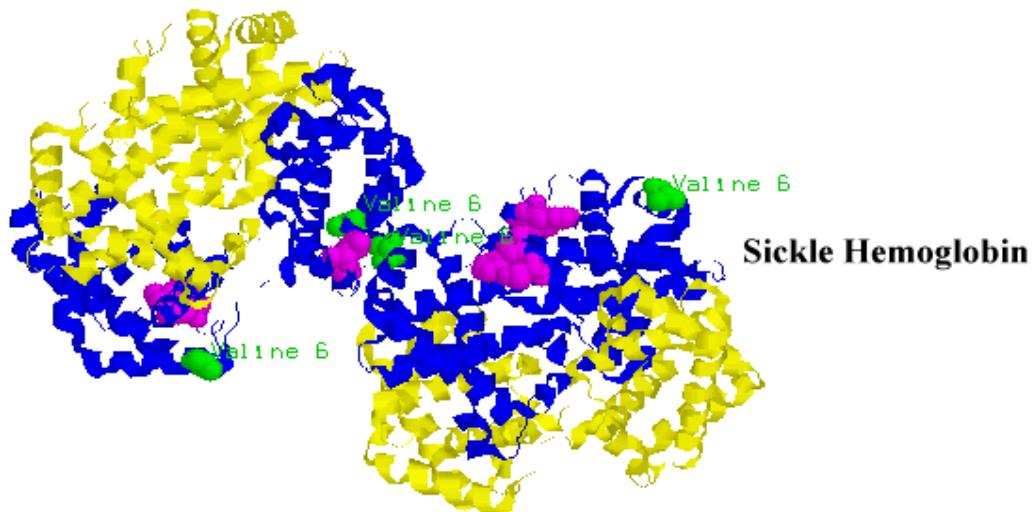
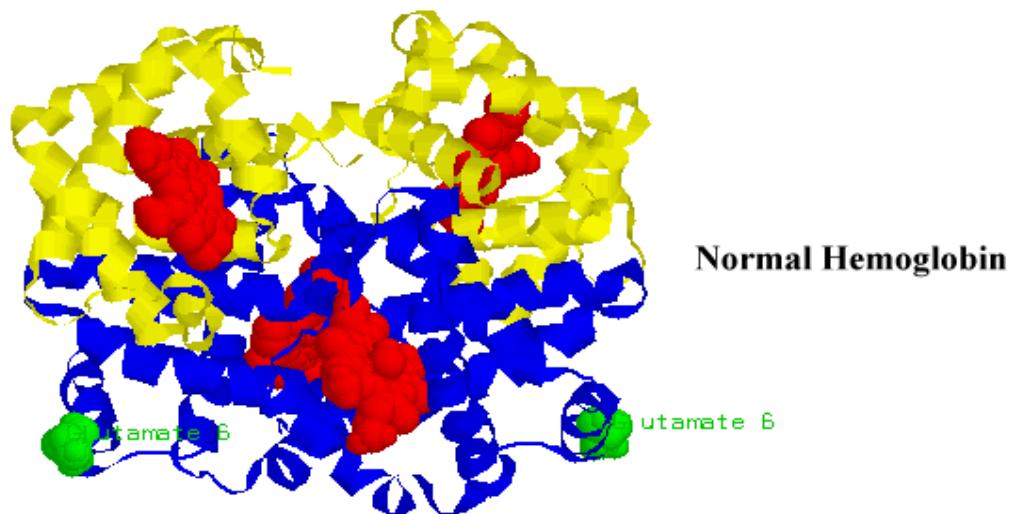
Representative for CATH code **1.10.490.10.1**

Mainly Alpha : Non-Bundle : Globin-like : 1hlm
PDB: [1hlm](#)

- PDB <http://pdb.tau.ac.il>
- PDB <http://www.rcsb.org/pdb/>
- CATH
- <http://www.biochem.ucl.ac.uk/bsm/cath/>
- SCOP <http://scop.mrc-lmb.cam.ac.uk/scop/>

Restriction enzymes

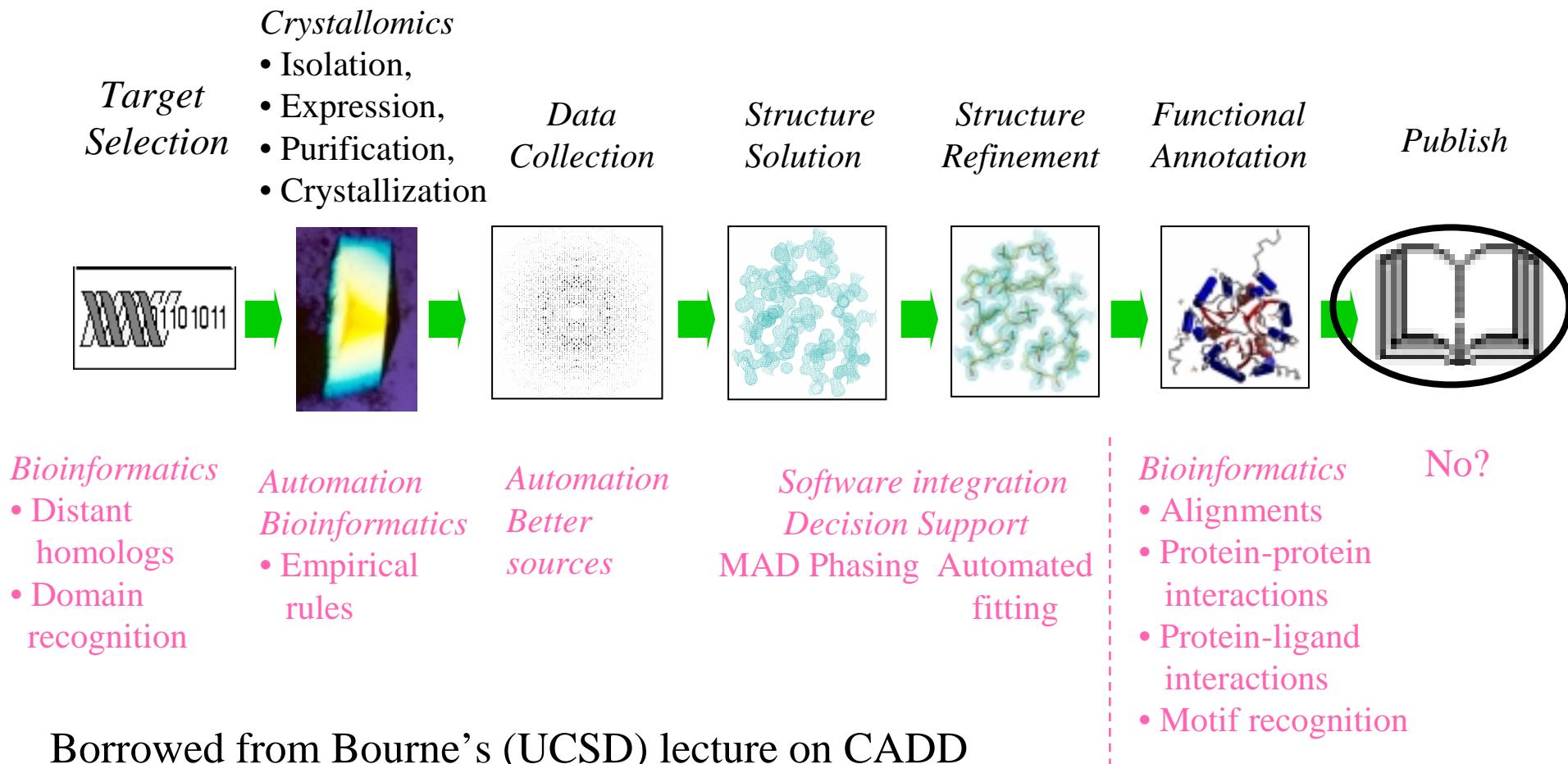




Note: The Sickle hemoglobin image is drawn at 50% of the size
of the Normal hemoglobin

The Structural Genomics Pipeline (X-ray Crystallography)

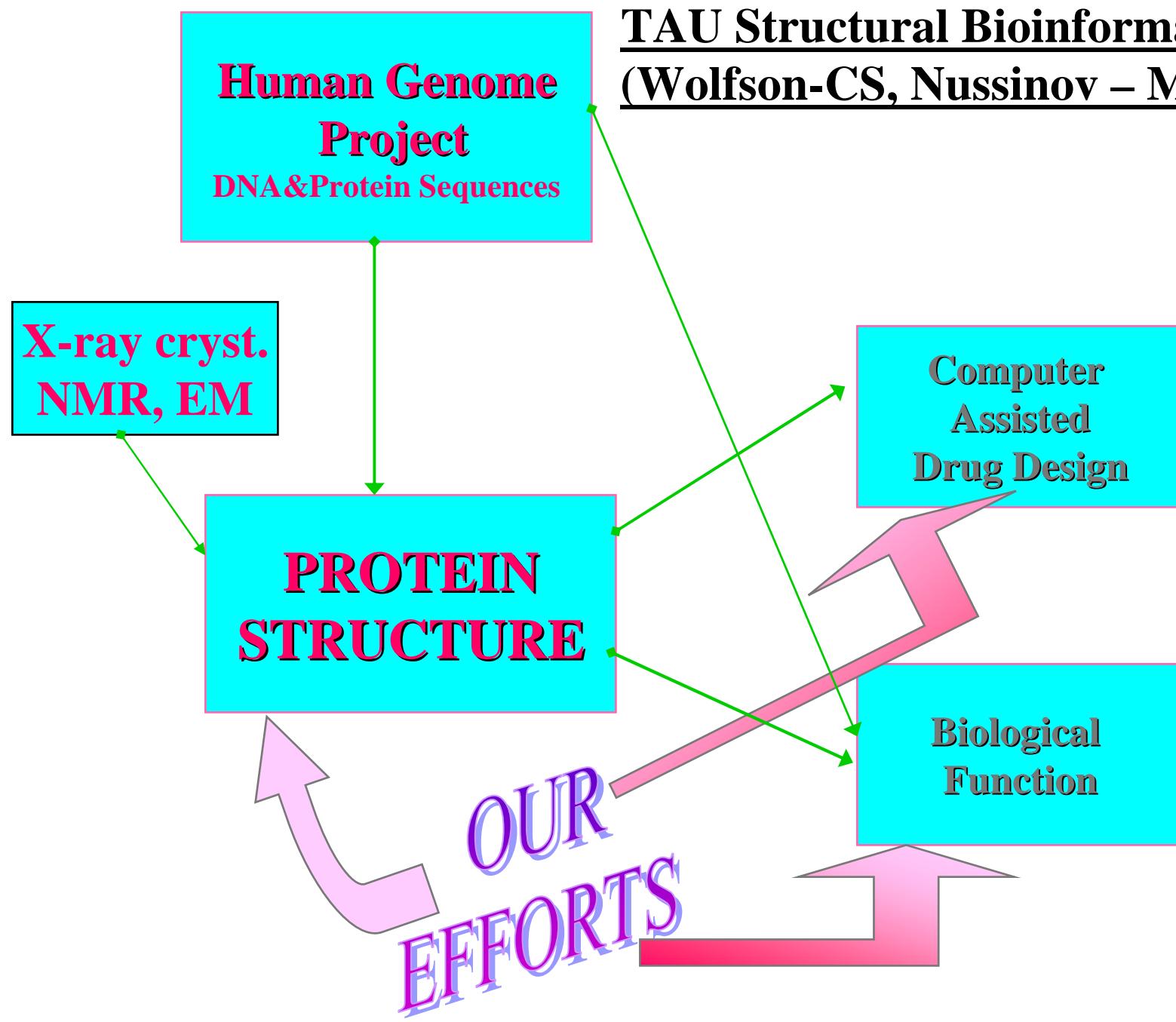
Basic Steps



Borrowed from Bourne's (UCSD) lecture on CADD

TAU Structural Bioinformatics Lab

(Wolfson-CS, Nussinov – MB)



Structural Bioinformatics Lab Goals

Development of *state of the art* algorithmic methods to tackle major computational tasks in protein structure analysis, biomolecular recognition, and *Computer Assisted Drug Design*.

Establish truly *interdisciplinary* collaboration between Life and Computer Sciences.

Bioinformatics and Genomics - Economic Impact

- Medicine and public health.
- Pharmaceuticals.
- Agriculture.
- Food industry.
- Biological Computers (?).

Bioinformatics and Genomics - the Computational Viewpoint

- Molecular Biology is becoming a Computational Science.
- The emergence of large databases of DNA, proteins, small molecules and drugs requires computational techniques to analyze the data.
- Efficient CPU and memory intensive algorithms are being developed.
- Many of the computational tasks have analogs in other well established fields of Computer Science allowing cross-fertilization of ideas.

Bioinformatics - Computational Genomics

- DNA mapping.
- Protein or DNA sequence comparisons ,
primary structure.
- Exploration of huge textual databases.
- In essence one- dimensional methods
and intuition.
- Graph - theoretic methods.

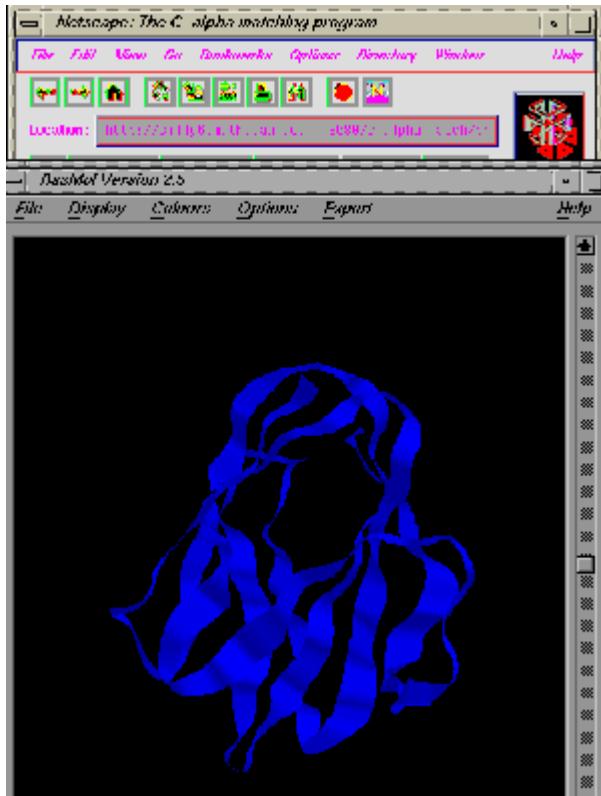
Structural Bioinformatics - Structural Genomics

- Elucidation of the 3D structures of biomolecules.
- Analysis and comparison of biomolecular structures.
- Prediction of biomolecular recognition.
- Handles three-dimensional (3-D) structures.
- *Geometric Computing.*

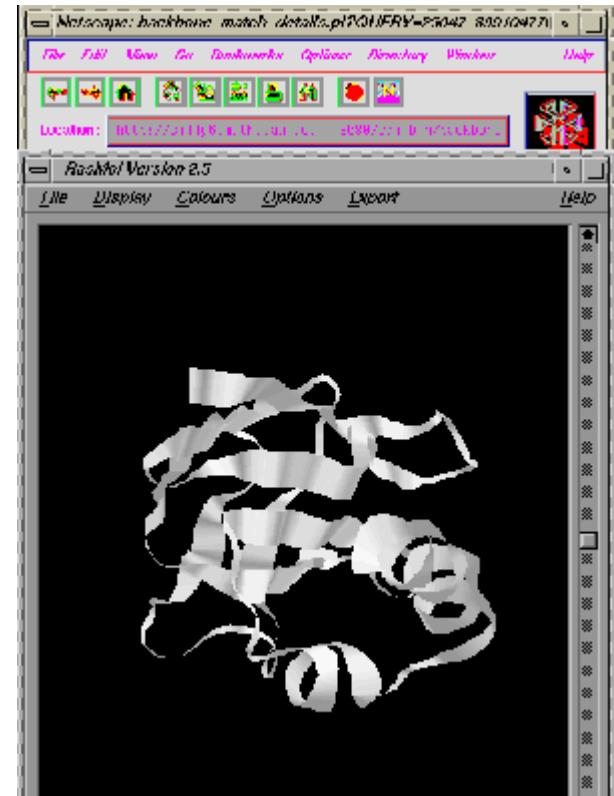
Why bother with structures when we have sequences ?

- In evolutionary related proteins structure is much better preserved than sequence.
- Structural motifs may predict similar biological function .
- Getting insight into protein folding. Recovering the limited (?) number of protein folds.

Case in Point : *Protein Structural Comparison*



ApoAmicyanin - 1aaJ



Pseudoazurin - 1pmy

Geometric Task :

**Given two configurations of points in the three dimensional space,
find those rotations and translations of one of the point sets which produce “large” superimpositions of corresponding 3-D points.**

Remarks :

The superimposition pattern is not known a-priori – *pattern discovery* .

The matching recovered can be *inexact*.

We are looking not necessarily for the largest superimposition, since other matchings may have *biological meaning*.

Algorithmic Solution

File Edit View Go Bookmarks Options Directory Window Help

Location: http://silly6.math.tau.ac.il:8088/cgi-bin/backbone_match.cgi?PRI

Mail What's New? What's Cool? Destinations Net Search Welcome

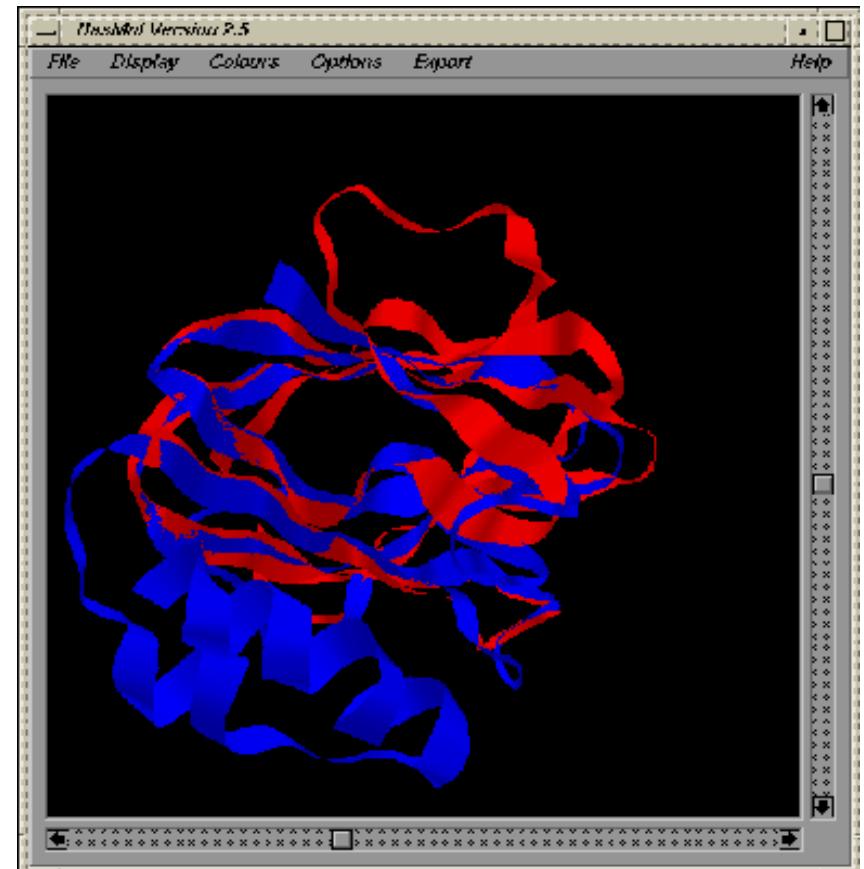
Results for matching 1PMY with 1AAJ

Results

#	Score	Match Size	RMS	Rotation			Translation		
Result 1	78.00	78	1.44	1.178	-0.059	-2.615	30.230	14.864	17.912
Result 2	61.00	61	2.05	-0.952	0.393	0.832	-1.717	7.031	-8.936
Result 3	60.00	60	1.82	1.999	-0.582	0.353	-4.668	19.664	30.651
Result 4	48.00	48	2.14	1.270	0.128	-2.818	31.964	9.542	13.826
Result 5	47.00	47	1.86	1.544	-1.136	-0.213	6.628	25.553	27.876
Result 6	45.00	45	1.82	-1.323	-0.208	0.324	-14.846	6.148	-1.035
Result 7	43.00	43	1.64	-1.133	0.232	0.291	-7.763	7.696	-11.996
Result 8	42.00	42	1.57	1.535	0.090	-3.050	26.773	7.332	12.549
Result 9	41.00	41	1.89	-2.113	0.924	-2.010	6.091	33.671	-16.447
Result 10	41.00	41	1.82	-2.066	0.590	-1.669	6.098	37.415	-6.752

• Rotations are given in radians for X-Y-Z axes. Rotating space around the X-axis, then around the Y-axis and finally around the Z-axis would give the required rotation.
• X-Y-Z translation coordinates are given in Angstrom units.

http://silly6.math.tau.ac.il:8088/cgi-bin/backbone_match.cgi?PRI

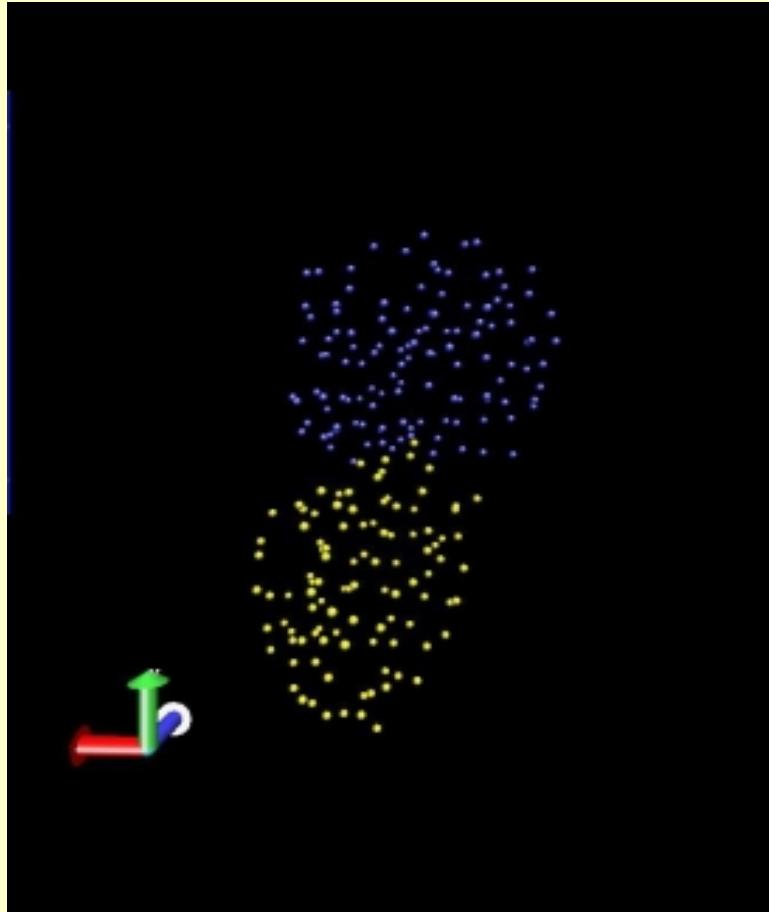


About 1 sec. Fischer, Nussinov, Wolfson ~ 1990.

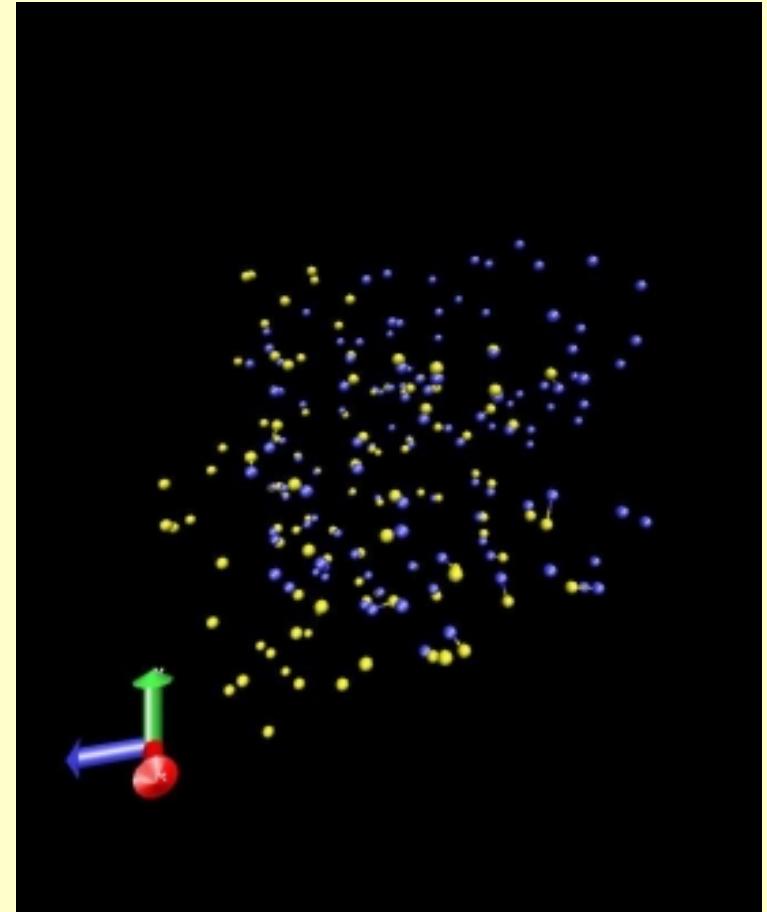
Applications

- Classification of protein databases by structure.
- Search of **partial and disconnected** structural patterns in large databases.
- Detection of structural pharmacophores in an ensemble of drugs.
- Comparison and detection of drug receptor *active sites*.

Geometric Matching task = Geometric Pattern Discovery

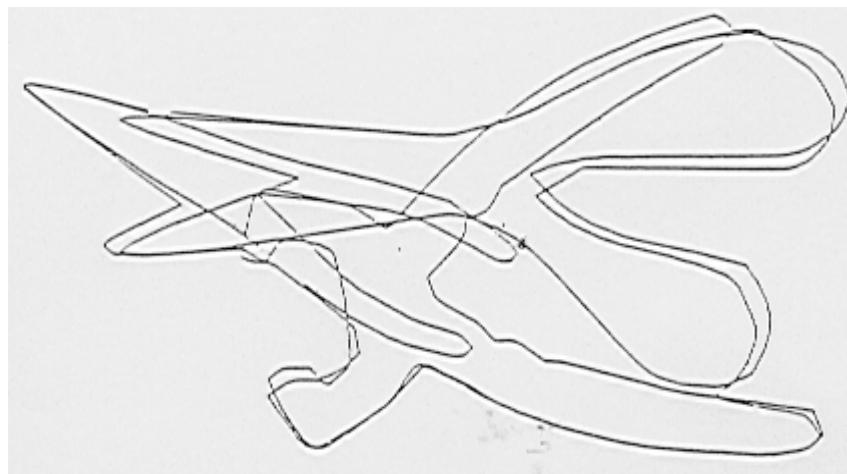
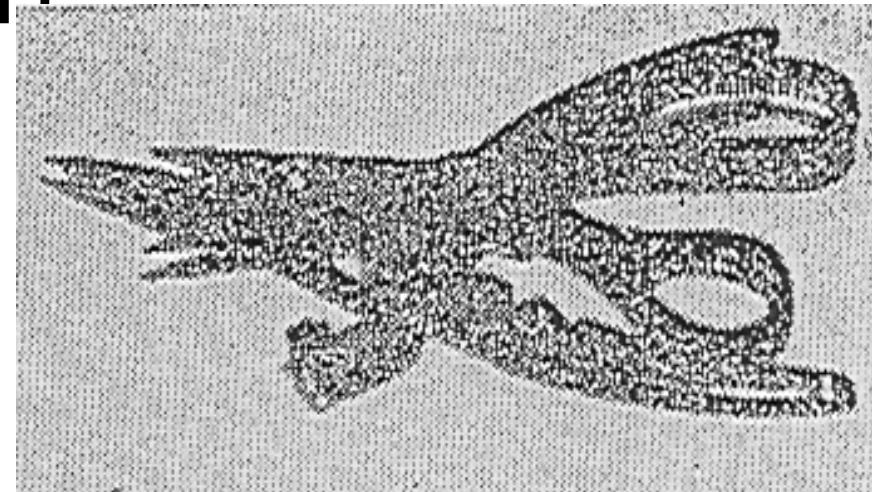
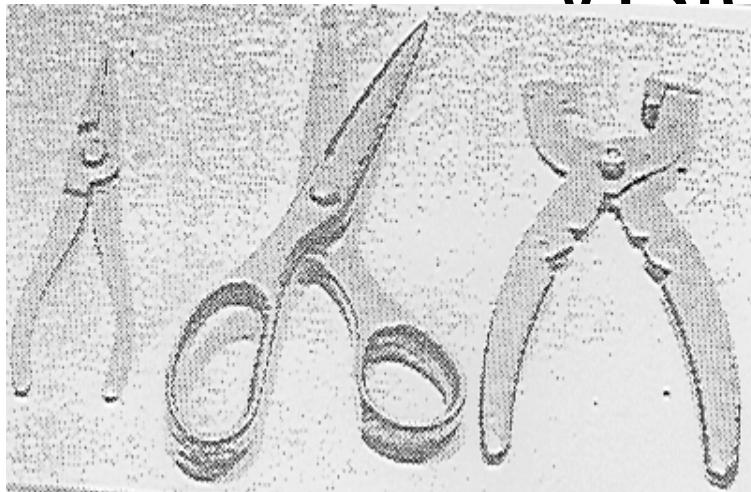


C_α constellations - before



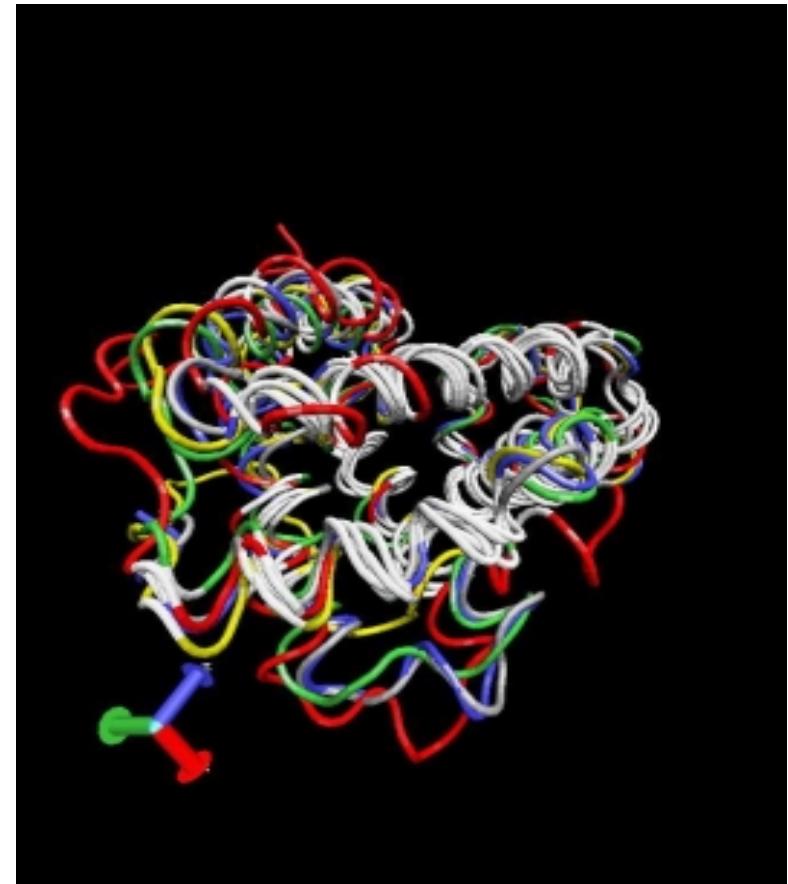
Superimposed constellations

Analogy with Object Recognition in Computer Vision



Wolfson, “Curve
Matching”, 1987.

Multiple Structural Alignment (Globin example)



Leibowitz, Fligelman, Nussinov, Wolfson, - ISMB'99 – Heidelberg.

Biomolecular Recognition - docking

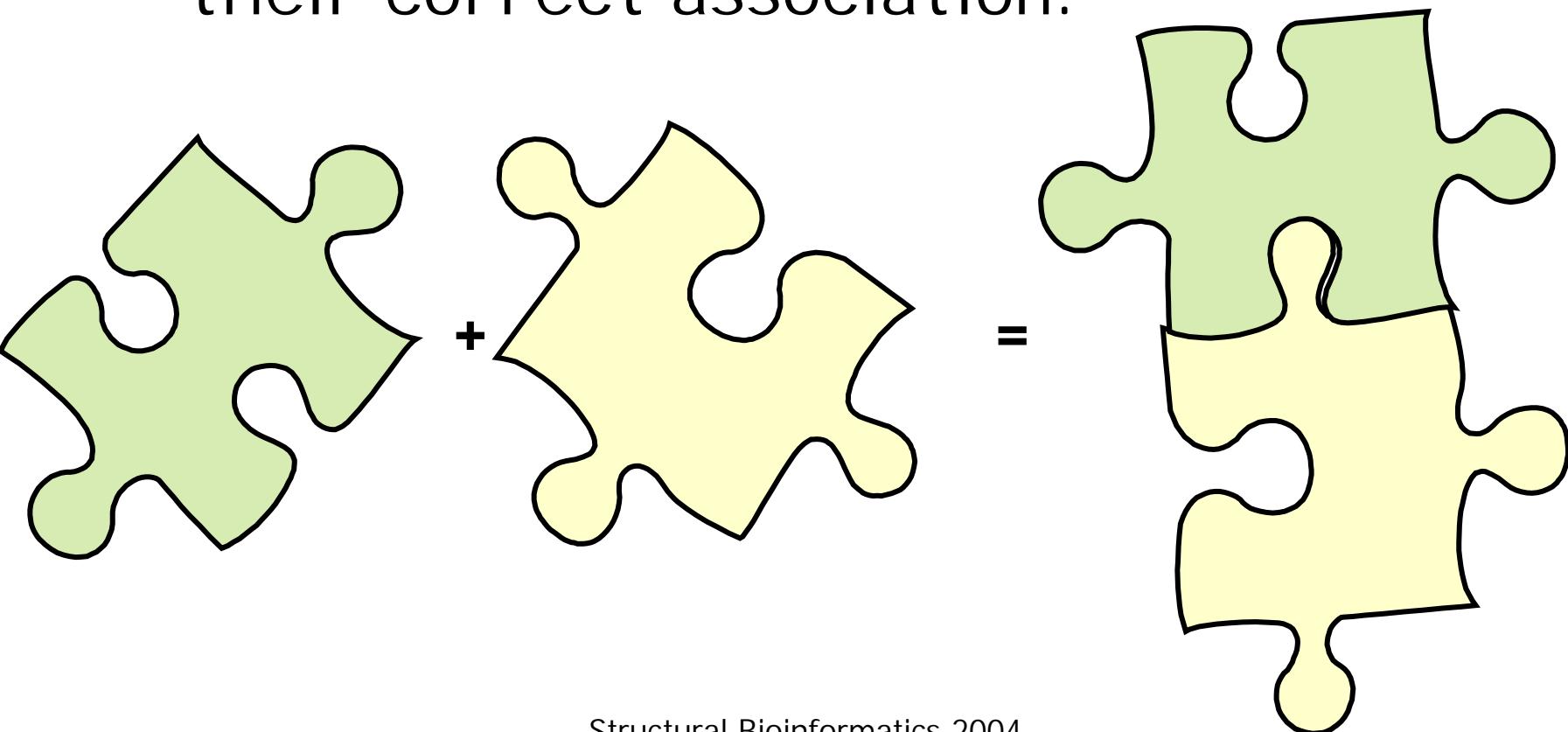
- Predict association of protein molecules.
- Predict binding of a protein molecule with a potential drug.
- Scan libraries of drugs to detect a suitable inhibitor for a target molecule.

Docking Algorithms

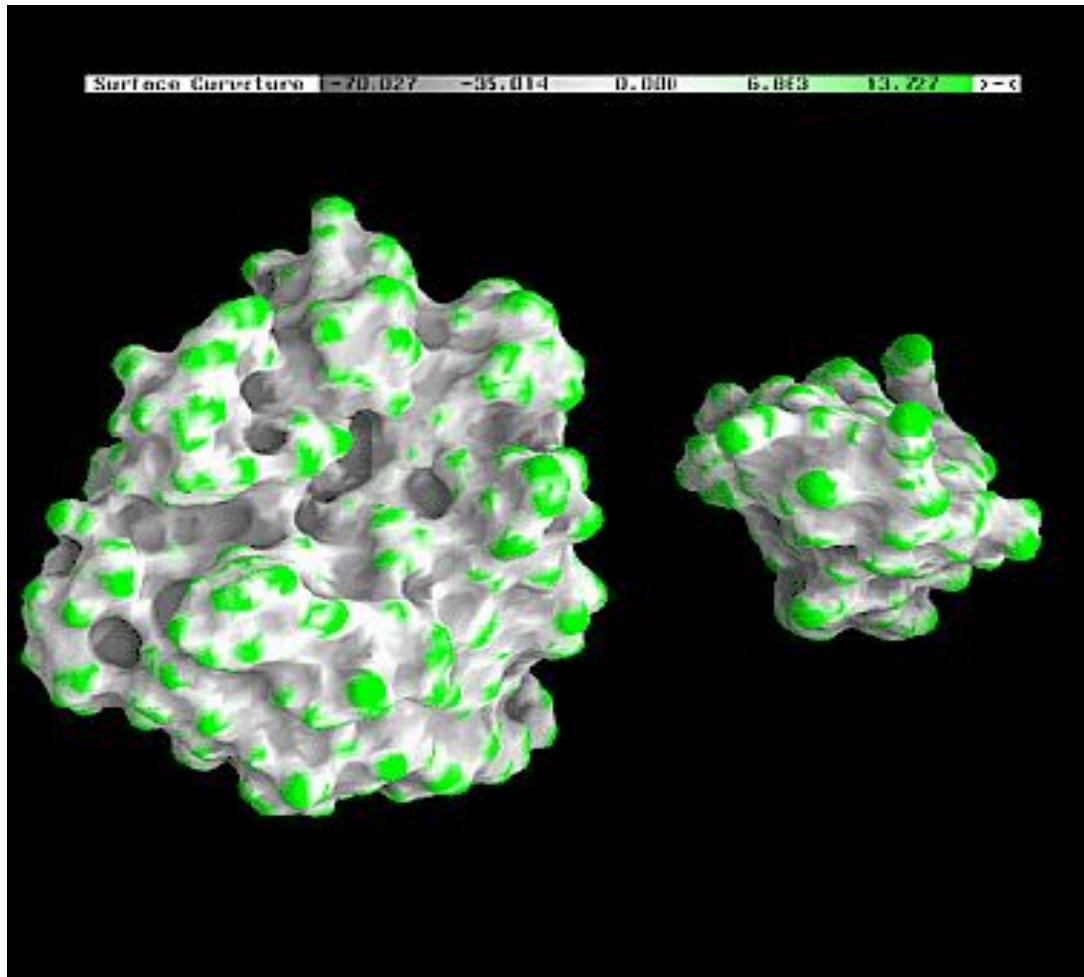
- Rigid receptor-ligand and protein-protein docking.
- Flexible receptor-ligand docking allowing a small number of hinges either in the ligand or the receptor.

Docking – Problem Definition

- Given a pair of molecules find their correct association:



Docking - Trypsin and BPTI



Docking – Relevance

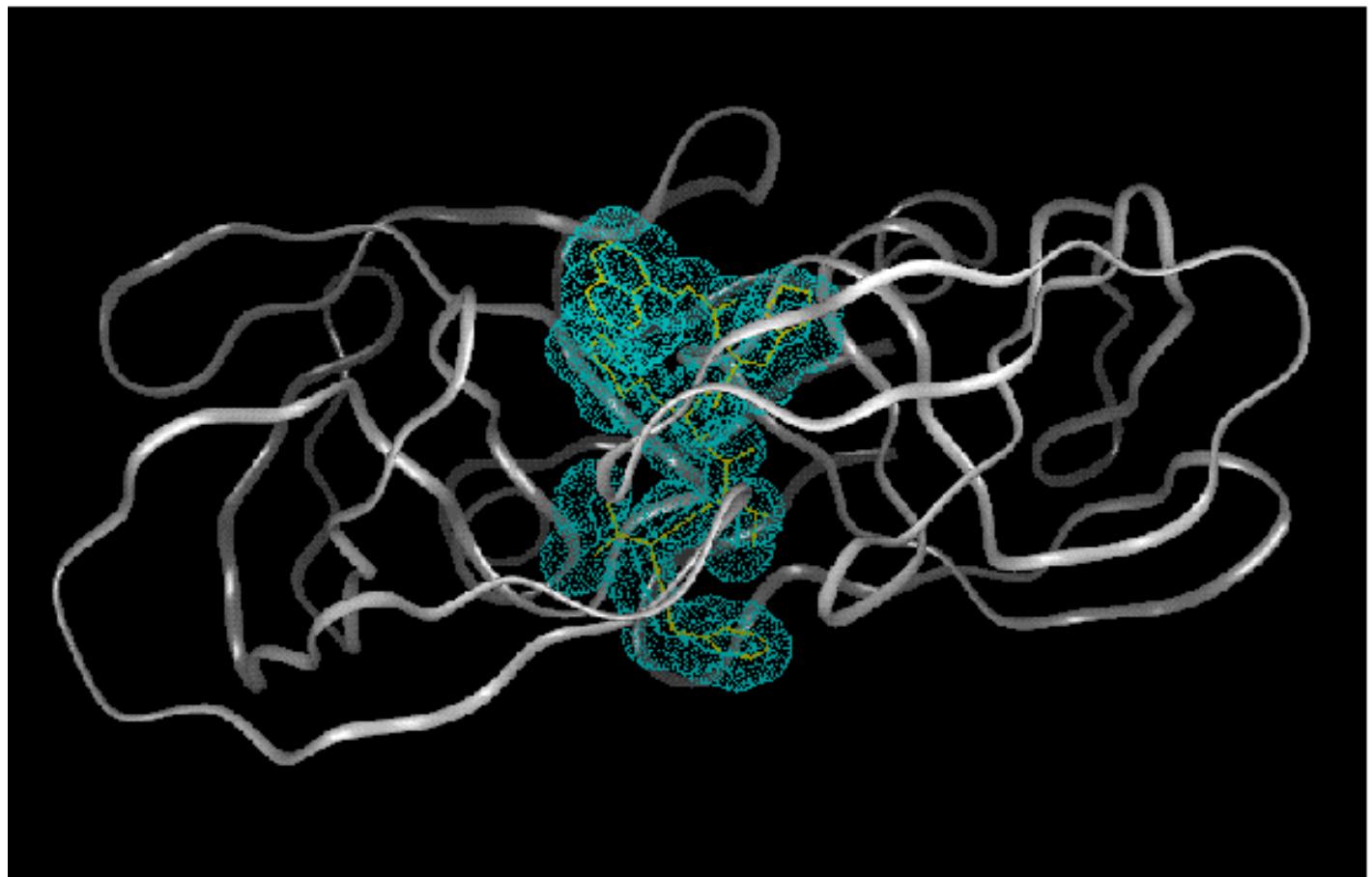
- Computer aided drug design – a new drug should fit the active site of a specific receptor.
- Understanding of the biochemical pathways - many reactions in the cell occur through interactions between the molecules.
- Crystallizing large complexes and finding their structure is difficult.

Flexible Docking Calmodulin with M13 ligand



Sandak, Nussinov, Wolfson - JCB 1998.

Flexible Docking HIV Protease Inhibitor



Sandak, Nussinov, Wolfson - CABIOS 1995.

Software Infrastructure

- Development of a software infrastructure for Geometric Computing in Molecular Biology.
- Object oriented, C++ library.
- Speed up development of new and re-usability of old software.
- Development of building blocks for fast testing of new ideas.

Cross - fertilization 1

- Analogous tasks appear in Computer Vision, Medical Imaging, Structural Bioinformatics, Target Recognition.
- Similar software and hardware can handle all of these *Geometric Computing* tasks - *method based cross fertilization.*

Cross - fertilization 2

- Bioinformatics brings together Computer Scientists, Molecular Biologists, Chemists etc. to tackle major problems in Computational Biology and Computer Assisted Drug Design - *task based cross-fertilization.*

Conclusions 1

- Molecular Biology and Biotechnology have entered a stage in which advanced algorithmic methods make the difference between theory and practice.
- Only true interdisciplinary collaboration among Computer and Life scientists can deliver **biologically relevant** computational techniques.

Conclusions 2

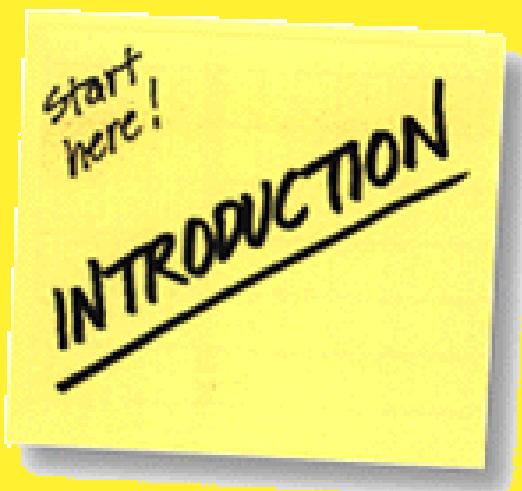
- The **b.c.** (before Computer Science) algorithms in Computational Biology/Biotechnology, which have been mostly developed by chemists and physicists, are analogous to the first generation CS algorithms. The current state-of-the-art of CS (~fifth generation) provides a quantum leap.

Sample of Topics to be covered

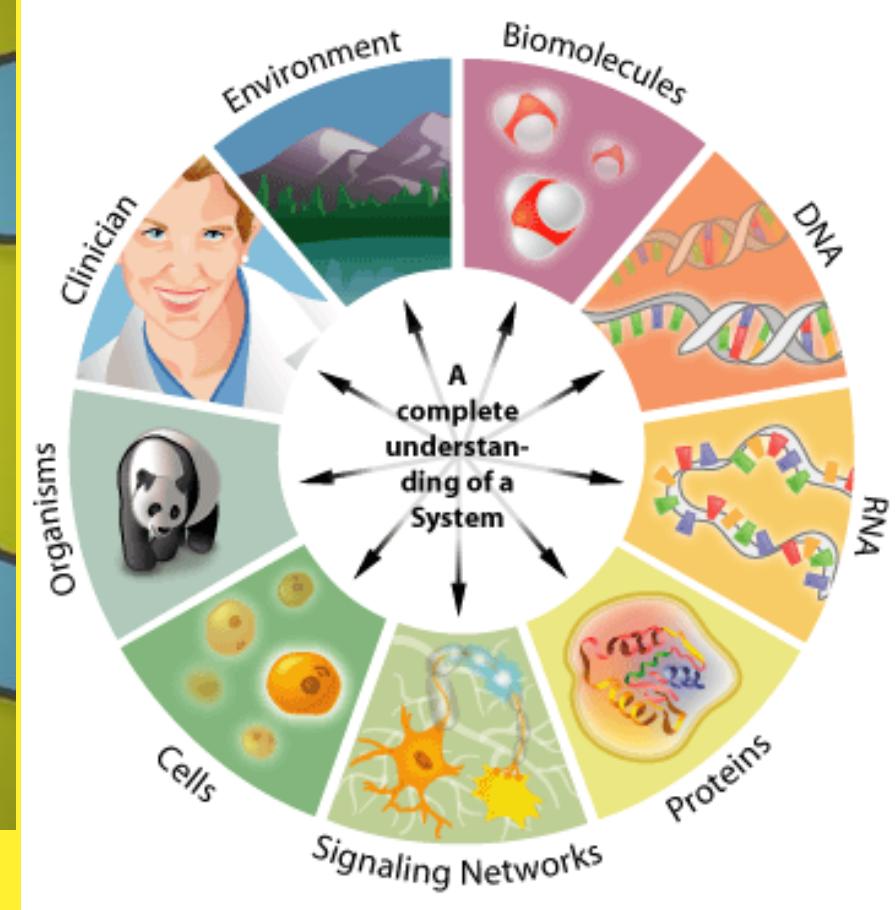
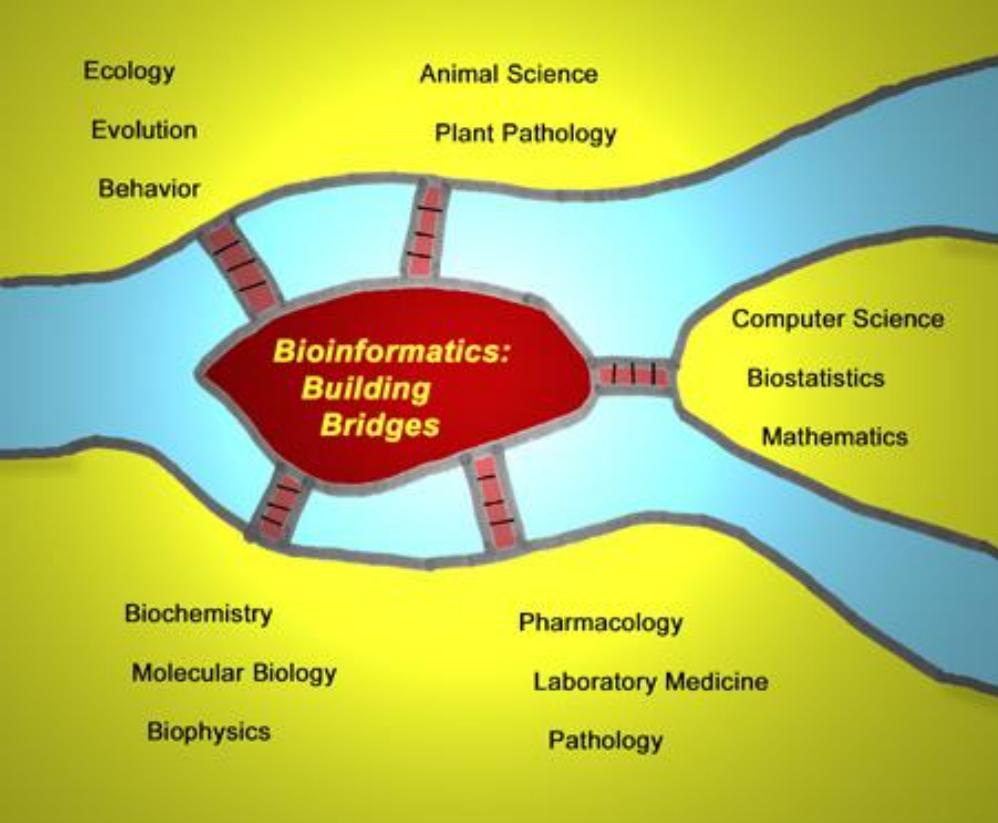
- Protein and DNA sequence alignment.
- Protein structural alignment and classification.
- Biomolecular recognition prediction – docking.
- Folding (homology modelling, threading, ab-initio).
- Distance Geometry for structure calculation from NMR data (?)
- Computer Assisted Structural Drug Design.

GRADING

- Exercises - 50%.
- Final (individual) Project, which involves heavy programming, based on the exercises – 50%.
- Most likely, all the students will get the same project assignment.
- The exact grading details will be supplied by the TA, Maxim Shatsky.



Introduction to Bioinformatics



What is Bioinformatics?

NIH – definitions

What is Bioinformatics? - Research, development,
and application of computational tools and on molecular
approaches for expanding the use of biological,
medical, behavioral, and health data, including the
means to acquire, store, organize, archive, analyze,
or visualize such data.

What is Computational Biology? The
development and application of analytical and
theoretical methods, mathematical modeling and
computational simulation techniques to the study of
biological, behavioral, and social data.

NSF – introduction

Large databases that can be accessed and analyzed with sophisticated tools have become central to biological research and education. The information content in the genomes of organisms, in the molecular dynamics of proteins, and in population dynamics, to name but a few areas, is enormous. Biologists are increasingly finding that the management of complex data sets is becoming a bottleneck for scientific advances. Therefore, **bioinformatics** is rapidly become a key technology in all fields of biology.

NSF – mission statement

The present bottlenecks in **bioinformatics** include the education of biologists in the use of advanced computing tools, the recruitment of computer scientists into this evolving field, the limited availability of developed databases of biological information, and the need for more efficient and intelligent search engines for complex databases.

NSF – mission statement

The present bottlenecks in **bioinformatics** include *the education of biologists in the use of advanced computing tools*, the recruitment of computer scientists into this evolving field, the limited availability of developed databases of biological information, and the need for more efficient and intelligent search engines for complex databases.

Molecular Bioinformatics

Molecular Bioinformatics involves the use of computational tools to discover new information in complex data sets (from the **one-dimensional** information of DNA through the **two-dimensional** information of RNA and the **three-dimensional** information of proteins, to the **four-dimensional** information of evolving living systems).

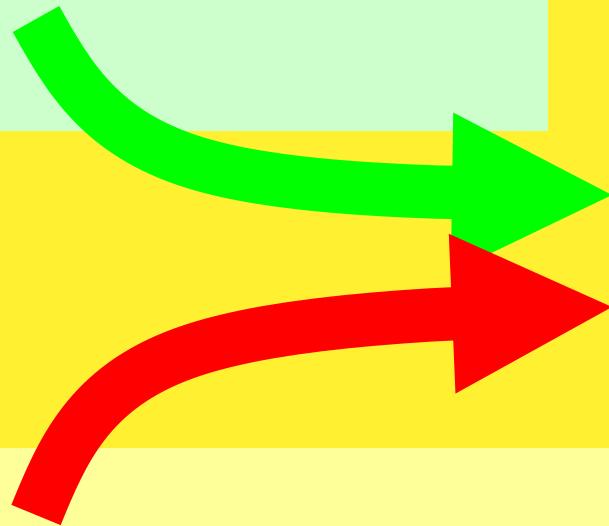
Bioinformatics (Oxford English Dictionary):

The branch of science concerned with information and information flow in biological systems, esp. the use of computational methods in genetics and genomics.

The field of science in which **biology**, **computer science** and **information technology** merge into a single discipline

Biologists

collect molecular data:
DNA & Protein sequences,
gene expression, etc.



Bioinformaticians

Study biological questions by
analyzing molecular data

Computer scientists

(+Mathematicians, Statisticians, etc.)
Develop tools, softwares, algorithms
to store and analyze the data.

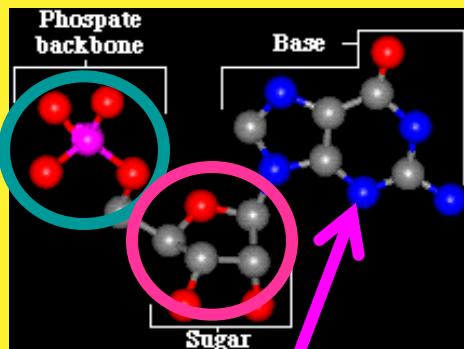
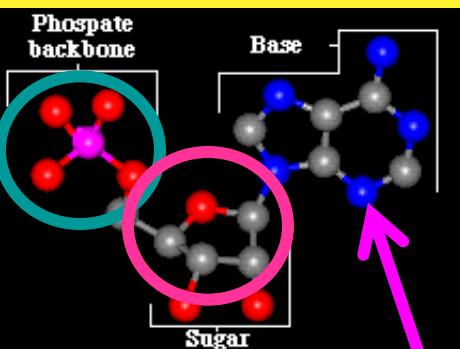
Some biological background....

A biologist



The hereditary information of all living organisms, with the exception of some viruses, is carried by **deoxyribonucleic acid (DNA)** molecules.

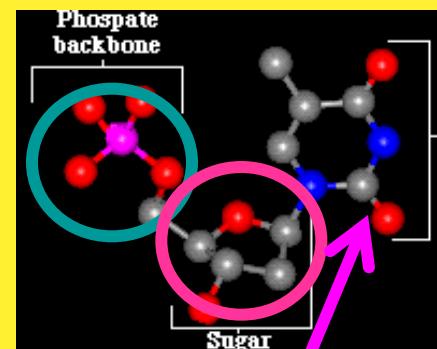
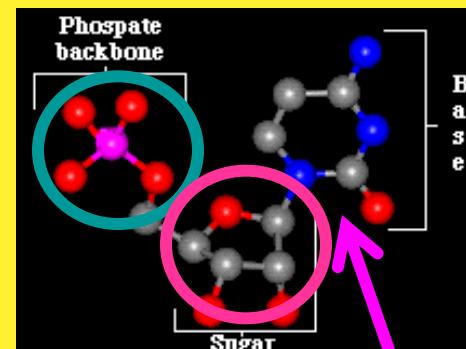
2 purines:



adenine (A)

two rings

2 pyrimidines:



guanine (G)

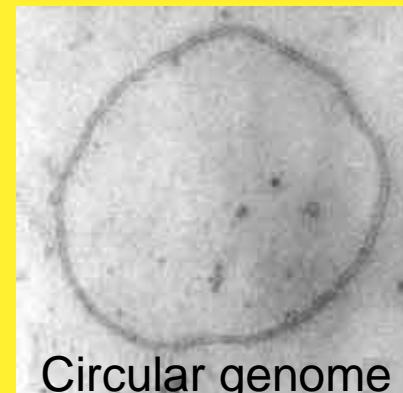
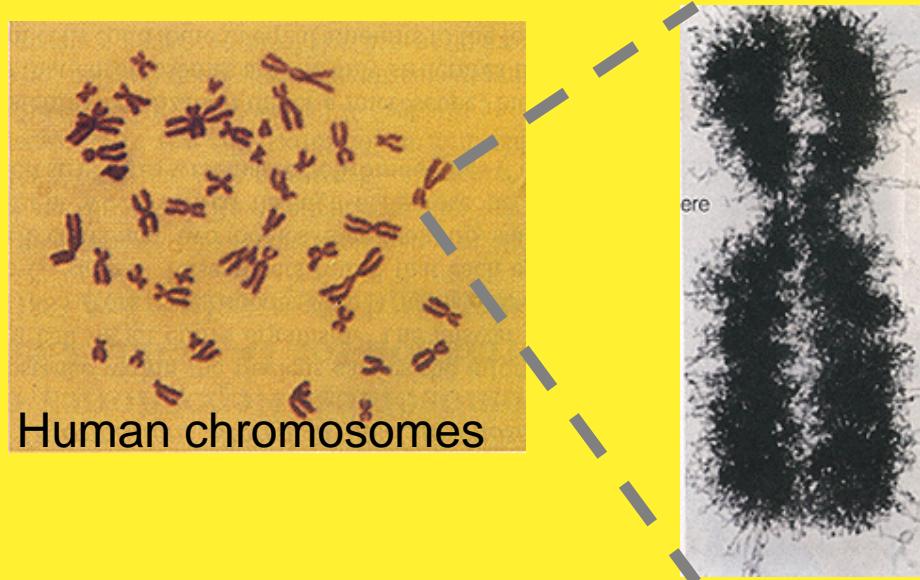
one ring

The entire complement of genetic material carried by an individual is called the **genome**.

Eukaryotes may have up to 3 subcellular genomes:

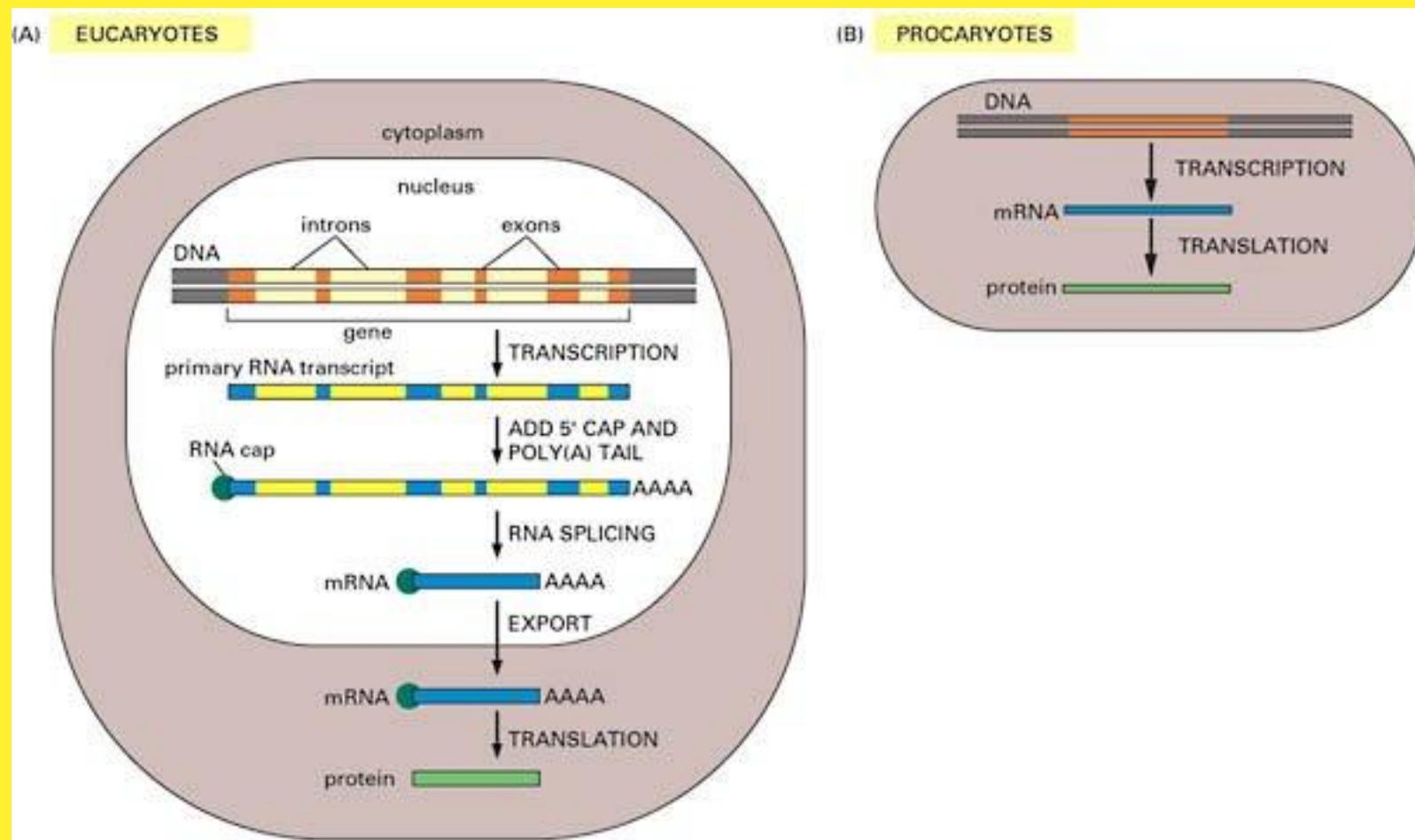
1. Nuclear
2. Mitochondrial
3. Plastid

Bacteria have either circular or linear genomes and may also carry plasmids



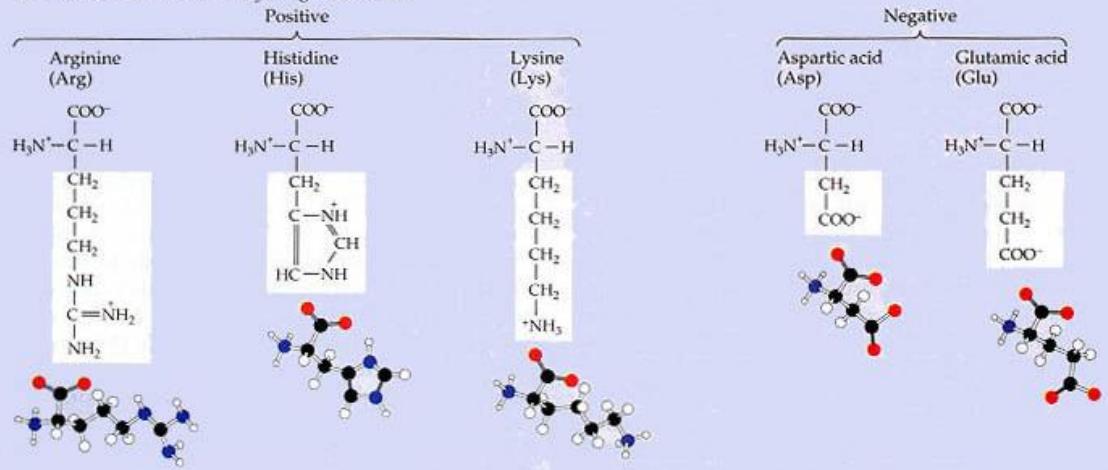
Central dogma: DNA makes RNA makes Protein

Modified dogma: DNA makes DNA and RNA, RNA makes DNA, RNA an Protein

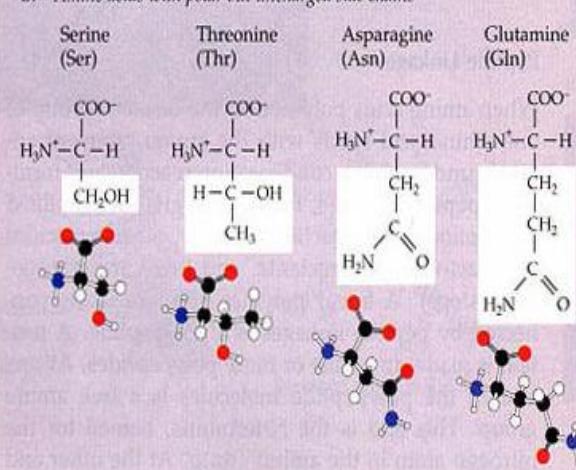


Amino acids - The protein building blocks

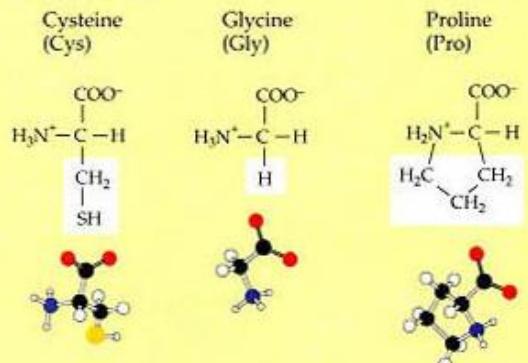
A. Amino acids with electrically charged side chains



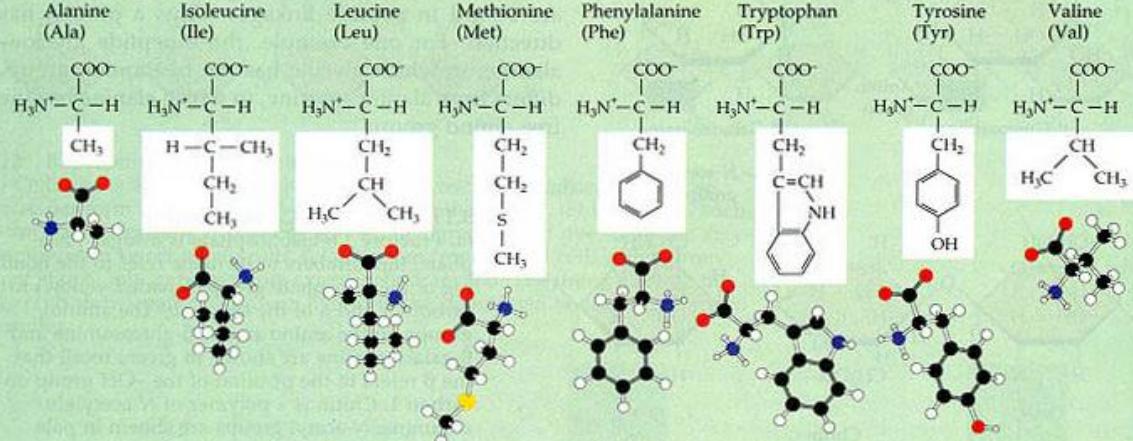
B. Amino acids with polar but uncharged side chains



C. Special cases

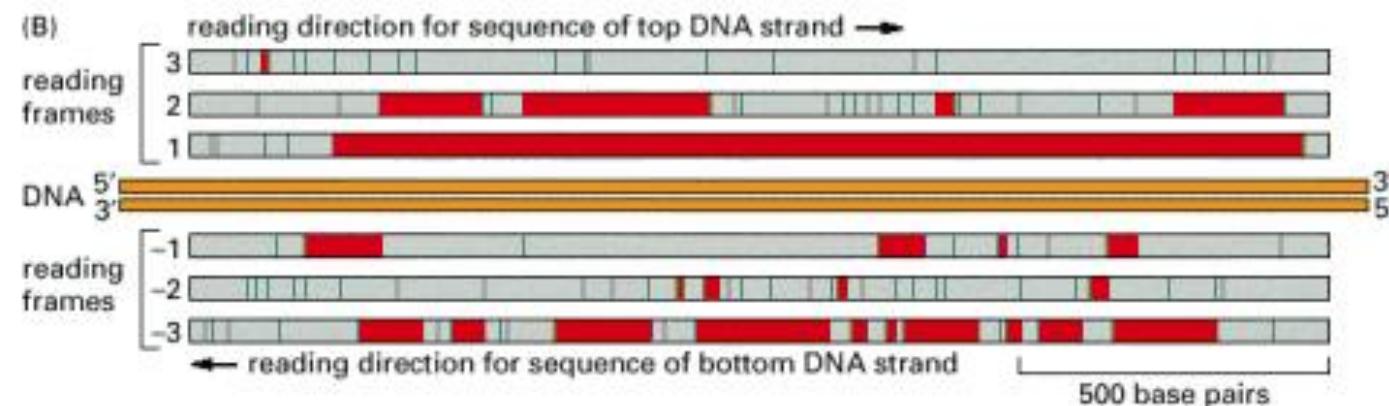
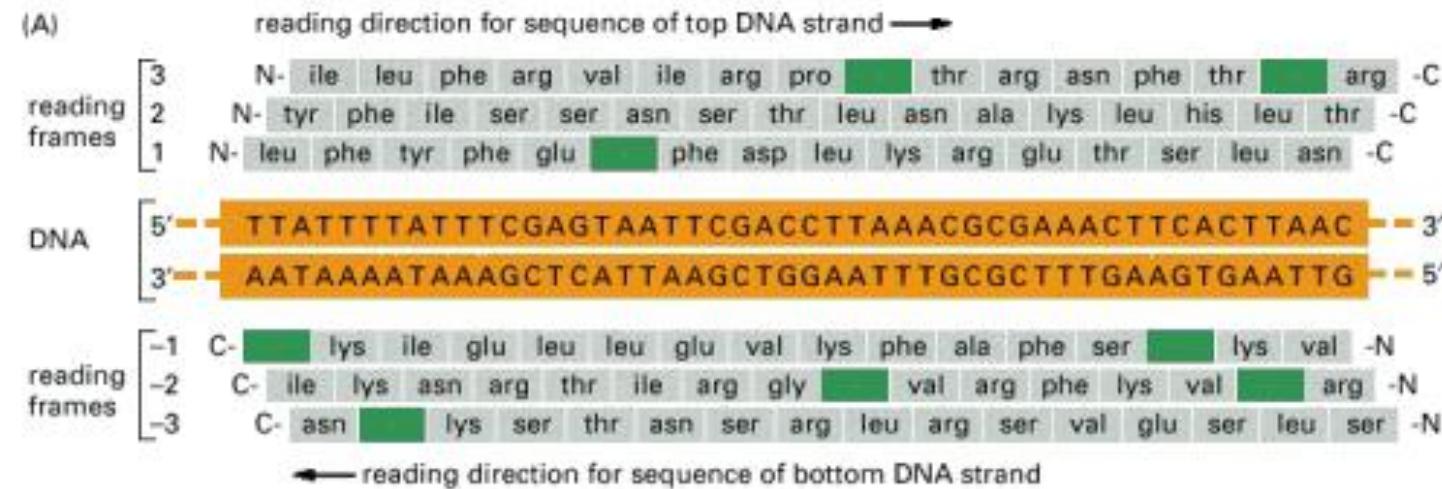


D. Amino acids with hydrophobic side chains



		Second letter						
		U	C	A	G			
First letter	U	UUU UUC UUA UUG	UCU UCC UCA UCG	UAU UAC UAA UAG	Tyr Stop Stop	UGU UGC UGA UGG	Cys Stop Trp	U C A G
	C	CUU CUC CUA CUG	CCU CCC CCA CCG	CAU CAC CAA CAG	His Gln	CGU CGC CGA CGG	Arg	U C A G
	A	AUU AUC AUA AUG	ACU ACC ACA ACG	AAU AAC AAA AAG	Asn Lys	AGU AGC AGA AGG	Ser Arg	U C A G
	G	GUU GUC GUA GUG	GCU GCC GCA GCG	GAU GAC GAA GAG	Asp Glu	GGU GGC GGA GGG	Gly	U C A G
Third letter								

Any region of the DNA sequence can, in principle, code for six different amino acid sequences, because any one of three different reading frames can be used to interpret each of the two strands.



Protein folding

A human Hemoglobin:



How does it all looks like on a computer monitor?

A cDNA sequence

>gi|14456711|ref|NM_000558.3| **Homo sapiens hemoglobin, alpha 1 (HBA1), mRNA**
ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTCTCCTGCCGACAAGACCAA
CGTCAAGGCCGCCTGGGTAAGGTGGCGCGACGCTGGCGAGTATGGTGGAGGCCCCTGGAG
AGGATGTT CCTGT CTT CCCC ACCACCA AGACCT ACTT CCCGC ACTTCGACCTGAGCCACGGCTCT
GCC CAGGTTAAGGCCACGGCAAGAAGGTGGCCACGCGCTGACCAACGCCGTGGCGACGTGG
ACGACATGCCAACCGCCTGTCCGCCCTGAGCGACCTGCACGCCACAAGCTTCGGGTGGACCCG
GTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTACCCCTGGCCGCCACCTCCCCGCCAGTT
ACCCCTGCGGTGCACGCCCTGGACAAGTTCTGGCTTCTGTGAGCACCGTGCTGACCTCCAAA
TACCGTTAAGCTGGAGCCTCGGTGGCCATGCTTCTTGGCCCTGGCCTCCCCCAGCCCCCTCCTC
CCCTTCCCTGCACCCGTACCCCCGTGGTCTTGAATAAGTCTGAGTGGCGGC

A cDNA sequence (reading frame)

>gi|14456711|ref|NM_000558.3| **Homo sapiens hemoglobin, alpha 1 (HBA1), mRNA**

ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACC **ATG**GTGCTGTCTCCTGCCGACAAGACCA
ACGTCAAGGCCGCCTGGGGTAAGGTCGGCGCGACGCTGGCGAGTATGGTGC GGAGGCCCTGGA
GAGGATGTT CCTGTCCTTCCCCACCAAGACCTACTTCCC GCACTTCGACCTGAGCCACGGCTC
TGCCCAGGTTAAGGGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGACGTG
GACGACATGCCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGACAAGCTTCGGGTGGACCC
GGTCAACTTCAAGCTCCTAAGCCACTGCCTGCTGGTACCCCTGGCCGCCACCTCCCCGCCGAGTT
CACCCCTGCGGTGCACGCCCTGGACAAAGTTCTGGCTTCTGTGAGCACC GTGCTGACCTCCAA
ATACCGT TAAGCTGGAGCCTCGGTGCCATGCTTCTTCCCCCTGGGCCTCCCCCAGCCCCCTCC
TCCCCCTCCTGCACCCGTACCCCCGTGGTCTTGAAATAAGTCTGAGTGGCGGC

A protein sequence

>gi|4504347|ref|NP_000549.1| **alpha 1 globin [Homo sapiens]**

MVLSPADKTNVKA AWGKVGA HAGEYGA EALERMF LS FPTTKTYFP HF DLSHGSAQVK GHG KKVA
DAL TNAVAHVDDMPNALSALSDLHAHKL RVDPVNFKLLSHCLLVTAAHLPAEFTP AVHASLDKFLA
SVSTVLTSKYR

And, a whole genome...

ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTCTGCCGACAAGACCAACGTC
AAGGCCGCCTGGGTAAGGTCGGCGCGACGCTGGCGAGTATGGTGC GGAGGCCCTGGAGAGGGATGTT
CTGTCCTTCCCCACCAAGACCTACTTCCCCTGACCTGAGCCACGGCTCTGCCCAAGGTTAAGG
GCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGCACGTGGACGACATGCCAACGCG
CTGTCGCCCTGAGCGACCTGCACGCGACAAGCTTCGGGTTGGACCCGGTCAACTCAAGCTCCTAACGCC
ACTGCCTGCTGGTGA CCCTGGCCGCCACCTCCCCGCCAGTTCACCCCTGC GGTCACGCCCTCCCTGG
ACAAGTTCCCTGGCTTCTGTGAGCACC GTGCTGACCTCCAAATACCGTTAAGCTGGAGCCTCGGTGGCCATG
CTTCTGCCCTGGCCTCCCCCAGGCCCTCCCTGCACCCGTACCCCGTGGTCTTGAAT
AAAGTCTGAGTGGCGGC ACTCTTCTGGTCCCCACAGACTCAGAGAGAACCCACCATGGTGCTGTCCCT
GCCGACAAGACCAACGTCAAGGCCCTGGGTAAGGTCGGCGCGACGCTGGCGAGTATGGTGC GGAG
GCCCTGGAGAGGATGTT CCTGTCCCTCCCCACCAAGACCTACTTCCCGACTTCGACCTGAGCCACG
GCTCTGCCCAAGGTTAAGGCCACGGCAAGAAGGTGGCCGACGCGCTGACCAACGCCGTGGCGACGTG
GACGACATGCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGACAAGCTTCGGGTTGGACCCGGTC
AACTTCAAGCTCCTAACGCCACTGCCTGCTGGTACCCCTGGCCCCACCTCCCCGCCAGTTAACCCCTG
CGGTGCACGCCCTGGACAAGTT CCTGGCTTCTGCCCCCTGGCCTCCCCCAGGCCCTCCCTGCACCCGT
GGAGCCTCGGTGGCCATGCTTCTGCCCCCTGGCCTCCCCCAGGCCCTCCCTGCACCCGT
ACCCCGTGGTCTTGAATAAGTCTGAGTGGCGGC ACTCTTCTGGTCCCCACAGACTCAGAGAGAAC
CACCATGGTGCTGTCTCCTGCCGACAAGACCAACGTCAAGGCCCTGGGTAAGGTCGGCGCGACGC
TGGCGAGTATGGTGC GGAGGCCCTGGAGAGGATGTT CCTGTCCCTCCCCACCAAGACCTACTTCCCG
CACTTCGACCTGAGCCACGGCTCTGCCCAAGGTTAAGGCCACGGCAAGAAGGTGGCCGACGCGCTGACC
AACGCCGTGGCGCACGTGGACGACATGCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGACAAG
CTTCGGGTGGACCCGGTCAACTTCAAGCTCCTAACGCCACTGCCTGCTGGTACCCCTGGCCGCCACCTCC
CCGCCAGTTACCCCTGCCGTGCACGCCCTGGACAAGTT CCTGGCTTCTGTGAGCACC GTGCTGAC
CTCCAAATACCGTTAAGCTGGAGCCTCGGTGGCCATGCTTCTGCCCTGGCCTCCCCCAGGCCCTCC
TCCCCTCCCTGCACCCGTACCCCGTGGTCTTGAATAAGTCTGAGTGGCGGCCGTGGCGACGTG
GACGACATGCCAACGCGCTGTCCGCCCTGAGCGACCTGCACGCGACAAGCTTCGGGTTGGACCCGGTC
AACTTCAAGCTCCTAACGCCACTGCCTGCTGGTACCCCTGGCCGCCACCTCCCCGCCAGTTCAC
CGGTGCACGCCCTCCCTGGACAAGTT CCTGGCTTCTGTGAGCACC GTGCTGACCTCCAAATACCGTTAAGCT
GGAGCCTCGGTGGCCATGCTTCTGCCCTGGCCTCCCCCAGGCCCTCCCTGGCCTGCCGT

How big are whole genomes?



E. coli 4.6×10^6 nucleotides
– Approx. 4,000 genes



Yeast 15×10^6 nucleotides
– Approx. 6,000 genes

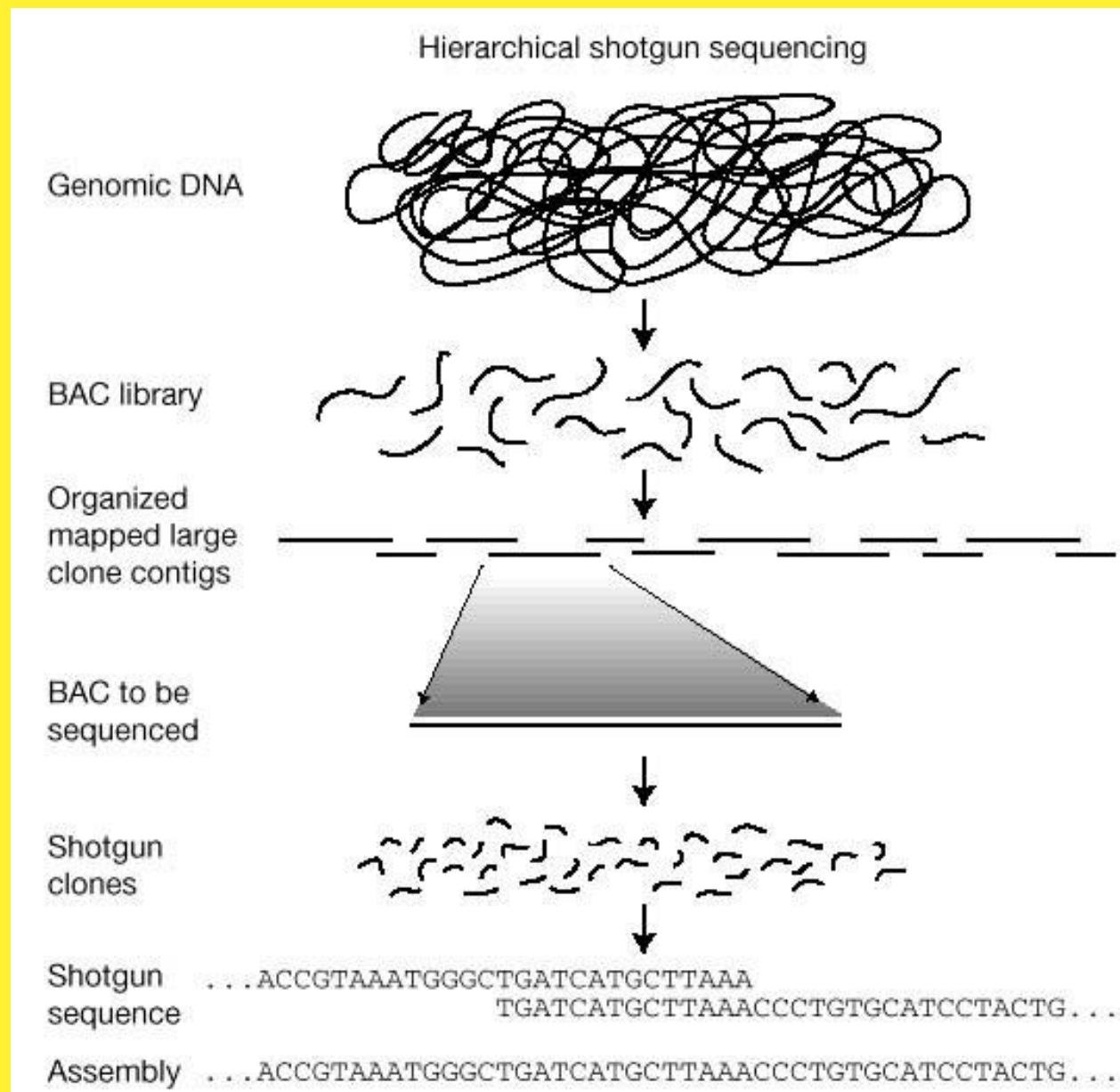


Human 3×10^9 nucleotides
– Approx. 30,000 genes

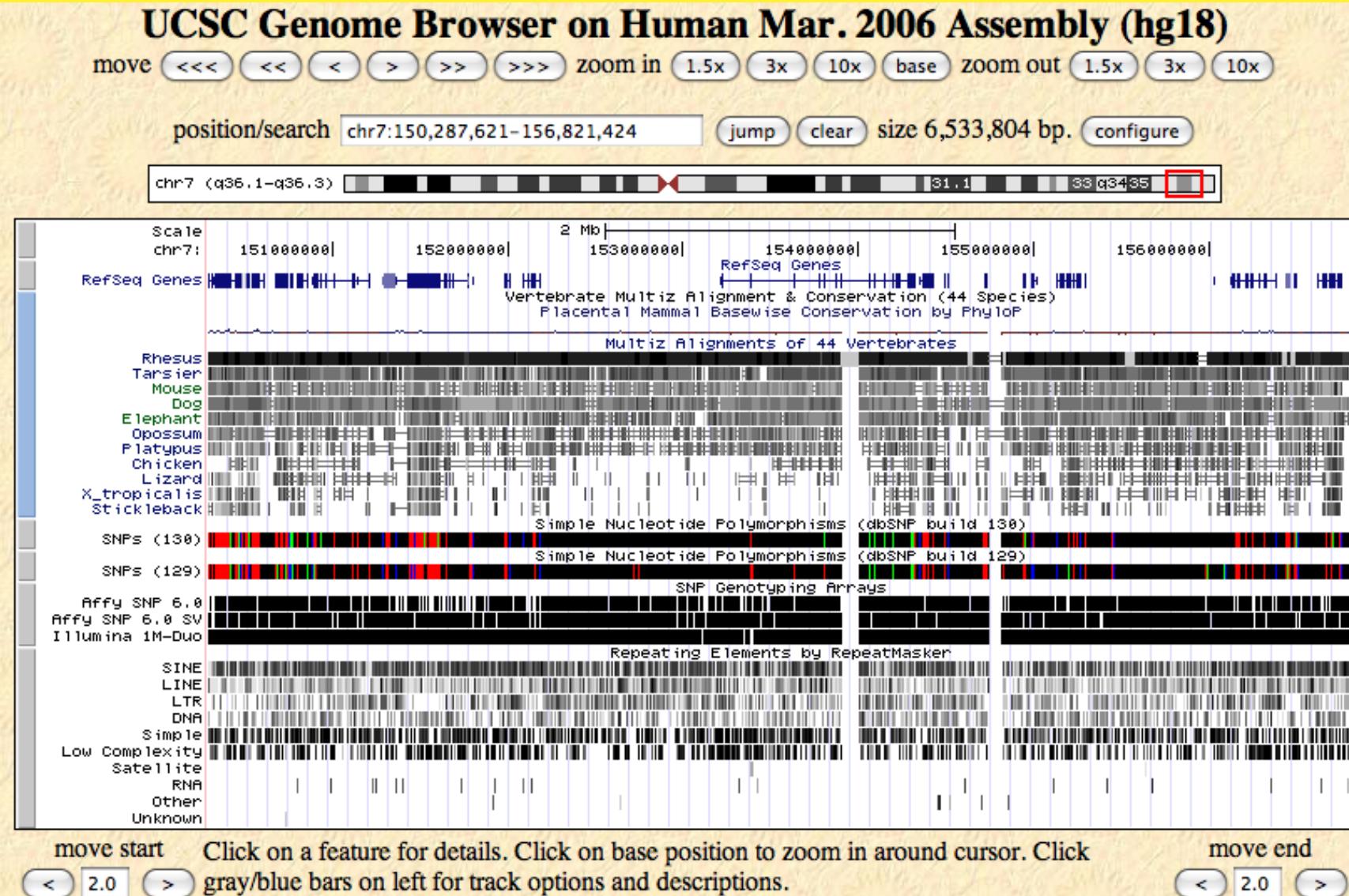
Smallest human chromosome 50×10^6 nucleotides

What do we actually do with bioinformatics?

Sequence assembly

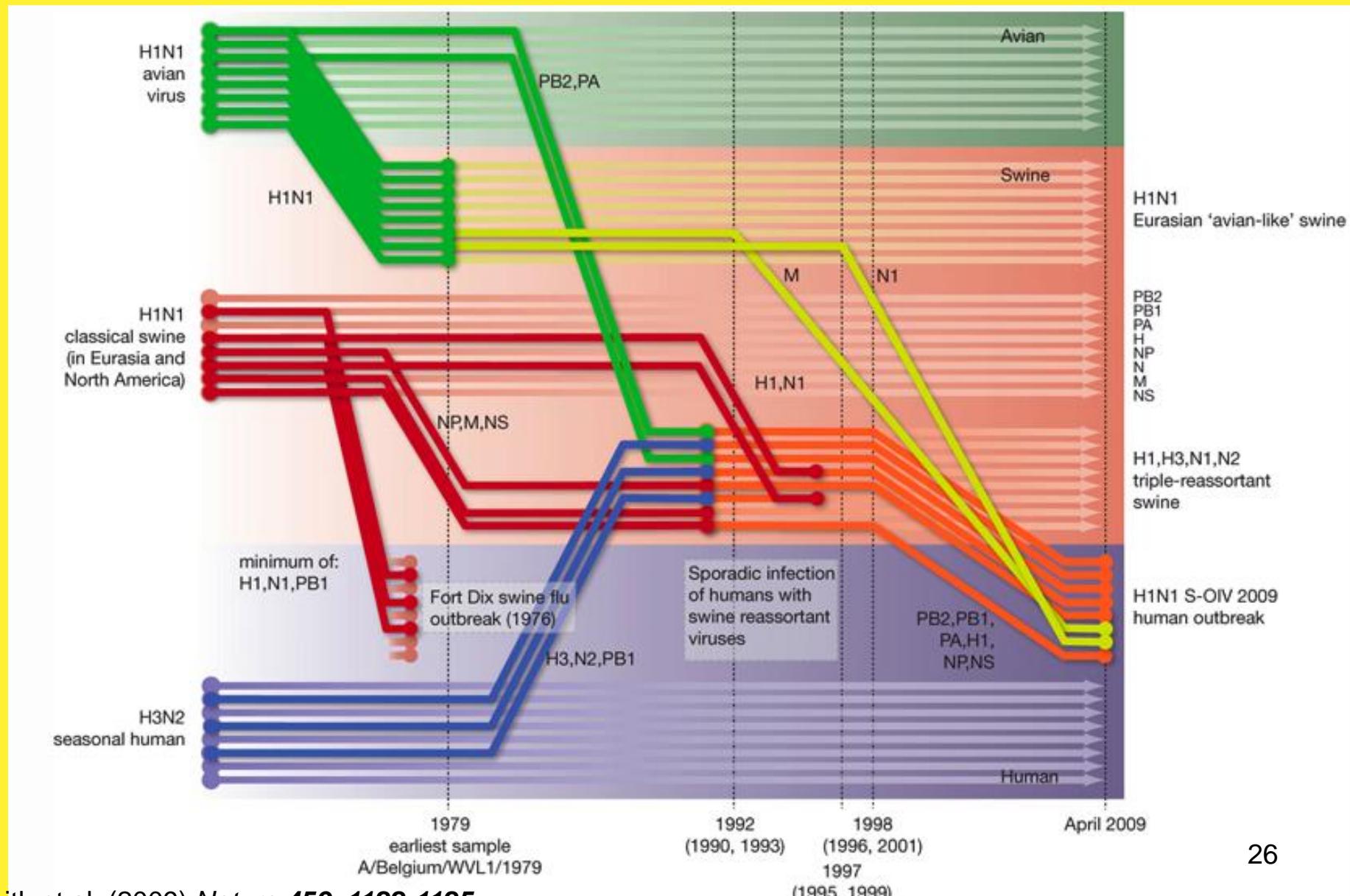


Genome annotation

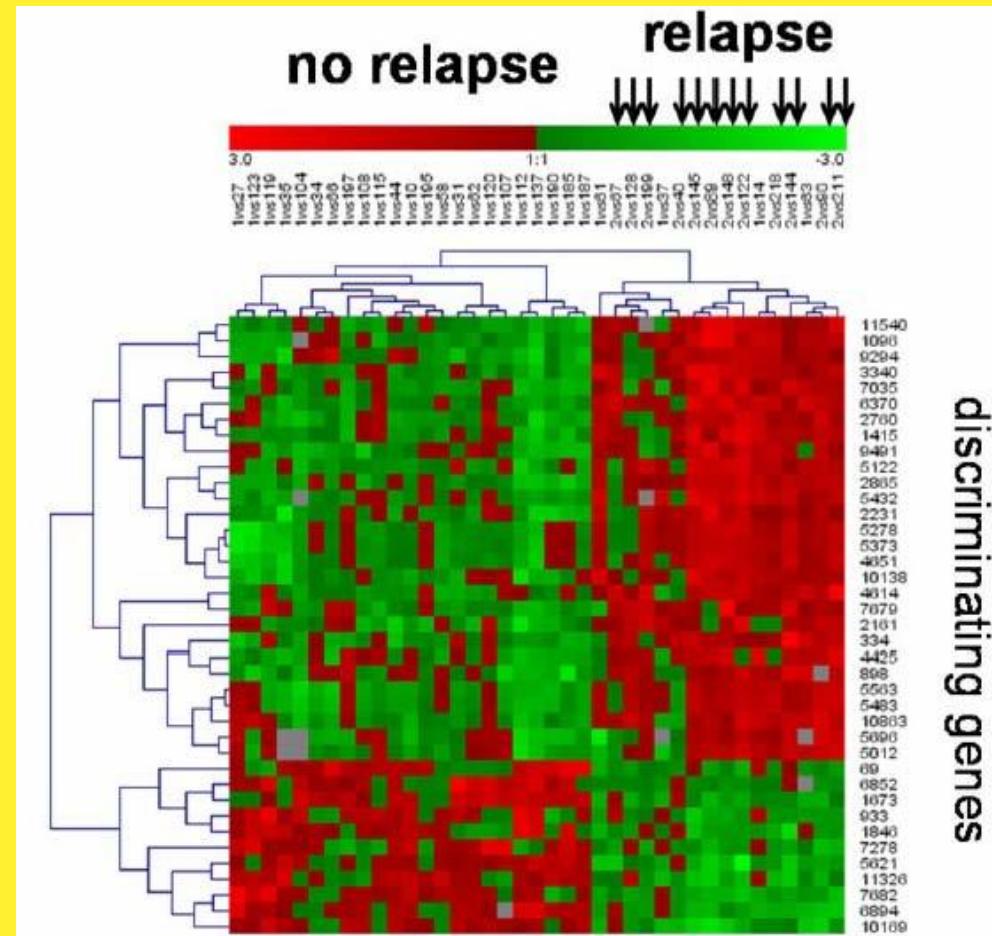


Molecular evolution

Origins and evolutionary genomics of the 2009 swine-origin H1N1 influenza A epidemic

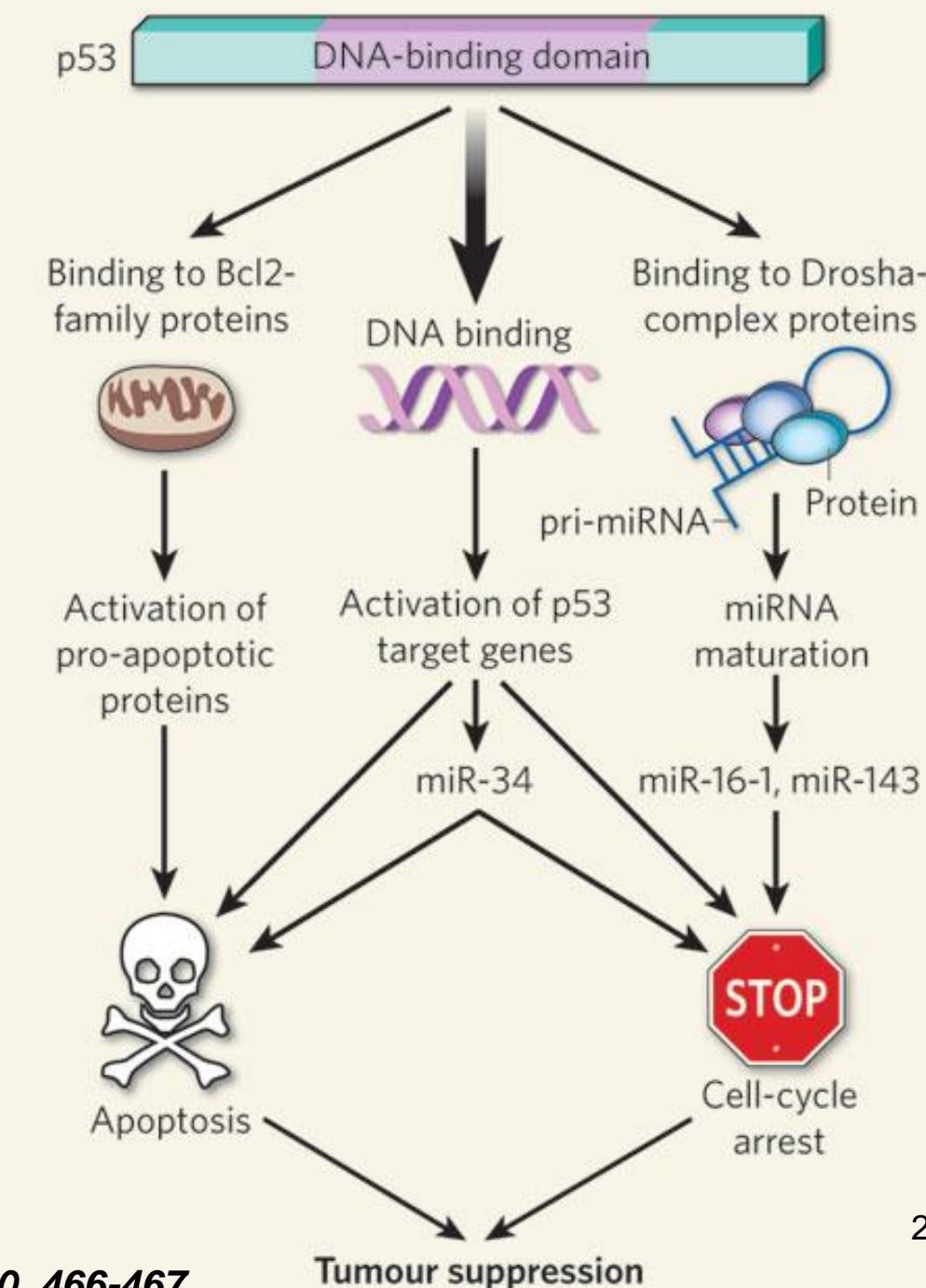


Analysis of gene expression

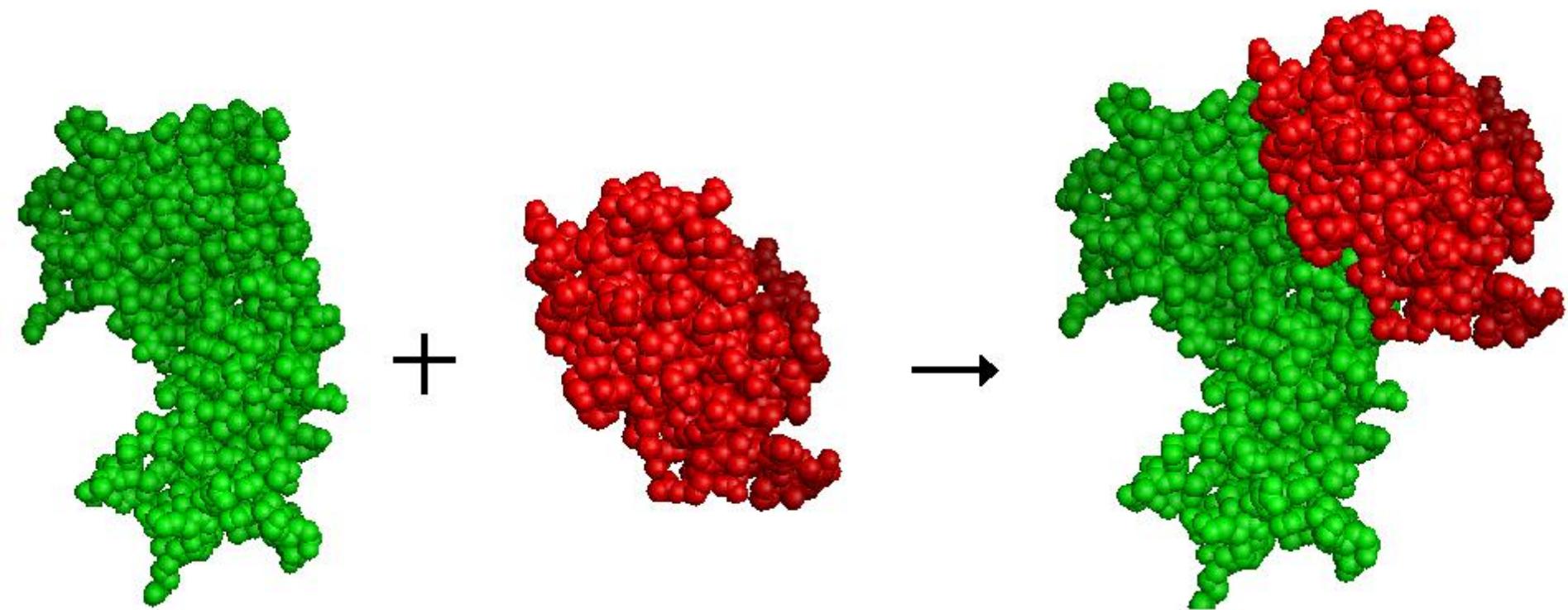


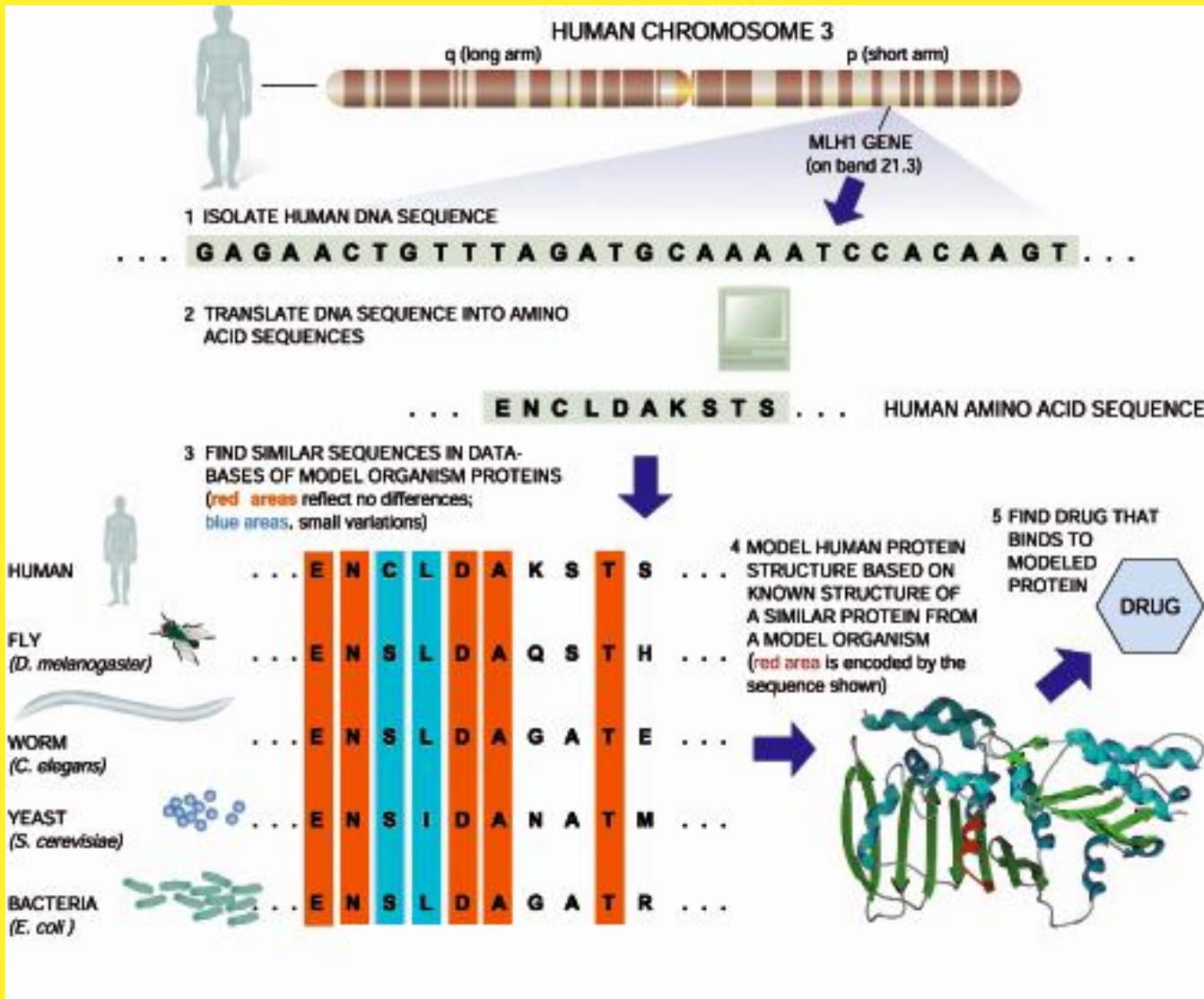
Gene expression profile of relapsing versus non-relapsing Wilms tumors.
A set of 39 genes discriminates between the two classes of tumors.

Analysis of regulation



Protein structure prediction
Protein docking





From DNA to Genome

Watson and Crick
DNA model

Sequence
alignment

PDB (Protein
Data Bank)

GenBank
database

195
5

196
5

197
5

198
5

196
0

197
0

198
0

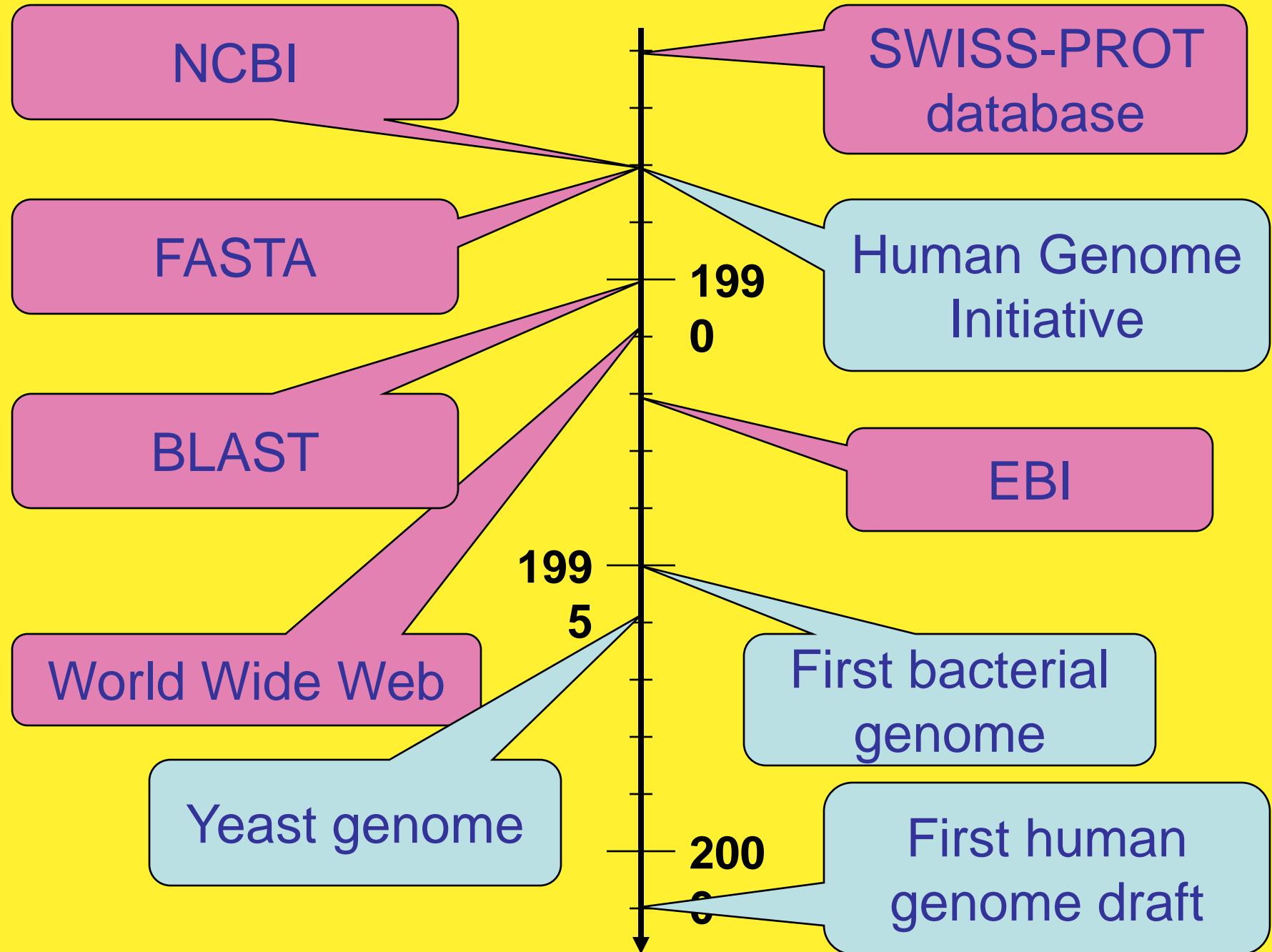
Sanger
sequences
insulin protein

Dayhoff's Atlas

ARPANET
(early
Internet)

Sanger dideoxy
DNA sequencing

PCR
(Polymerase
Chain Reaction)



Origin of bioinformatics and biological databases:

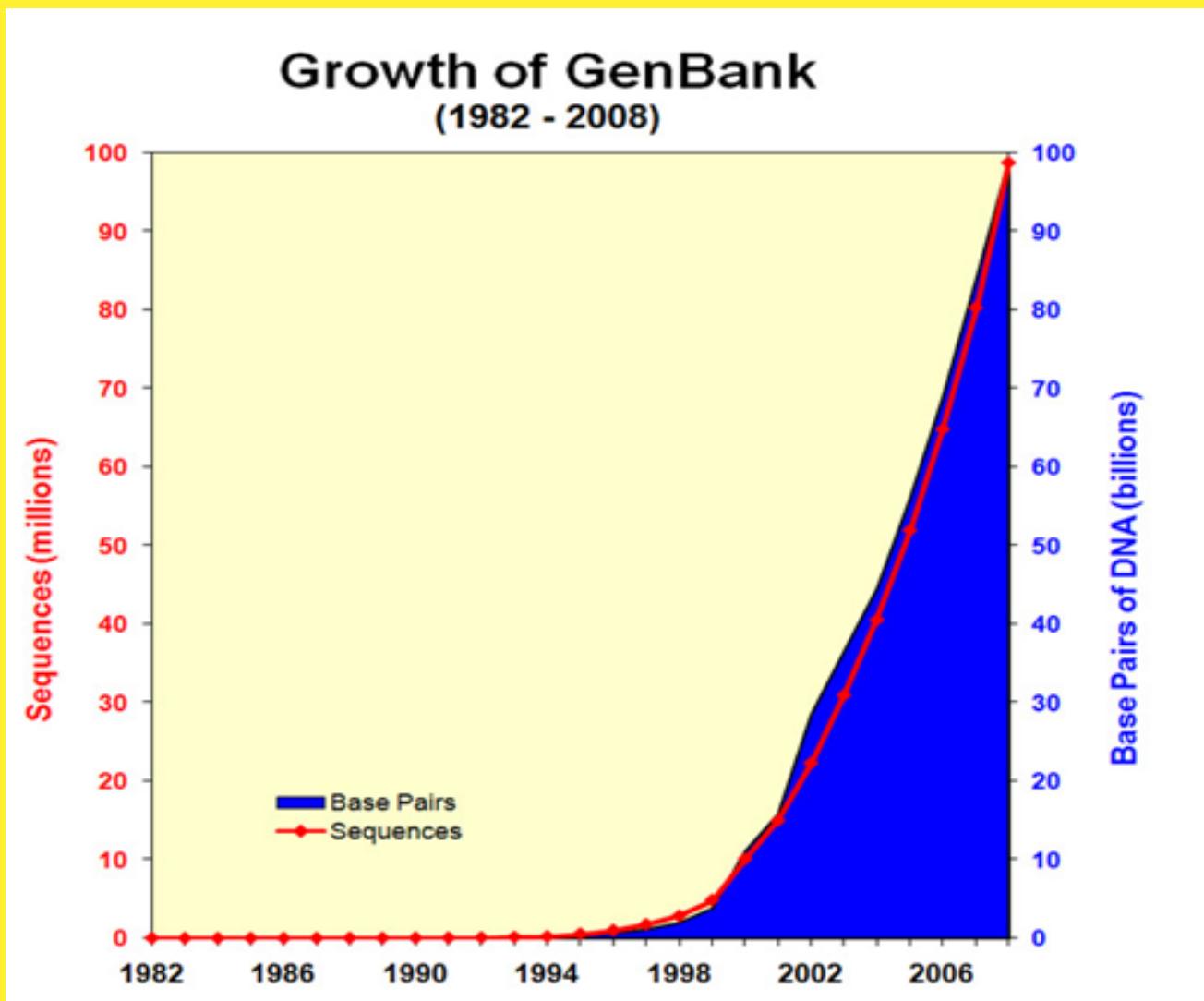
The first protein sequence reported was that of bovine insulin in **1956**, consisting of 51 residues.

Nearly a decade later, the first nucleic acid sequence was reported, that of yeast tRNA^{alanine} with 77 bases.

In 1965, Dayhoff gathered all the available sequence data to create the first bioinformatic database (*Atlas of Protein Sequence and Structure*).

The Protein DataBank followed in 1972 with a collection of ten X-ray crystallographic protein structures. The SWISSPROT protein sequence database began in 1987.

Nucleotides



Complete Genomes

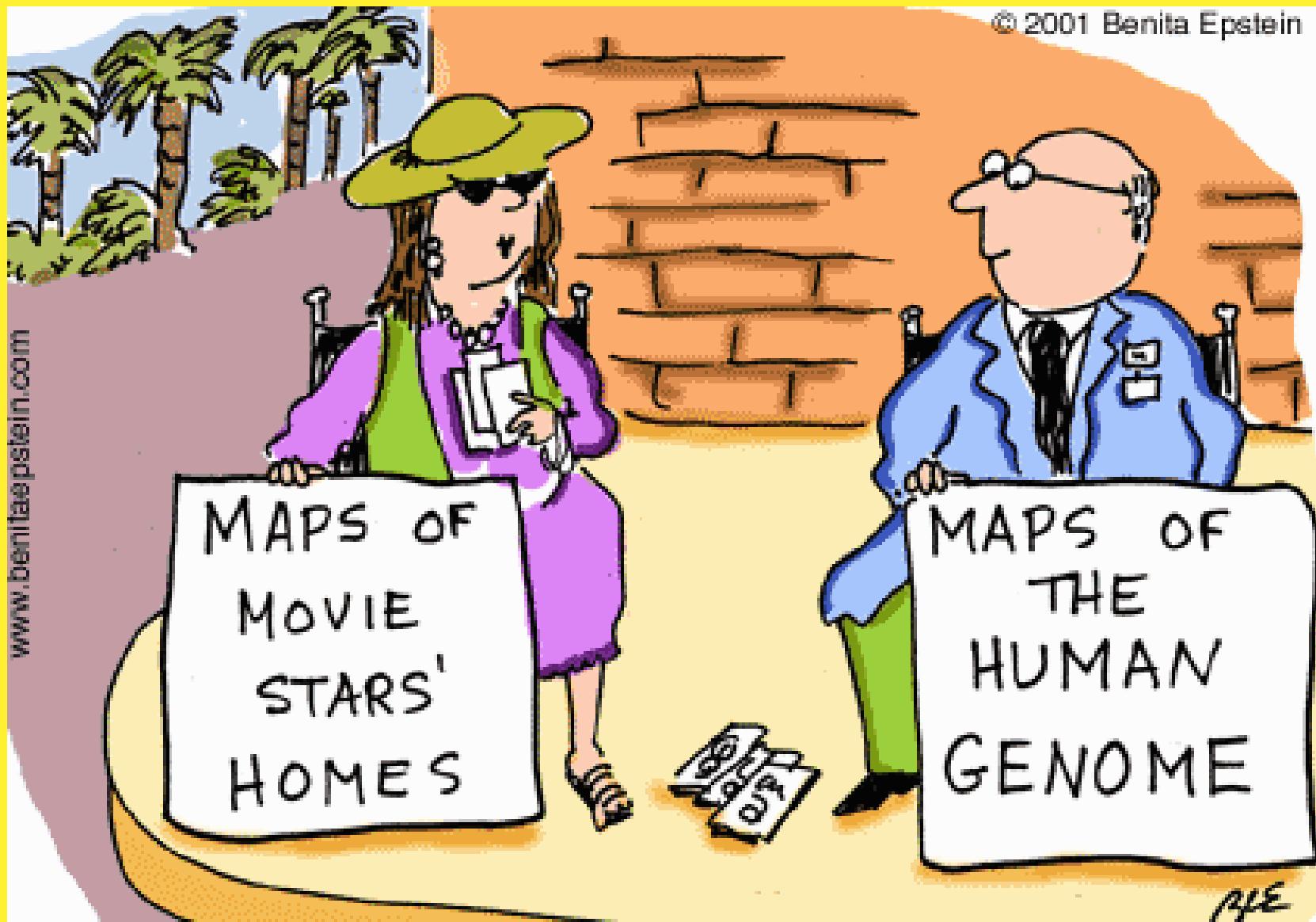
as of August 2011:

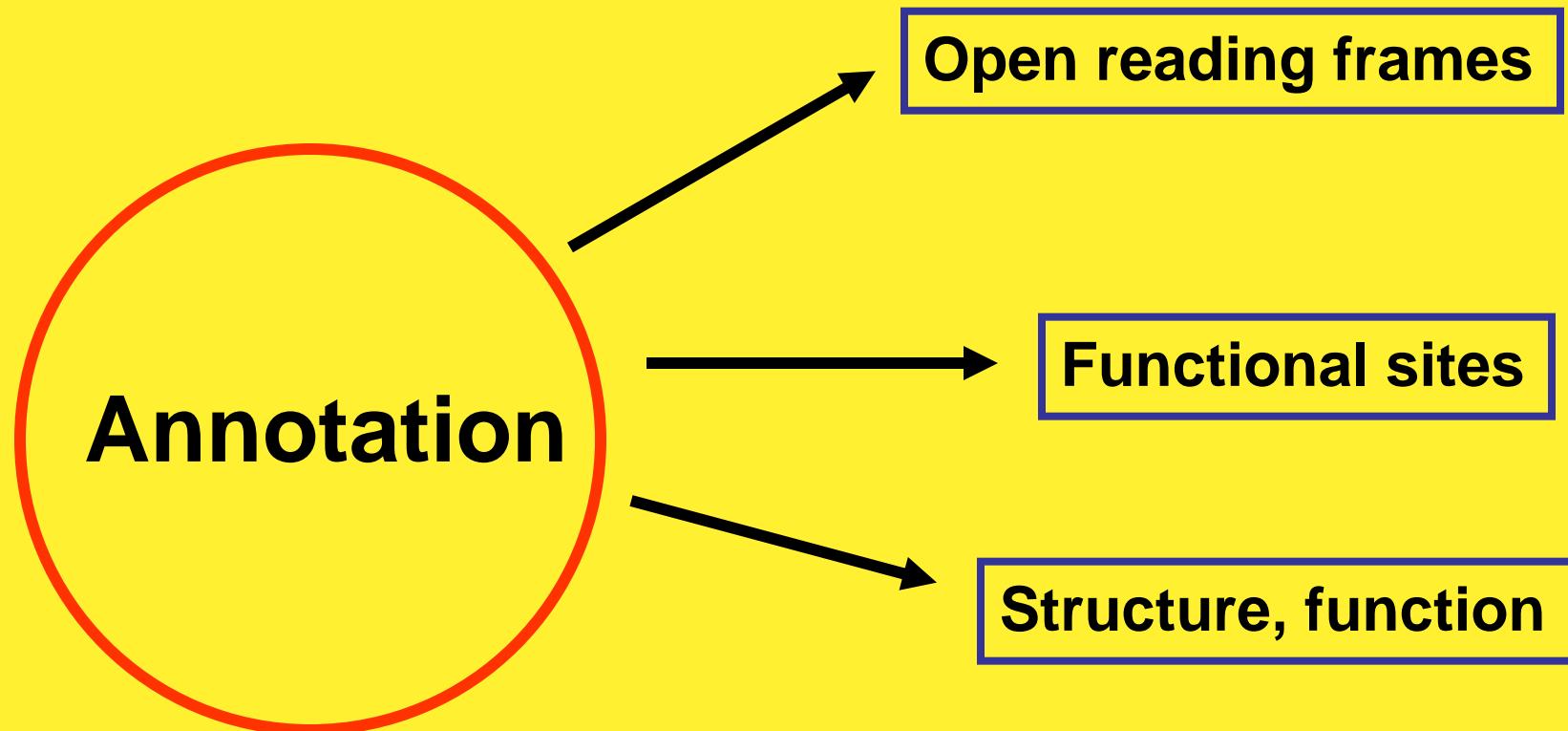
Eukaryotes 37

Prokaryotes 1708

Total 1745

What can we do with sequences and other type of molecular information?





CCTGACAAATTGACGTGCGGCATTGCATGCAGACGTGCATG
CGTGCAAATAATCAATGTGGACTTTCTGCGATTATGGAAGAA
CTTGTTACGCGTTTGTCAATGGCTTGGTCCCGCTTGTTC
AGAATGCTTTAATAAGCGGGGTTACCGGTTGGTAGCGAGA
AGAGCCAGTAAAAGACGCAGTGACGGAGATGTCTGATG CAA
TAT GGA CAA TTG GTT TCT TCT CTG AAT,
..... TGAAAAACGTA

promoter

TF binding site

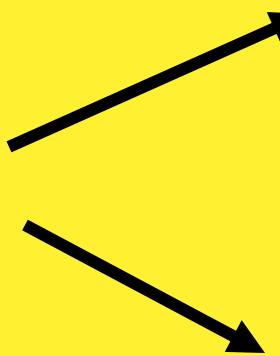
Transcription Start Site

CCTGACAAATTGACGTGC~~GG~~CATTGCATGC**AGACGTGCATG**
CGTGCAAA**TATCA**ATGTGGACTTTCTGC**GATTAT**GGAAAGAA
CTTGTTCACGCGTTTGTCAATGGCTTGGTCCCGCTTGTTC
AGAACGCTTTAATAAGCGGGGTTACCGGTTGGTAGCGAGA
AGAGCCAGTAAAAGACGCAGTGACGGAGATGTCTGATG CAA
TAT GGA CAA TTG GTT TCT TCT CTG AAT
..... TGAAAAACGTA

Ribosome binding Site

ORF = Open Reading Frame
CDS = Coding Sequence

Comparative genomics



Comparing ORFs

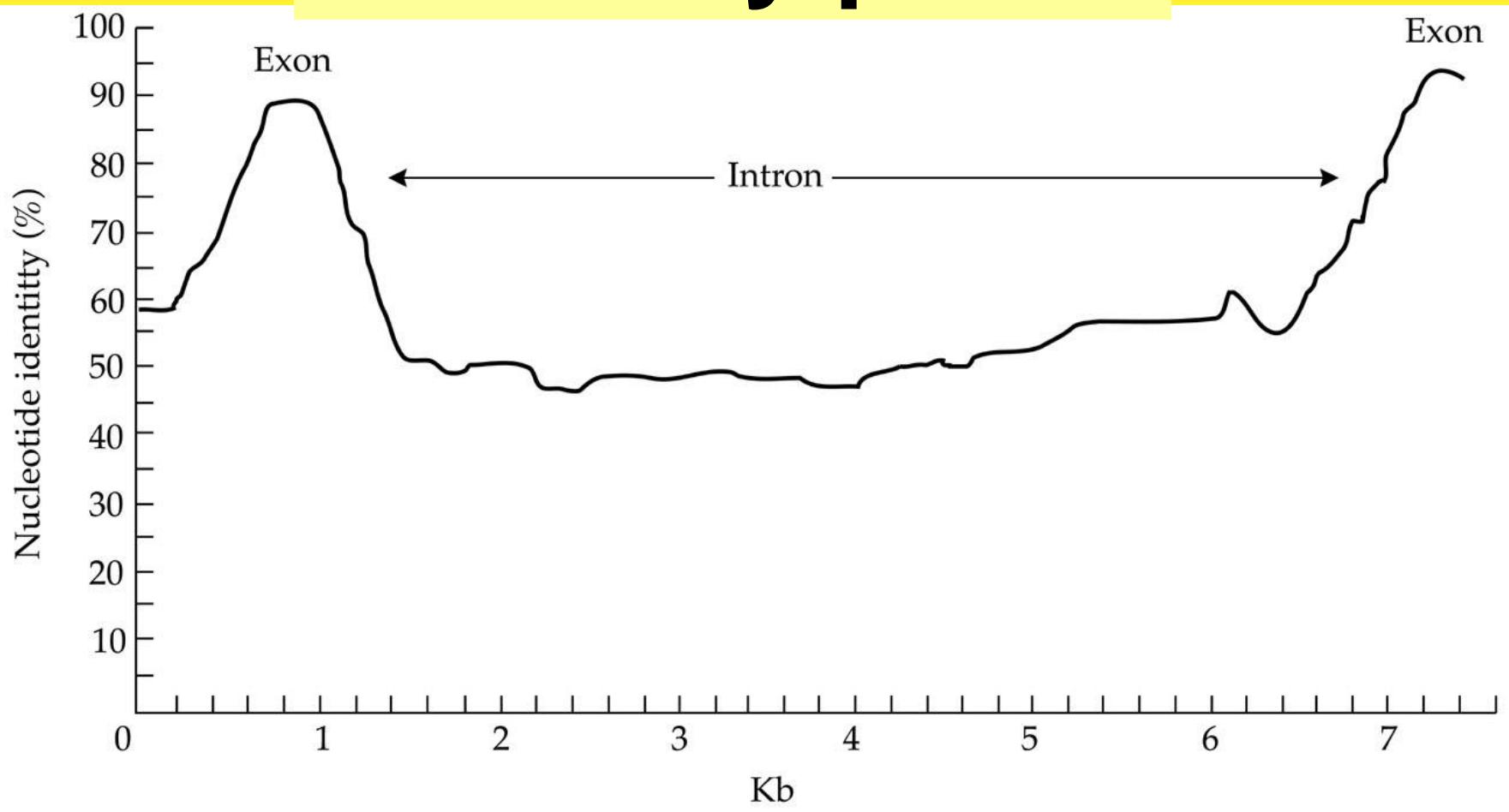
↓
Identifying orthologs

↓
**Inferences on structure
and function**

Comparing functional sites

↓
**Inferences on regulatory
networks**

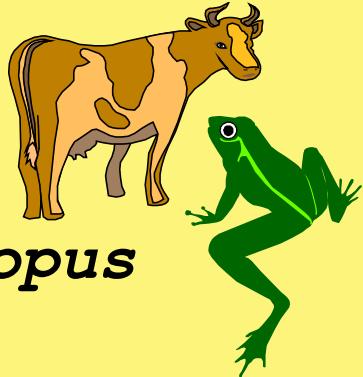
Similarity profiles



Researchers can learned a great deal about the structure and function of human genes by examining their counterparts in model organisms.

Alignment preproinsulin

Xenopus



MALWMQCLP-LVLVLLFSTPNTEALANQHL

MALWTRLRPLLALLALWPPPPARAFVNQHL

**** : * *.*: * : . * : . * : ****

Bos

Xenopus

Bos

CGSHLVEALYLVCGDRGFFYYPKIKRDIEQ

CGSHLVEALYLVCGERGFFYTPKARREVEG

***** : ***** * : * : : * : : :

Xenopus

Bos

AQVNGPQDNELDG-MQFQPQEYQKMKRGIV

PQVG---ALELAGGPGAGGLEGPPQKRGIV

. ** . * * * * ****

Xenopus

Bos

EQCCHSTCSLFQLENYCN

EQCCASVCSLYQLENYCN

**** * . *** : *****

Preproinsulin

Proinsulin

Signal

B chain

C peptide

A chain

B

S

S

A

Insulin

1.20×10^{-9}
substitutions
per site per year

0.20×10^{-9}
substitutions
per site per year

1.07×10^{-9}
substitutions
per site per year

Ultraconserved Elements in the Human Genome

Gill Bejerano, Michael Pheasant, Igor Makunin, Stuart Stephen, W. James Kent, John S. Mattick, & David Haussler
(Science 2004. 304:1321-1325)

There are 481 segments longer than 200 base pairs (bp) that are absolutely conserved (100% identity with no insertions or deletions) between orthologous regions of the human, rat, and mouse genomes. Nearly all of these segments are also conserved in the chicken and dog genomes, with an average of 95 and 99% identity, respectively. Many are also significantly conserved in fish. These ultraconserved elements of the human genome are most often located either overlapping exons in genes involved in RNA processing or in introns or nearby genes involved in the regulation of transcription and development.

There are 156 intergenic, untranscribed, ultraconserved segments

Megabase deletions of gene deserts result in viable mice

Marcelo A. Nóbrega*, Yiwen Zhu*, Ingrid Plajzer-Frick, Veena Afzal & Edward M. Rubin

**Junk:
Supporting evidence**

the Institute Walnut Creek, California 94598, USA, and
n Lawrence Berkeley National Laboratory Berkeley, California

contributed equally to this work

The importance of the roughly 98% of mammalian genomes not corresponding to protein coding sequences remains largely undetermined¹. Here we show that some large-scale deletions of the non-coding DNA referred to as gene deserts²⁻⁴ can be well tolerated by an organism. We deleted two large non-coding intervals, 1,511 kilobases and 845 kilobases in length, from the mouse genome. Viable mice homozygous for the deletions were generated and were indistinguishable from wild-type littermates with regard to morphology, reproductive fitness, growth, longevity and a variety of physiological parameters, including homeostasis. Further detailed analysis revealed expression differences in multiple genes bracketing the deleted segments in the two mouse strains. Together, the two deleted segments harbour 1,243 non-coding sequences conserved between humans and rodents (more than 100 base pairs, 70% identity). Some of the deleted sequences might encode for functions unidentified in our screen; nonetheless, these findings support the existence of potentially important regulatory elements in the non-coding genomes of mammals.

Junk is real!

Nature **431**, 988-993 (21 October 2004)

The genome of an organism is frequently referred to as its 'book

Functional genomics

Genome-wide profiling of:

- mRNA levels
- Protein levels

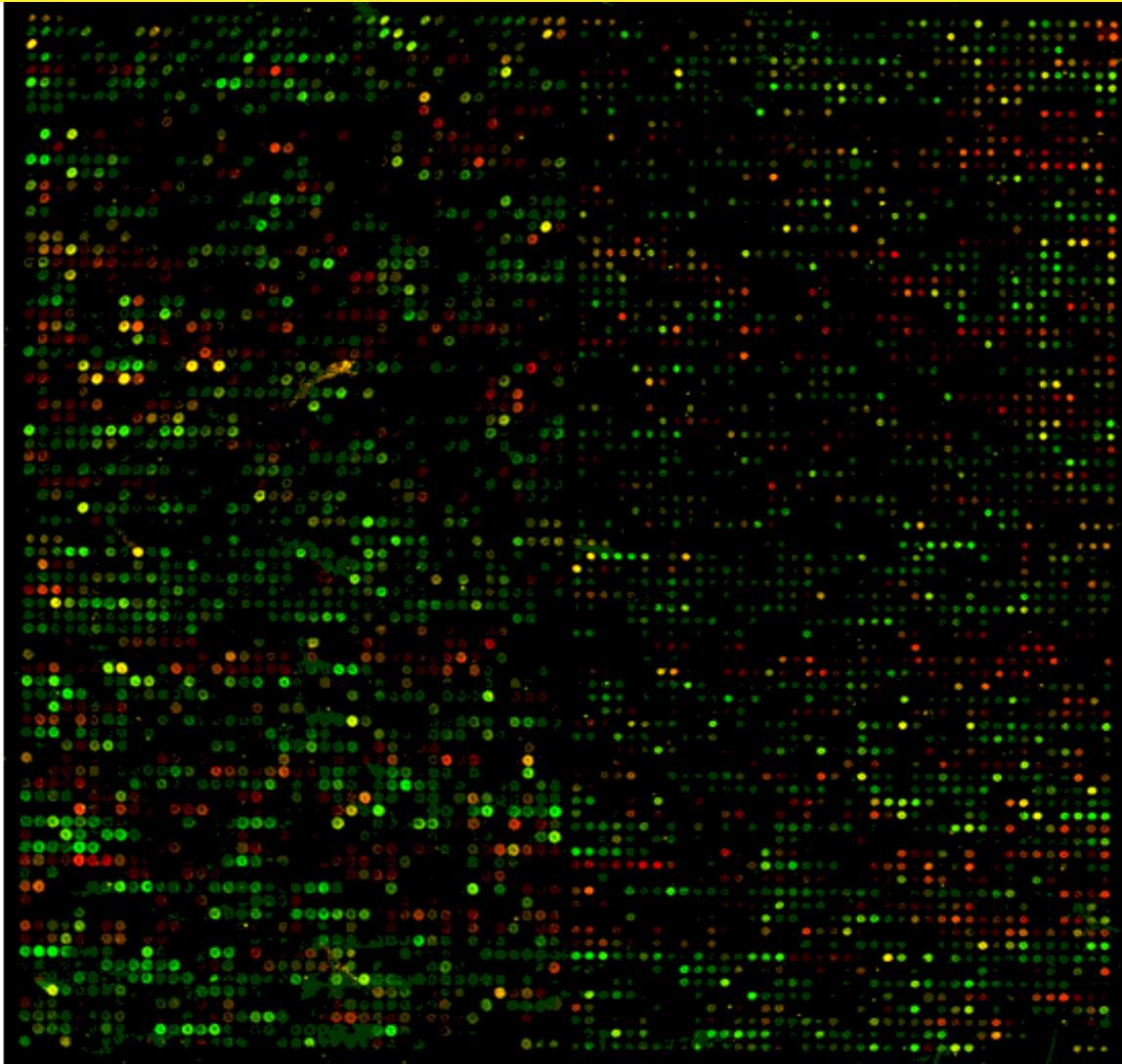
**↓
Co-expression of genes
and/or proteins**

**Identifying protein-protein
interactions**

**↓
Networks of interactions**

Understanding the function of genes and other parts of the genome

The complete *S. cerevisiae* genome on a microarray chip hybridised to RNA from cultures in anaerobic and aerobic stationary phase

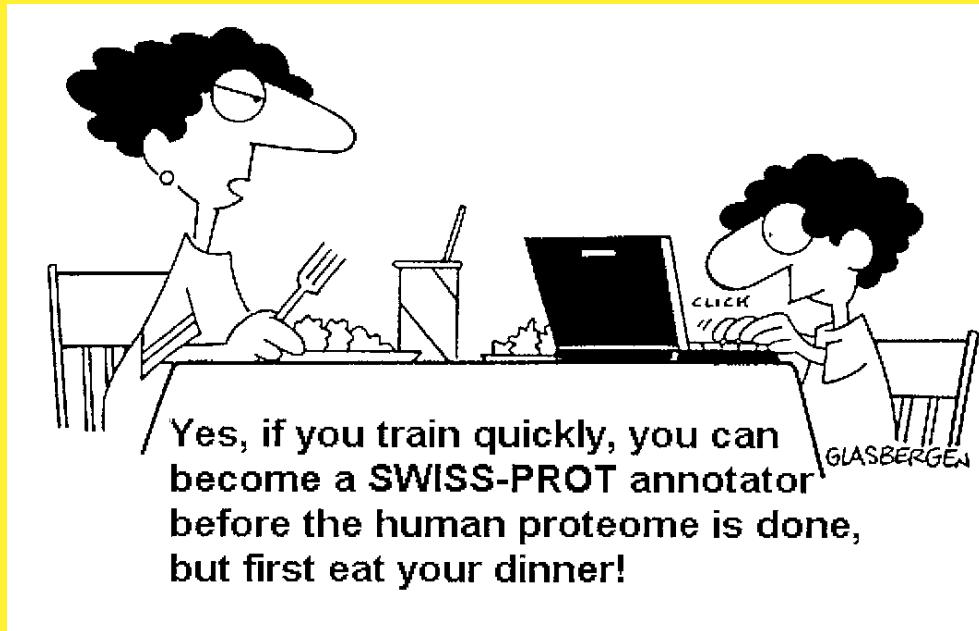


**Structural
genomics**



**Assign structure to all
proteins encoded in
a genome**





databases

Database or databank?

Initially

- Databank (in UK)
- Database (in the USA)

Solution

- The abbreviation *db*

What is a Database?

A **structured collection** of data held in computer storage; esp. one that incorporates software to make it accessible in a variety of ways; *transf.*, any **large collection** of information.

database management: the organization and manipulation of data in a database.

database management system (DBMS): a software package that provides all the functions required for database management.

database system: a database together with a database management system.

What is a database?

- A collection of data
 - structured
 - searchable (index)
-> table of contents
 - updated periodically (release)
-> new edition
 - cross-referenced ([hyperlinks](#))
-> links with other db
- Includes also associated tools (software) necessary for access, updating, information insertion, information deletion....
- Data storage management: flat files, relational databases...

Database: a « flat file » example

Flat-file database (« flat file, 3 entries »):

Accession number: 1

First Name: Amos

Last Name: Bairoch

Course: Pottery 2000; Pottery 2001;

//

Accession number: 2

First Name: Dan

Last name: Graur

Course: Pottery 2000, Pottery 2001; Ballet 2001, Ballet 2002

//

Accession number 3:

First Name: John

Last name: Travolta

Course: Ballet 2001; Ballet 2002;

//

- Easy to manage: all the entries are visible at the same time !

Database: a « relational » example

Relational database (« table file »):

Teacher	Accession number	Education
Amos	1	Biochemistry
Dan	2	Genetics
John	3	Scientology



Course	Year	Involved teachers
Advanced Pottery	2000; 2001	1; 2
Ballet for Fat People	2001; 2002	2; 3

Why biological databases?

- Exponential growth in biological data.
- Data (genomic sequences, 3D structures, 2D gel analysis, MS analysis, Microarrays....) are no longer published in a conventional manner, but directly submitted to databases.
- Essential tools for biological research. The only way to publish massive amounts of data without using all the paper in the world.

Distribution of sequences

- Books, articles 1968 -> 1985
- Computer tapes 1982 -> 1992
- Floppy disks 1984 -> 1990
- CD-ROM 1989 ->
- FTP 1989 ->
- On-line services 1982 -> 1994
- WWW 1993 ->
- DVD 2001 ->

Some statistics

- More than 1000 different ‘biological’ databases
- Variable size: <100Kb to >20Gb
 - DNA: > 20 Gb
 - Protein: 1 Gb
 - 3D structure: 5 Gb
 - Other: smaller
- Update frequency: **daily to annually to seldom to forget about it.**
- Usually accessible through the web (some free, some not)

■ Some databases in the field of molecular biology...

AATDB, AceDb, ACUTS, ADB, AFDB, AGIS, AMSdb,
ARR, AsDb, BBDB, BCGD, Beanref, Biolimage,
BioMagResBank, BIOMDB, BLOCKS, BovGBASE,
BOVMAP, BSORF, BTKbase, CANSITE, CarbBank,
CARBHYD, CATH, CAZY, CCDC, CD4OLbase, CGAP,
ChickGBASE, Colibri, COPE, CottonDB, CSNDB, CUTG,
CyanoBase, dbCFC, dbEST, dbSTS, DDBJ, DGP, DictyDb,
Picty_cDB, DIP, DOGS, DOMO, DPD, DPInteract, ECDC,
ECGC, ECO2DBASE, EcoCyc, EcoGene, EMBL, EMD db,
ENZYME, EPD, EpoDB, ESTHER, FlyBase, FlyView,
GCRDB, GDB, GENATLAS, Genbank, GeneCards,
Genline, GenLink, GENOTK, GenProtEC, GIFTS,
GPCRDB, GRAP, GRBase, gRNAsdb, GRR, GSDB,
HAEMB, HAMSTERS, HEART-2DPAGE, HEXAdb, HGMD,
HIDB, HIDC, HIVdb, HotMolecBase, HOVERGEN, HPDB,
HSC-2DPAGE, ICN, ICTVDB, IL2RGbase, IMGT, Kabat,
KDNA, KEGG, KloTho, LGIC, MAD, MaizeDb, MDB,
Medline, Mendel, MEROPS, MGDB, MGI, MHCPEP5
Micado, MitoDat, MITOMAP, MJDB, MmtDB, Mol-R-Us,
MPDB, MRR, MutBase, MycDB, NDB, NRSub, O-lycBase,
OMIA, OMIM, OPD, ORDB, OWL, PAHdb, PatBase, PDB,
PDD, Pfam, PhosphoBase, PigBASE, PIR, PKR, PMD,
PPDB, PRESAGE, PRINTS, ProDom, Prolysis, PROSITE,
PROTOMAP, RatMAP, RDP, REBASE, RGP, SBASE,
SCOP, SeqAnaiRef, SGD, SGP, SheepMap, Soybase,
SPAD, SRNA db, SRPDB, STACK, StyGene, Sub2D,
SubtiList, SWISS-2DPAGE, SWISS-3DIMAGE, SWISS-
MODEL Repository, SWISS-PROT, TelDB, TGN, tmRDB,
TOPS, TRANSFAC, TRR, UniGene, URNADB, V BASE,
VDRR, VectorDB, WDCM, WIT, WormPep, YEPD, YPD,
YPM, etc !!!!

Categories of databases for Life Sciences

- Sequences (DNA, protein)
- Genomics
- Mutation/polymorphism
- Protein domain/family
- Proteomics (2D gel, Mass Spectrometry)
- 3D structure
- Metabolic networks
- Regulatory networks
- **Bibliography**
- Expression (Microarrays,...)
- Specialized

NCBI:

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

EBI:

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

DDBJ:

<http://www.ddbj.nig.ac.jp/>

Literature Databases:

[Bookshelf](#): A collection of searchable biomedical books linked to PubMed.

[PubMed](#): Allows searching by author names, journal titles, and a new Preview/Index option. PubMed database provides access to over 12 million MEDLINE citations back to the mid-1960's. It includes History and Clipboard options which may enhance your search session.

[PubMed Central](#): The U.S. National Library of Medicine digital archive of life science journal literature.

[OMIM](#): Online Mendelian Inheritance in Man is a database of human genes and genetic disorders (also OMIA).

Search

All Databases

Search

Clear

NCBI Home**Site Map (A-Z)**

All Resources

Chemicals & Bioassays

Data & Software

DNA & RNA

Domains & Structures

Genes & Expression

Genetics & Medicine

Genomes & Maps

Homology

Literature

Proteins

Sequence Analysis

Taxonomy

Training & Tutorials

Variation

Welcome to NCBI

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

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

Get Started

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

Genomic Structural Variation

dbVar archives large scale genomic variation data and associates defined variants with phenotypic information.



Popular Resources

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

NCBI News

Preliminary genomic assemblies from two isolates from the European *E. coli* outbreak now available

07 Jun 2011

Preliminary genomic assemblies of two isolates are in the Nucleotide

New version of Cn3D (v.4.3) now available

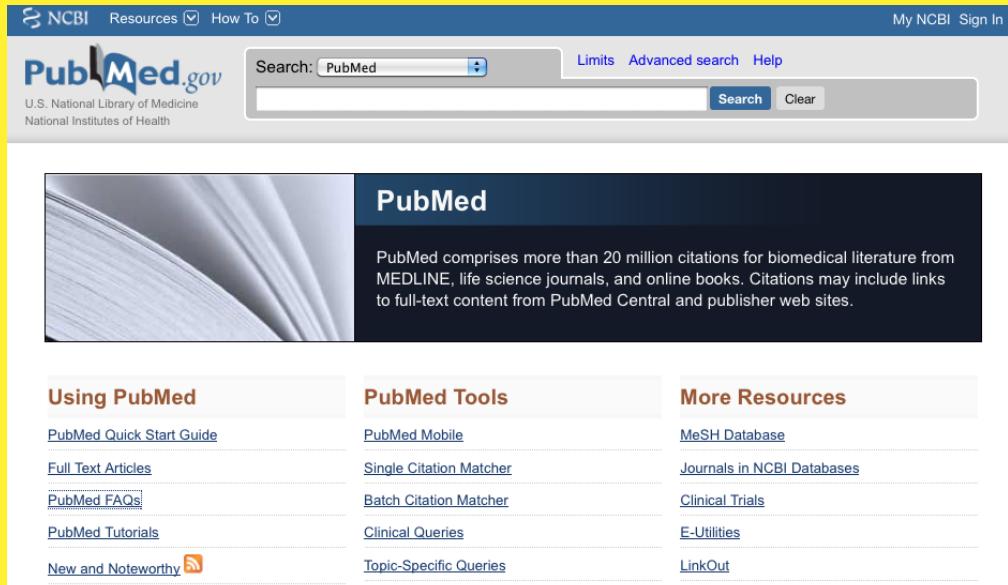
07 Jun 2011

A new version of this popular 3D molecular visualization program

More...

PubMed (Medline)

- MEDLINE covers the fields of medicine, nursing, dentistry, veterinary medicine, public health, and **preclinical sciences**
- Contains citations from approximately 5,200 worldwide journals in 37 languages; 60 languages for older journals.
- Contains over 20 million citations since 1948
- Contains links to biological db and to some journals
- New records are added to PreMEDLINE daily!



The screenshot shows the PubMed homepage. At the top, there's a blue header bar with the NCBI logo, 'Resources', 'How To', 'My NCBI', and 'Sign In' links. Below the header is the PubMed logo and text: 'U.S. National Library of Medicine' and 'National Institutes of Health'. A search bar has 'PubMed' in it, with 'Search' and 'Clear' buttons. Below the search bar, there's a large image of a stack of papers and a dark sidebar with the word 'PubMed' in white. The main content area describes PubMed's scope: 'PubMed comprises more than 20 million citations for biomedical literature from MEDLINE, life science journals, and online books. Citations may include links to full-text content from PubMed Central and publisher web sites.' At the bottom, there are three columns: 'Using PubMed' (with links to Quick Start Guide, Full Text Articles, PubMed FAQs, PubMed Tutorials, and New and Noteworthy), 'PubMed Tools' (with links to PubMed Mobile, Single Citation Matcher, Batch Citation Matcher, Clinical Queries, and Topic-Specific Queries), and 'More Resources' (with links to MeSH Database, Journals in NCBI Databases, Clinical Trials, E-Utilities, and LinkOut).

- Alerting services
 - <http://www.pubcrawler.ie/>
 - <http://www.biomail.org>

A search by subject:

“mitochondrion evolution”

NCBI Resources How To My NCBI Sign In

PubMed.gov
U.S. National Library of Medicine
National Institutes of Health

Search: PubMed RSS Save search Limits Advanced search Help

mitochondrion evolution Search Clear

Display Settings: Summary, 20 per page, Sorted by Recently Added Send to: Filter your results:

Results: 1 to 20 of 3990 << First < Prev Page 1 of 200 Next > Last >>

All (3990)
[Free Full Text \(1255\)](#)
[Review \(739\)](#)
[Manage Filters](#)

Energetics and genetics across the prokaryote-eukaryote divide.
1. Lane N.
Biol Direct. 2011 Jun 30;6(1):35. [Epub ahead of print]
PMID: 21714941 [PubMed - as supplied by publisher] [Free Article](#)
[Related citations](#)

A horizontally transferred tRNA(Cys) gene in the sugar beet mitochondrial genome: evidence that the gene is present in diverse angiosperms and its transcript is aminoacylated.
2. Kitazaki K, Kubo T, Kagami H, Matsumoto T, Fujita A, Matsuhira H, Matsunaga M, Mikami T.
Plant J. 2011 Jun 23. doi: 10.1111/j.1365-313X.2011.04684.x. [Epub ahead of print]
PMID: 21699590 [PubMed - as supplied by publisher]
[Related citations](#)

[Sequence and phylogeny analysis of the complete mitochondrial genome of Pelteobagrus vachelli.]
3. Li L, Liang HW, Li Z, Luo XZ, Hu GF, Zhang ZW, Zhu YY, Zou GW.
Yi Chuan. 2011 Jun;33(6):627-635. Chinese.
PMID: 21684869 [PubMed - as supplied by publisher]
[Related citations](#)

The complete mitochondrial genome of two recently derived species of the fish genus Nannoperca (Perciformes, Percichthyidae).
4. Prosdocimi F, de Carvalho DC, de Almeida RN, Beheregaray LB.
Mol Biol Rep. 2011 Jun 17. [Epub ahead of print]
PMID: 21681429 [PubMed - as supplied by publisher]
[Related citations](#)

Save Results in Collections Tutorial

My NCBI — Collections

Save: All None 0 items selected Merge Delete

Name	Type	Last Modified
My Bibliography	Standard	last month
Other Citations	Standard	never
pancreatic_cancer	PubMed	last month
ltd3_acetylation_reviews	PubMed	last month
interactive_video_across_in_children	PubMed	today

See larger video at YouTube
See all NCBI YouTube video channel videos

Titles with your search terms

Evolution and disease converge in the mitochondrion. [Biochim Biophys Acta. 2010]
On the evolution of programmed cell death: apoptosis of the unicellular [Cell Death Differ. 2002]
Fungal evolution: the case of the vanishing mitochondrion. [Curr Opin Microbiol. 2005]

See more...



Search: PubMed

Limits Advanced search Help

Search Clear

Display Settings: AbstractSend to: [Biol Direct. 2011 Jun 30;6\(1\):35. \[Epub ahead of print\]](#)

Energetics and genetics across the prokaryote-eukaryote divide.

[Lane N.](#)

Abstract

ABSTRACT: BACKGROUND: All complex life on Earth is eukaryotic. All eukaryotic cells share a common ancestor that arose just once in 4 billion years of **evolution**. Prokaryotes show no tendency to evolve greater morphological complexity, despite their metabolic virtuosity. Here I argue that the eukaryotic cell originated in a unique prokaryotic endosymbiosis, a singular event that transformed the selection pressures acting on both host and endosymbiont. Results: The reductive **evolution** and specialisation of endosymbionts to **mitochondria** resulted in an extreme genomic asymmetry, in which the residual mitochondrial genomes enabled the expansion of bioenergetic membranes over several orders of magnitude, overcoming the energetic constraints on prokaryotic genome size, and permitting the host cell genome to expand (in principle) over 200,000-fold. This energetic transformation was permissive, not prescriptive; I suggest that the actual increase in early eukaryotic genome size was driven by a heavy early bombardment of genes and introns from the endosymbiont to the host cell, producing a high mutation rate. Unlike prokaryotes, with lower mutation rates and heavy selection pressure to lose genes, early eukaryotes without genome-size limitations could mask mutations by cell fusion and genome duplication, as in allopolyploidy, giving rise to a proto-sexual cell cycle. The side effect was that a large number of shared eukaryotic basal traits accumulated in the same population, a sexual eukaryotic common ancestor, radically different to any known prokaryote. Conclusions: The combination of massive bioenergetic expansion, release from genome-size constraints, and high mutation rate favoured a protosexual cell cycle and the accumulation of eukaryotic traits. These factors explain the unique origin of eukaryotes, the absence of true evolutionary intermediates, and the **evolution** of sex in eukaryotes but not prokaryotes. Reviewers: This article was reviewed by: Eugene Koonin, William Martin, Ford Doolittle and Mark van der Giezen. For complete reports see the Reviewers' Comments section.

PMID: 21714941 [PubMed - as supplied by publisher] [Free full text](#) [LinkOut - more resources](#)

Related citations

The origin of introns and their role in eukaryogenesis: a compromise [Biol Direct. 2006]

The energetics of genome complexity. [Nature. 2010]

On the need for widespread horizontal gene transfers under genome size cc [Biol Direct. 2009]

Review The neomuran origin of archaebacteria, the negibacterial i [Int J Syst Evol Microbiol. 2002]

Review The phagotrophic origin of eukaryotes and phylogenetic [Int J Syst Evol Microbiol. 2002]

[See reviews...](#)[See all...](#)

Recent activity

[Turn Off](#) [Clear](#)

Energetics and genetics across the prokaryote-eukaryote divide. [PubMed](#)

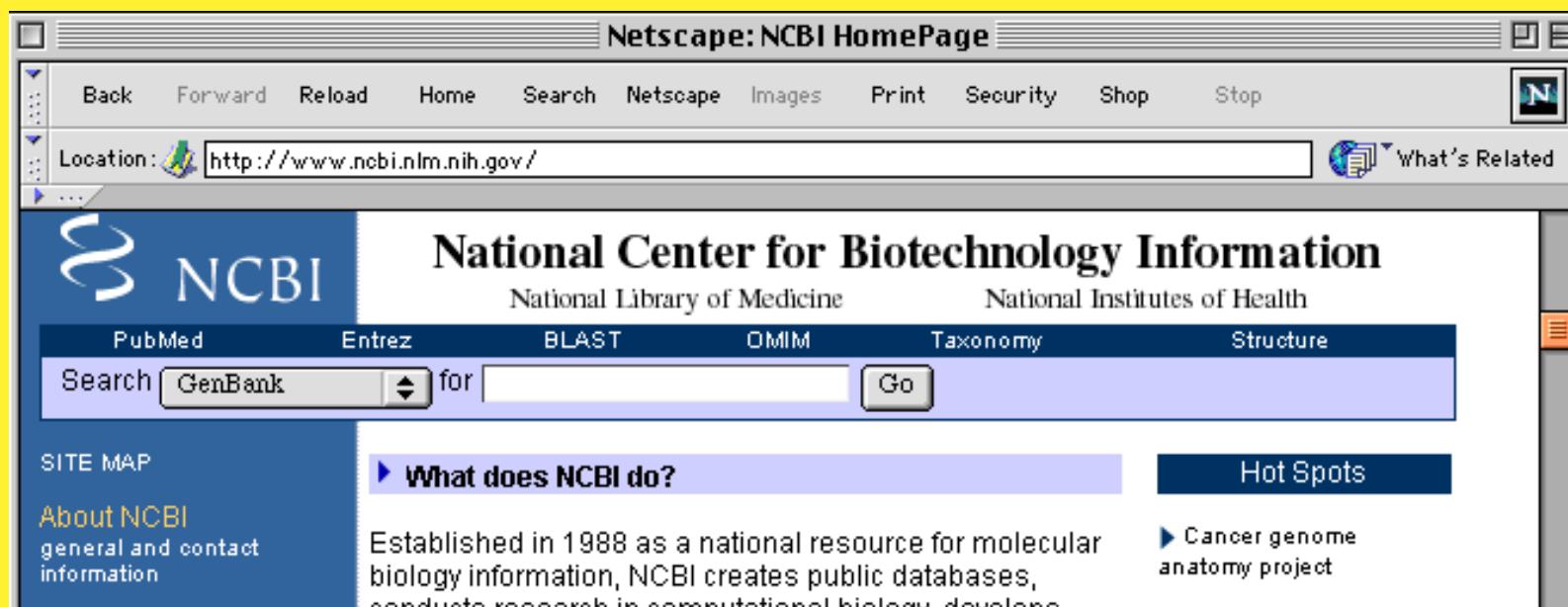
mitochondrion evolution (3990) [PubMed](#)

PubMed Help - PubMed Help [Bookshelf](#)

[See more...](#)

Type in a Query term

- Enter your search words in the query box and hit the “Go” button



<http://www.ncbi.nlm.nih.gov/entrez/query/static/help/helpdoc.html#Searchfaq>

The Syntax ...

1. **Boolean operators:** AND, OR, NOT must be entered in UPPERCASE (e.g., promoters OR response elements). The default is AND.
2. Entrez processes all Boolean operators in a left-to-right sequence. The order in which Entrez processes a search statement can be changed by enclosing individual concepts in parentheses. The terms inside the parentheses are processed first. For example, the search statement: g1p3 OR (response AND element AND promoter).
3. Quotation marks: The term inside the quotation marks is read as one phrase (e.g. “public health” is different than public health, which will also include articles on public latrines and their effect on health workers).
4. Asterisk: Extends the search to all terms that start with the letters before the asterisk. For example, dia* will include such terms as ⁶⁹ diaphragm, dial, and diameter.

Refine the Query

- Often a search finds too many (or too few) sequences, so you can go back and try again with more (or fewer) keywords in your query
- The “History” feature allows you to combine any of your past queries.
- The “Limits” feature allows you to limit a query to specific organisms, sequences submitted during a specific period of time, etc.
- [Many other features are designed to search for literature in MEDLINE]

Related Items

You can search for a text term in sequence annotations or in **MEDLINE** abstracts, and find all articles, DNA, and protein sequences that mention that term.

Then from any article or sequence, you can move to "related articles" or "related sequences".

- Relationships between sequences are computed with [BLAST](#)
- Relationships between articles are computed with "MESH" terms (shared keywords)
- Relationships between DNA and protein sequences rely on accession numbers
- Relationships between sequences and **MEDLINE** articles rely on both shared keywords and the mention of accession numbers in the articles.

Search PubMed

for human HIV 1

Go Clear

Limits

Preview/Index

History

Clipboard

Details

- Use All Fields pull-down menu to specify a field.
- Boolean operators AND, OR, NOT must be in upper case.
- If search fields tags are used enclose in square brackets, e.g., rubella [ti].
- Search limits may exclude in process and publisher supplied citations.

Limited to:

All Fields

All Fields

Author

Corporate Author

EC/RN Number

Entrez Date

Filter

First Author

Full Author Name

Grant Number

Issue

Journal

 only items with abstracts

Languages

Humans or Animals

To
month and day are optional.

Subsets

Subsets

AIDS

Bioethics

Cancer

Complementary Medicine

Core clinical journals

Dental journals

History of Medicine

MEDLINE

Nursing journals

OLDMEDLINE for Pre1966

Search PubMed

PubMed

for human HIV 1

Preview

Go

Clear

[Limits](#)[Preview/Index](#)[History](#)[Clipboard](#)[Details](#)

- Enter terms and click Preview to see only the number of search results.
- To combine searches use # before search number, e.g., (#2 OR #3) AND asthma.
- Click on query # to add to strategy

[Search](#)[Most Recent Queries](#)[#2 Search human HIV 1](#)

04

Add Term(s) to Query or View Index:

- Enter a term in the text box; use the pull-down menu to specify a search field.
- Click Preview to add terms to the query box and see the number of search results, or click Index to view terms within

All Fields

Affiliation

All Fields

Author

Corporate Author

EC/RN Number

Entrez Date

Filter

First Author

Full Author Name

Grant Number

Issue

add a term to the query box.

Preview

Index

Search Field Descriptions and Tags

Affiliation [AD]	Issue [IP]	Place of Publication [PL]
All Fields [ALL]	Journal Title [TA]	Publication Date [DP]
Author [AU]	Language [LA]	Publication Type [PT]
Comment Corrections	Last Author [LASTAU]	Publisher Identifier [AID]
Corporate Author [CN]	MeSH Date [MHDA]	Secondary Source ID [SI]
EC/RN Number [RN]	MeSH Major Topic [MAJR]	Subset [SB]
Entrez Date [EDAT]	MeSH Subheadings [SH]	Substance Name [NM]
Filter [FILTER]	MeSH Terms [MH]	Text Words [TW]
First Author Name [1AU]	NLM Unique ID [JID]	Title [TI]
Full Author Name [FAU]	Other Term [OT]	Title/Abstract [TIAB]
Full Investigator Name [FIR]	Owner	Transliterated Title [TT]
Grant Number [GR]	Pagination [PG]	UID [PMID]
Investigator [IR]	Personal Name as Subject [PS]	Volume [VI]
	Pharmacological Action MeSH Terms [PA]	

A search by authors:

“Esser” [au] AND “martin” [au]

NCBI Resources How To My NCBI Sign In

PubMed.gov U.S. National Library of Medicine National Institutes of Health

Search: PubMed RSS Save search Limits Advanced search Help

"Esser C" [au] AND "Martin W" [au] Search Clear

Display Settings: Summary, Sorted by Recently Added Send to:

Results: 4

[Supertrees and symbiosis in eukaryote genome evolution.](#)

1. **Esser C, Martin W.**
Trends Microbiol. 2007 Oct;15(10):435-7. Epub 2007 Sep 19.
PMID: 17884500 [PubMed - indexed for MEDLINE]
[Related citations](#)

[The origin of mitochondria in light of a fluid prokaryotic chromosome model.](#)

2. **Esser C, Martin W**, Dagan T.
Biol Lett. 2007 Apr 22;3(2):180-4.
PMID: 17251118 [PubMed - indexed for MEDLINE] [Free PMC Article](#)
[Free full text](#) [Related citations](#)

[A genome phylogeny for mitochondria among alpha-proteobacteria and a predominantly eubacterial ancestry of yeast nuclear genes.](#)

3. **Esser C, Ahmadinejad N, Wiegand C, Rotte C, Sebastiani F, Gellius-Dietrich G, Henze K, Kretschmann E, Richly E, Leister D, Bryant D, Steel MA, Lockhart PJ, Penny D, Martin W.**
Mol Biol Evol. 2004 Sep;21(9):1643-60. Epub 2004 May 21.
PMID: 15155797 [PubMed - indexed for MEDLINE] [Free Article](#)
[Related citations](#)

[Phylogenomics of the reproductive parasite Wolbachia pipiensis wMel: a streamlined genome overrun by mobile genetic elements.](#)

4. Wu M, Sun LV, Vamathevan J, Riegler M, Deboy R, Brownlie JC, McGraw EA, **Martin W, Esser C, Ahmadinejad N, Wiegand C, Madupu R, Beanan MJ, Brinkac LM, Daugherty SC, Durkin AS, Kolonay JF, Nelson WC, Mohamoud Y, Lee P, Berry K, Young MB, Utterback T, Weidman J, Nierman WC, Paulsen IT, Nelson KE, Tettelin H, O'Neill SL, Eisen JA.**
PLoS Biol. 2004 Mar;2(3):E69. Epub 2004 Mar 16.
PMID: 15024419 [PubMed - indexed for MEDLINE] [Free PMC Article](#)
[Free full text](#) [Related citations](#)

Filter your results:
All (4)
[Free Full Text \(3\)](#)
Review (0)
[Manage Filters](#)

Save Results in Collections Tutorial

My NCBI — Collections

Select All None 0 items selected Merge Delete

Name	Type	Last Modified	Type
My Bibliography	Private	last month	Standard
Other Citations	Private	never	Standard
pancreatic cancer	Public	last month	PubMed
p62 acetylation, reviews	Public	last month	PubMed
obstructive sleep apnea in children	Public	today	PubMed

See larger video at YouTube
See all NCBI YouTube video channel videos

2 free full-text articles in PubMed Central

The origin of mitochondria in light of a fluid prokaryotic chromosome model. [Biol Lett. 2007]
Phylogenomics of the reproductive parasite Wolbachia pipiensis wMel: a str [PLoS Biol. 2004]
See all (2)...

Display Settings: Summary, Sorted by Recently Added Send to:

A search by title word:

“Wolbachia pipiensis” [ti]

Entrez PubMed

<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?CMD=search&DB=pubmed>

Fukille Gmail uni-duss mail Wörterbuch LEO מורפיקה Haaretz Ynet NCBI UCSC ISI eJournal MBE JME Nature NAR Gen. Res

My NCBI [Sign In] [Register]

All Databases PubMed Nucleotide Protein Genome Structure OMIM PMC Journals Books

Search PubMed for "Wolbachia pipiensis" [ti] Go Clear Save Search

Limits Preview/Index History Clipboard Details

Display Summary Show 20 Sort by Send to

All: 29 Review: 3

Items 1 - 20 of 29 Page 1 of 2 Next

1: [Makepeace BL, Rodgers L, Trees AJ.](#) Related Articles, Links
Rate of elimination of Wolbachia pipiensis by doxycycline in vitro increases following drug withdrawal.
Antimicrob Agents Chemother. 2006 Mar;50(3):922-7.
PMID: 16495252 [PubMed - indexed for MEDLINE]

2: [Casiraghi M, Bordenstein SR, Baldo L, Lo N, Beninati T, Werren JH, Bandi C.](#) Related Articles, Links
Phylogeny of Wolbachia pipiensis based on gltA, groEL and ftsZ gene sequences: clustering of arthropod and nematode symbionts in the F supergroup, and evidence for further diversity in the Wolbachia tree.
Microbiology. 2005 Dec;151(Pt 12):4015-22.
PMID: 16339946 [PubMed - indexed for MEDLINE]

3: [Mavingui P, Van VT, Labeyrie E, Rances E, Vavre F, Simonet P.](#) Related Articles, Links
Efficient procedure for purification of obligate intracellular Wolbachia pipiensis and representative amplification of its genome by multiple-displacement amplification.
Appl Environ Microbiol. 2005 Nov;71(11):6910-7.
PMID: 16269725 [PubMed - indexed for MEDLINE]

4: [Iturbe-Ormaetxe I, Burke GR, Riegler M, O'Neill SL.](#) Related Articles, Links
Distribution, expression, and motif variability of ankyrin domain genes in Wolbachia pipiensis.
J Bacteriol. 2005 Aug;187(15):5136-45.
PMID: 16030207 [PubMed - indexed for MEDLINE]

NCBI PubMed Services Related Resources Order Documents NLM Mobile NLM Catalog NLM Gateway TOXNET Consumer Health Clinical Alerts ClinicalTrials.gov PubMed Central

Database Search Strategies

- General search principles - not limited to sequence (or to biology).
- Start with broad keywords and narrow the search using more specific terms.
- Try variants of spelling, numbers, etc.
- Search many databases.
- **Be persistent!!**

Searching PubMed

- Structureless searches
 - Automatic term mapping
- Structured searches
 - Tags, e.g. [au], [ta], [dp], [ti]
 - Boolean operators, e.g. AND, OR, NOT, ()
- Additional features
 - Subsets, limits
 - Clipboard, history

Start working:

Search PubMed

- 1. cuban cigars**
- 2. cuban OR cigars**
- 3. “cuban cigars”**
- 4. cuba* cigar***
- 5. (cuba* cigar*) NOT smok***
- 6. Fidel Castro**
- 7. “fidel castro”**
- 8. #6 NOT #7**

“Details” and “History” in PubMed

The screenshot shows the PubMed search interface. In the search bar, the query 'for castro j' is entered. Below the search bar, there are tabs for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. The 'Details' tab is currently selected. Underneath, there are buttons for 'Display' (set to 'Summary') and 'Show' (set to 20). There are also 'Sort By' and 'Send to' dropdown menus. The search results are displayed as a list of items. Item 1 is a study titled 'Elimination of Onchocercia volvulus Transmission in the Santa Rosa Focus of Guatemala.' by Lindblade KA et al. Item 2 is a study titled 'Diabetes care in hospitalized noncritically ill patients: More evidence for clinical inertia and negative therapeutic momentum.' by Cook CB et al.

A service of the National Library of Medicine
and the National Institutes of Health
www.pubmed.gov

All Databases PubMed Nucleotide Protein Genome Structure OMIM PMC Journals

Search PubMed for castro j Go Clear Save Search

About Entrez Text Version Entrez PubMed Overview Help | FAQ Tutorials New/Noteworthy E-Utilities PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation

Limits Preview/Index History Clipboard Details

Display Summary Show 20 Sort By Send to

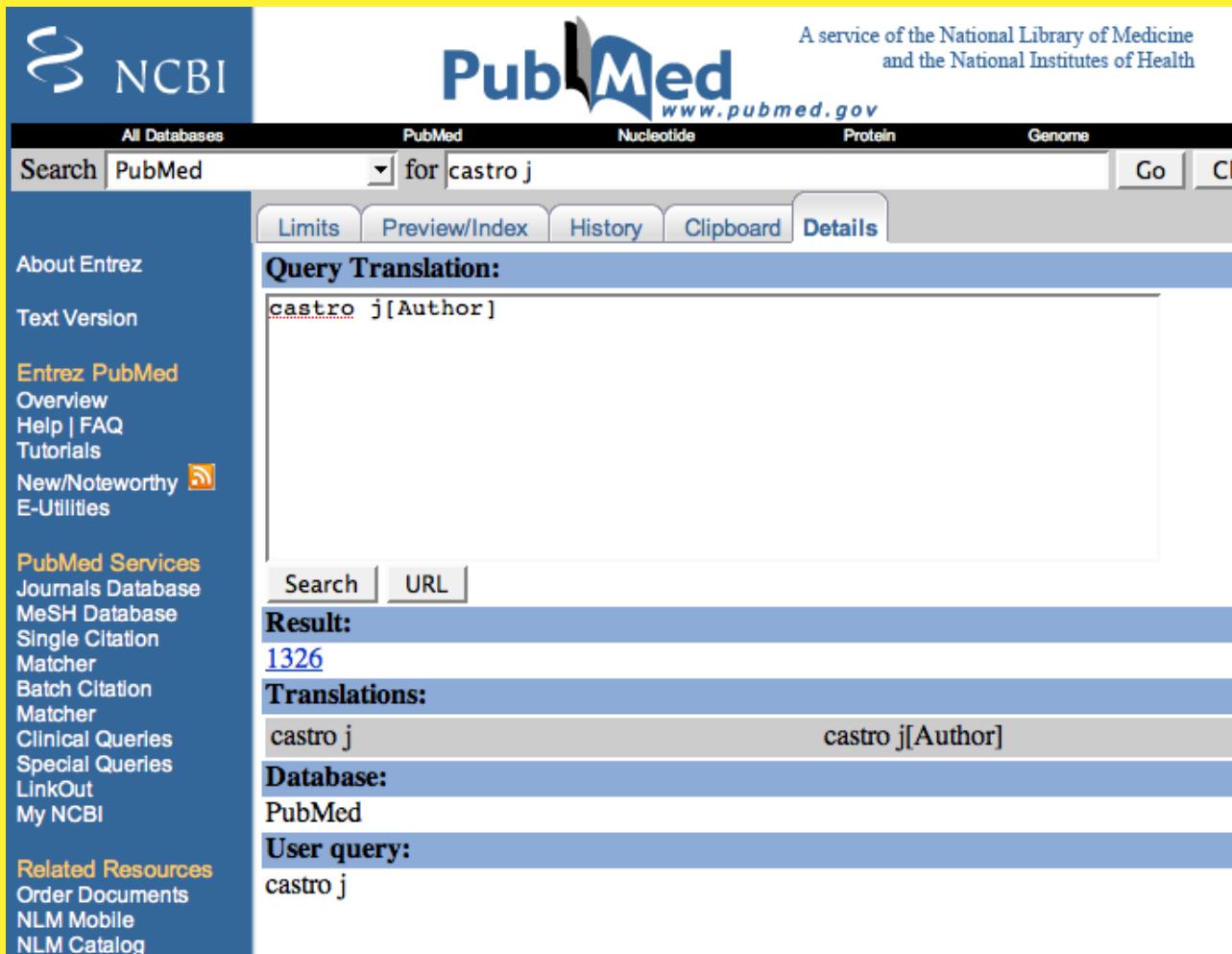
All: 1326 Review: 57

Items 1 - 20 of 1326

1: [Lindblade KA, Arana B, Zea-Flores G, Rizzo N, Porter CH, Dominguez A, Cruz-Ortiz N, Unnasch TR, Punkosdy GA, Richards J, Sauerbrey M, Castro J, Catu E, Oliva O, Richards FO Jr.](#)
Elimination of Onchocercia volvulus Transmission in the Santa Rosa Focus of Guatemala.
Am J Trop Med Hyg. 2007 Aug;77(2):334-341.
PMID: 17690408 [PubMed - as supplied by publisher]

2: [Cook CB, Castro JC, Schmidt RE, Gauthier SM, Whitaker MD, Roust LR, Argueta R, Hull BP, Zimmerman RS.](#)
Diabetes care in hospitalized noncritically ill patients: More evidence for clinical inertia and negative therapeutic momentum.
J Hosp Med. 2007 Aug;7(2):203-211 [Epub ahead of print]
PMID: 17683100 [PubMed - as supplied by publisher]

“Details” and “History” in PubMed



The screenshot shows the PubMed search interface. The search bar contains the query "castro j". Below the search bar, a navigation bar includes tabs for Limits, Preview/Index, History, Clipboard, and Details. The "Details" tab is currently selected. The main results area displays the following information:

- Query Translation:** castro j[Author]
- Result:** 1326
- Translations:** castro j castro j[Author]
- Database:** PubMed
- User query:** castro j

The OMIM (Online Mendelian Inheritance in Man)

- Genes and genetic disorders
- Edited by team at Johns Hopkins
- Updated daily

MIM Number Prefixes

- * gene with known sequence
- + gene with known sequence and phenotype
- # phenotype description, molecular basis known
- % mendelian phenotype or locus, molecular basis unknown
- no prefix other, mainly phenotypes with suspected mendelian basis

Searching OMIM

- Search Fields
 - Name of trait, e.g., hypertension
 - Cytogenetic location, e.g., 1p31.6
 - Inheritance, e.g., autosomal dominant
 - Gene, e.g., coagulation factor VIII

OMIM search tags

All Fields	[ALL]
Allelic Variant	[AV] or [VAR]
Chromosome	[CH] or [CHR]
Clinical Synopsis	[CS] or [CLIN]
Gene Map	[GM] or [MAP]
Gene Name	[GN] or [GENE]
Reference	[RE] or [REF]



All Databases

PubMed

Nucleotide

Protein

Genome

Structure

Search OMIM

for f8c

Go

Clear

[Limits](#)[Preview/Index](#)[History](#)[Clipboard](#)[Details](#)

Entrez

OMIM

Search OMIM

Search Gene Map

Search Morbid Map

Help

OMIM Help

How to Link

FAQ

Numbering System

Symbols

How to Print

Citing OMIM

Download

OMIM Facts

Statistics

Update Log

Restrictions on Use

Allied Resources

Genetic Alliance

Databases

HGMD

Locus-Specific

Model Organisms

MitoMap

Phenotype

- To Search all fields, leave the following boxes unchecked.
- To narrow the search, check the boxes with specific fields' names, or use [search field tags](#) enclosed in square brackets, e.g. aaa[title].
- [Boolean operators](#) AND, OR, NOT must be in upper case.

Search in Field(s):[clear](#)

<input type="checkbox"/> Title	<input type="checkbox"/> MIM number	<input type="checkbox"/> Allelic Variants
<input type="checkbox"/> Text	<input type="checkbox"/> References	<input type="checkbox"/> Clinical Synopsis
<input type="checkbox"/> Gene Map Disorder		<input type="checkbox"/> Contributors

Chromosome(s):[clear](#)

<input type="checkbox"/> 1	<input type="checkbox"/> 2	<input type="checkbox"/> 3	<input type="checkbox"/> 4	<input type="checkbox"/> 5	<input type="checkbox"/> 6	<input type="checkbox"/> 7	<input type="checkbox"/> 8
<input type="checkbox"/> 9	<input type="checkbox"/> 10	<input type="checkbox"/> 11	<input type="checkbox"/> 12	<input type="checkbox"/> 13	<input type="checkbox"/> 14	<input type="checkbox"/> 15	<input type="checkbox"/> 16
<input type="checkbox"/> 17	<input type="checkbox"/> 18	<input type="checkbox"/> 19	<input type="checkbox"/> 20	<input type="checkbox"/> 21	<input type="checkbox"/> 22	<input type="checkbox"/> X	<input type="checkbox"/> Y
<input type="checkbox"/> mitochondrial <input type="checkbox"/> unknown							

MIM Number Prefix:[clear](#)

<input type="checkbox"/> * gene with known sequence
<input type="checkbox"/> + gene with known sequence and phenotype
<input type="checkbox"/> # phenotype description, molecular basis known
<input type="checkbox"/> % mendelian phenotype or locus, molecular basis unknown
<input type="checkbox"/> base other, mainly phenotypes with suspected mendelian basis

Only Records with:[clear](#)

<input type="checkbox"/> Allelic Variants
<input type="checkbox"/> Clinical Synopsis
<input type="checkbox"/> Gene map locus

Creation Date

[From](#)[To](#)

Last Modification

[From](#)[To](#)

Use the format YYYY/MM/DD; month and day are optional.

Start working:

Search OMIM

How many types of hemophilia are there?

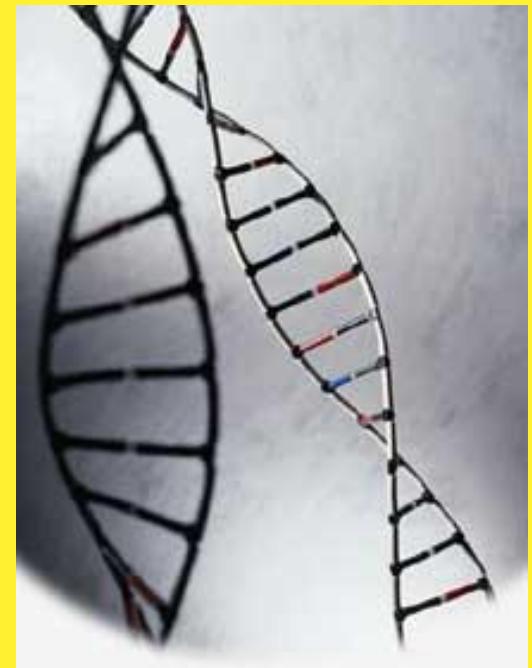
For how many is the affected gene known?

What are the genes involved in hemophilia A?

What are the mutations in hemophilia A?

Online Literature databases

1. How to use the UH online Library?
2. Online glossaries
3. Google Scholar
4. Google Books
5. Web of Science



How to use the online UH Library? <http://info.lib.uh.edu/>

UNIVERSITY of **HOUSTON** LIBRARIES

M.D. Anderson Library
Today's Hours 7:00am - 12:45am

Search Services About Help

OneSearch
e-Journals
Databases Catalog
Research Guides
Site Search

Databases

Search for database by name or keyword.
Top Databases: [Academic Search Complete](#), [JSTOR](#), [PUBMED](#)

12 A B C D E F G H I J K L M N O P R S T U V W
[New Databases](#) | [Browse by Subject](#) | [Browse by Type](#)



Find Computer Availability

Services »

- [Interlibrary Loan](#)
- [View & Renew](#)
- [Computers & Printing](#)
- [For Faculty & Graduate Students](#)
- [Course Reserves](#)
- [More...](#)

About »

- [Hours](#)
- [Campus Libraries & Collections](#)
- [Maps & Directions](#)
- [Staff Directory](#)
- [Employment](#)
- [News & Events](#)
- [More...](#)

Help »

- [Contact Us](#)
- [Ask a Librarian](#)
- [Live Chat](#)
- [Research Guides](#)
- [Electronic Resources Issues](#)
- [More...](#)

© 2010 The University of Houston, 114 University Libraries, Houston, TX 77204-2000  713-743-1050 
[Website Feedback](#) | [Policies](#) | [Giving to the Libraries](#) | [Contact Us](#) | [Maps & Directions](#) | [Site Map](#) | [Mobile Site](#)
[UH Home](#) | [UH System](#) | [State of Texas](#) | [Emergency Site](#)

UH Library User Authentication

Access to this electronic resource is limited to current University of Houston main campus students, faculty and person to ensure that he or she uses these resources only for individual, noncommercial use without systematically download portions of the information provided. Furthermore, copying or tampering with the software that organizes and displays info agreements.

Misuse of these resources can result in the termination of the licensing agreement and the loss of use of the resource by the campus. Penalties for misuse can and will be imposed by the University of Houston. Please read the policy for Appropriate details: <http://www.uh.edu/infotech/refguide/appropriate.html>

Providers and the University of Houston Libraries offer these resources on an "as is basis." There are no warranties for completeness, availability of storage and delivery devices. The University of Houston Libraries do not assume and hereby disclaim any liability resulting from the use of information in licensed resources.

By entering your login information, you agree to the above terms of use for the electronic resources provided by the UH Libraries.

Please login by entering your:

Last Name:

Your Library ID/Barcode:

[What is my Library ID/Barcode?](#)

Online Glossaries

Bioinformatics :

<http://www.geocities.com/bioinformaticsweb/glossary.html>

<http://big.mcw.edu/>

Genomics:

<http://www.geocities.com/bioinformaticsweb/genomicglossary.html>

Molecular Evolution:

<http://workshop.molecularevolution.org/resources/glossary/>

Biology dictionary:

http://www.biology-online.org/dictionary/satellite_cells

Other glossaries, e.g., the list of phobias:

<http://www.phobialist.com/class.html>

4. Google Scholar

<http://www.scholar.google.com/>

The screenshot shows the Google Scholar homepage in a web browser. The title bar reads "Google Scholar" and the address bar shows "scholar.google.com". The page features the classic Google logo with "scholar" written below it. A search bar is centered, with "Search" and "Advanced Scholar Search" buttons. Below the search bar are two radio button options: one selected for "Articles" (with a checked checkbox) and one for "Legal opinions and journals". A green banner below the search bar says "Stand on the shoulders of giants". At the bottom, links include "Go to Google Home - About Google - About Google Scholar" and the copyright notice "©2011 Google". The browser's toolbar at the top includes links like "Web", "MBE", "SH JEB", "Evol & Bioinfo", "UH", "Mail", "News", "Libraries", "Search", "Houston", "Art+Jewelry", "Dreidels", "Services+Shopping", "fun", "Musicals", "Google Ngram Viewer", and "Other Bookmarks". The status bar at the bottom right shows the email address "dgraur@gmail.com".



What is Google Scholar?

Enables you to search specifically for scholarly literature, including peer-reviewed papers, theses, books, preprints, abstracts and technical reports from all broad areas of research.

Use Google Scholar to find articles from a wide variety of academic publishers, professional societies, preprint repositories and universities, as well as scholarly articles available across the web.

A screenshot of the Google Scholar search interface. A red arrow points to the search bar containing the text "grauf". To the right of the search bar is a "Search" button. Below the search bar are three blue hyperlinks: "Advanced Scholar Search", "Scholar Preferences", and "Scholar Help". At the bottom of the interface, the phrase "Stand on the shoulders of giants" is displayed in green text.

grauf

Search

[Advanced Scholar Search](#)
[Scholar Preferences](#)
[Scholar Help](#)

Stand on the shoulders of giants

**Google Scholar
orders your
search results by
how relevant they
are to your query,
so the most
useful references
should appear at
the top of the
page**

This relevance ranking takes into account the: full text of each article, the article's author, the publication in which the article appeared and how often it has been cited in scholarly literature.

Scholar

[book] Fundamentals of molecular evolution
WH Li, D Graur - 1991 - Sunderland, Mass.: Sinauer Associates
[Cited by 372](#) [Web Search](#) - [Library Search](#)

[book] Fundamentals of molecular evolution
D Graur, WH Li - 2000 - Sunderland, Mass.: Sinauer Associates
[Cited by 186](#) - [Web Search](#) - [Library Search](#)

[Patterns of nucleotide substitution in pseudogenes and functional genes](#)
T Gojobori, WH Li, D Graur - J. Mol. Evol, 1982 - ncbi.nlm.nih.gov
Patterns of nucleotide substitution in pseudogenes and functional genes. Gojobori T, Li WH, Graur D. MeSH Terms: Base Sequence; Codon; DNA/genetics*; Evolution* ...
[Cited by 116](#) - [Web Search](#)

[CITATION] Extent of protein polymorphism and the neutral mutation theory
M Nei, D Graur - Evol. Biol, 1984
[Cited by 86](#) - [Web Search](#)

[Is the guinea-pig a rodent?](#)
D Graur, WA Hide, WH Li - Nature, 1991 - ncbi.nlm.nih.gov
The guinea-pig (*Cavia porcellus*), traditionally classified as a New World hystricomorph rodent, often shows anomalous morphological and molecular ...
[Cited by 86](#) - [Web Search](#)

[Phylogenetic position of the order Lagomorpha\(rabbits, hares and allies\)](#)
D Graur, L Duret, M Gouy - Nature, 1996 - ncbi.nlm.nih.gov
Ever since they have been classified as ruminants in the Old Testament (Leviticus 11:6, Deuteronomy 14:7) and equated with hyraxes in the vulgate ...
[Cited by 78](#) - [Web Search](#)

What other DATA can we retrieve from the record?

Scholar

[book] Fundamentals of molecular evolution

WH Li, D Graur - 1991 - Sunderland, Mass.: Sinauer Associates

Cited by 372 - Web Search - Library Search



Scholar

Results 1 - 10 of about 359 citing Li: Fundamentals of molecular evolution. (0.01 sec)

[Archaea and the prokaryote-to-eukaryote transition](#)

JR Brown, WF Doolittle - *Microbiology and Molecular Biology Reviews*, 1997 - mmbr.asm.org

Page 1. MICROBIOLOGY AND MOLECULAR BIOLOGY REVIEWS, 1092-2172/97/\$04.00

0 Dec. 1997, p. 456-502 Vol. 61, No. 4 Copyright © 1997 ...

[Cited by 171](#) - Web Search - [uprm.edu](#) - [plantbio.berkeley.edu](#) - [pubmedcentral.nih.gov](#) - all 7 versions »

[METALLOTHIONEIN: An Intracellular Protein to Protect Against Cadmium Toxicity](#)

CD Klaassen, J Liu, S Choudhuri - *Annual Review of Pharmacology and Toxicology*, 1999 - pharmtox.annualreviews.org

Page 1. Annu. Rev. Pharmacol. Toxicol. 1999. 39:267-94 Copyright © 1999 by Annual

Reviews. All rights reserved METALLOTHIONEIN: An Intracellular ...

[Cited by 142](#) - Web Search - [unc.edu](#) - [ncbi.nlm.nih.gov](#) - [csa.com](#) - all 6 versions »

Scholar

[book] Fundamentals of molecular evolution

WH Li, D Graur - 1991 - Sunderland, Mass.: Sinauer Associates

Cited by 372 - Web Search - Library Search



Web Images Groups News Froogle Local Scholar [more »](#)
"Li" "Fundamentals * molecular evolution" Advanced Search Preferences

Web

Results 1 - 7 of about 4,460 for "[Li](#)" "[Fundamentals](#)" * [molecular evolution](#)".

[Li Lab -- Publications \(Full List\)](#)

Li, W.-H. and D. Graur (1991) **Fundamentals of Molecular Evolution**, Sinauer Associates

... Graur, D. And W.-H. Li (1999) **Fundamentals of Molecular Evolution**, ...

pondside.uchicago.edu/~lilab/pubs_full.htm - 65k - [Cached](#) - [Similar pages](#)

Scholar

[book] Fundamentals of molecular evolution

WH Li, D Graur - 1991 - Sunderland, Mass.: Sinauer Associates

Cited by 372 - Web Search - Library Search



Find in a Library

powered by  WorldCat

Google

Find in a Library

WWW

Local Libraries

Fundamentals of molecular evolution

By: [Wen-Hsiung Li; Dan Graur](#)

Type: English : Book : Non-fiction

Publisher: Sunderland, Mass. : Sinauer Associates, ©1991.

ISBN: 0878934529

Subjects: [Molecular evolution](#). | [Evolution](#). | [Molecular Biology](#). | [Molecular Sequence Data](#). | [évolution moléculaire](#).

Related: [Title/Author Search](#)

Find Libraries with item

Postal code, state, province, or country:

5. Google Book Search

The screenshot shows the Google Book Search homepage. At the top is the Google logo with "Book Search" in green and "BETA" below it. Below the logo is a navigation bar with links: Web, Images, Groups, News, Froogle, Maps, Scholar, and more ». To the right of the navigation bar are links for Advanced Book Search and Google Book Search Help. A search bar contains the placeholder "Search Books". Below the search bar is a dropdown menu showing "All books" selected. A large green banner at the bottom reads "Search the full text of books and discover new ones."

[Web](#) [Images](#) [Groups](#) [News](#) [Froogle](#) [Maps](#) [Scholar](#) [more »](#)

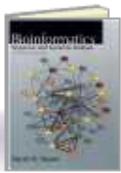
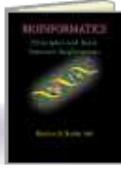
bioinformatics

Search: All books Full view books

[Search Books](#) [Advanced Book Search](#) [Google Book Search Help](#)

Book Search

Books 1 - 53 with 53 pages on bioinform

	Bioinformatics: Sequence and Genome Analysis - Page 1 by David W Mount - Science - 2004 - 692 pages ... 2 Book Guide for Computational Scientists, 3 Basics for Training Students in Bioinformatics, 5 Glossary Terms, 5 WHAT IS BIOINFORMATICS?, ... Limited preview - Table of Contents - Index - About this book	Spons
	Bioinformatics: A Practical Approach - Page xi by D Higgins, W Taylor - Science - 2000 - 249 pages Bioinformatics ... Limited preview - Table of Contents - Index - About this book	M.S. in Bioinform Johns Hopkins Ur School of Arts & S www.biotechnolog
	Bioinformatics: Principles and Basic Internet Applications by Hassan A Sadek - Technology - 2004 - 106 pages ... BASIC APPLICATION OF BIOINFORMATICS FOR PROTEIN 45 CHAPTER 5 FUNDAMENTAL ... OF BIOINFORMATICS FOR NUCLEOTIDE 75 REFERENCES 95 ADDITIONAL RESOURCES Limited preview - Table of Contents - About this book	Bioinformatics Bioinformatics N Bio-IT World - Sub www.Bio-ITworld.c
	Bioinformatics: Databases and Algorithms - Page 1 by N Gautham - 2006 - 248 pages Just as we may consider biochemistry as dealing with metabolic pathways, so we	bioinformatics Find out more ab applications for M www.apple.com
		Bioinformatics C Management Con Custom Software www.3rdmill.com
		Informatics Educa On-campus and d Grad programs in www.phsu.edu/dr

Start working:

Search Google Books

How many times is the tail of the giraffe mentioned in *On the Origin of Species* by Mr. Darwin?

6. Web of science

http://http://apps.webofknowledge.com.ezproxy.lib.uh.edu/WOS_GeneralSearch_input.do?product=WOS&search_mode=GeneralSearch&SID=4FB7LbbLgDMhG9fDiLh&preferencesSaved=



Go to mobile site

| Sign In | Marked List (0) | My EndNote Web | My ResearcherID | My Citation Alerts | My Saved Searches | Log Out | Help

All Databases Select a Database Web of Science Additional Resources

Search Author Finder Cited Reference Search Advanced Search Search History

Web of ScienceSM

Search

Example: oil spill* mediterranean in Topic

AND Example: O'Brian C* OR OBrian C* in Author

Need help finding papers by an author? Use [Author Finder](#).

AND Example: Cancer* OR Journal of Cancer Research and Clinical Oncology in Publication Name

Add Another Field >>

Searches must be in English

Current Limits: (To save these permanently, [sign in or register](#).)

Timespan
 All Years (updated 2011-08-21)
 From 1955 to 2011 (default is all years)

Citation Databases : Science Citation Index Expanded (SCI-EXPANDED); Social Sciences Citation Index (SSCI); Arts & Humanities Citation Index (A&HCI); Conference Proceedings Citation Index- Science (CPCI-S); Conference Proceedings Citation Index- Social Science & Humanities (CPCI-SSH)

Adjust your search settings
 Adjust your results settings

View in: 简体中文 | English | 日本語

© 2011 Thomson Reuters | Acceptable Use Policy | Please give us your [feedback](#) on using Web of Knowledge.

University of Houston



Maintenance Alert

We understand that you may have questions about differences you may be noting between the new Web of Knowledge and the previous version. Please take a look at our [Frequently Asked Questions](#) page for more information.

Support, Tools, Tips

Training & Support

- Download quick Recorded Training
- Access additional Training Resources
- More questions? Consult the [Help files](#).



What's new in Web of Knowledge?

- ResearcherID is now searchable from within Web of ScienceSM.
- Automatic spelling variations and all new Author Finder in Web of ScienceSM.
- More of What's New

[Access the previous version of Web of Knowledge]

Featured Tips

- Get a complete view of citation activity ([view demo](#)).
- Identify citation trends graphically with Citation Report ([view demo](#)).
- Find out about ResearcherID / Web of

[Go to mobile site](#)[Sign In](#) | [Marked List \(0\)](#) | [My EndNote Web](#) | [My ResearcherID](#) | [My Citation Alerts](#) | [My Saved Searches](#) | [Log Out](#) | [Help](#)[All Databases](#)[Select a Database](#)[Web of Science](#)[Additional Resources](#)[Search](#) | [Author Finder](#) | [Cited Reference Search](#) | [Advanced Search](#) | [Search History](#)

Web of ScienceSM

Cited Reference Search (Find the articles that cite a person's work)

[View our Cited Reference Search tutorial.](#)

Step 1: Enter information about the cited work. Fields are combined with the Boolean AND operator.

* Note: Entering the volume, issue, or page in combination with other fields may reduce the number of cited reference variants found.

<input type="text" value="graurd"/>	in	<input type="button" value="Cited Author"/>
-------------------------------------	----	---

Example: O'Brian C* OR OBrian C*



<input type="text"/>	in	<input type="button" value="Cited Work"/>
----------------------	----	---

Example: J Comp* Appl* Math* (journal abbreviation list)



<input type="text"/>	in	<input type="button" value="Cited Year(s)"/>
----------------------	----	--

Example: 1943 or 1943-1945

[Add Another Field >>](#)

Searches must be in English

Current Limits: (To save these permanently, [sign in or register](#).)

Timespan All Years (updated 2011-08-21) From 1955 to 2011 (default is all years)**Citation Databases :** Science Citation Index Expanded (SCI-EXPANDED); Social Sciences Citation Index (SSCI); Arts & Humanities Citation Index (A&HCI); Conference Proceedings Citation Index- Science (CPCI-S); Conference Proceedings Citation Index- Social Science & Humanities (CPCI-SSH)**Adjust your search settings****Adjust your results settings**

[Search](#)[Author](#)[Finder](#)[Cited Reference Search](#)[Advanced Search](#)[Search History](#)

Web of ScienceSM

[**<< Back to previous page**](#)

Cited Reference Search (Find the articles that cite a person's work)

[View our Cited Reference Search tutorial.](#)

Step 2: Select cited references and click "Finish Search."

Hint: Look for [cited reference variants](#) (sometimes different pages of the same article are cited or papers are cited incorrectly).

CITED REFERENCE INDEX

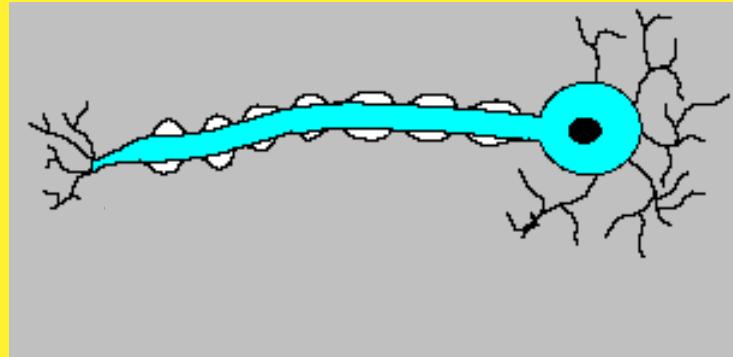
References: 1 - 50 of 193

◀◀ Page of 4 ▶▶

[Go](#)

[Select Page](#) [Select All*](#) [Clear All](#) [Finish Search](#)

Select References	Cited Author	Cited Work [SHOW EXPANDED TITLES]	Year	Volume	Page	Article ID	Citing Articles **	View Record
<input type="checkbox"/>	Abulafia, S...Graur, D	ISRAEL J PLANT SCI	1996	44	77		1	View Record
<input type="checkbox"/>	Armon, A...Graur, D	J MOL BIOL	2001	307	447	10.1006/jmbi.2000.4474	198	View Record
<input type="checkbox"/>	Arnason, U...Graur, D	J MOL EVOL	1996	43	41	10.1007/BF02352298	21	View Record
<input type="checkbox"/>	Barki, Y...Graur, D	MAR BIOL	2000	136	37	10.1007/s002270050005	13	View Record
<input type="checkbox"/>	Barki, Y...Graur, D	MAR ECOL-PROG SER	2002	231	91	10.3354/meps231091	17	View Record
<input type="checkbox"/>	Ben-Dor, A...Graur, D	J COMPUT BIOL	1998	5	377	10.1089/cmb.1998.5.377	21	View Record
<input type="checkbox"/>	BREIMAN, A...GRAUR, D	THEOR APPL GENET	1991	82	201		9	View Record
<input type="checkbox"/>	BREIMAN, A...GRAUR, D	ISRAEL J PLANT SCI	1995	43	85		10	View Record
<input type="checkbox"/>	Cohen, N...Graur, D	MOL BIOL EVOL	2005	22	1260	10.1093/molbev/msi115	37	View Record
<input type="checkbox"/>	Crozier, R. H....Graur, D.	MOL PHYLOGENET EVOL	1995	4	20	10.1006/mpev.1995.1003	29	View Record
<input type="checkbox"/>	Dagan, T...Graur, D	MOL BIOL EVOL	2006	23	310	10.1093/molbev/msj036	28	View Record
<input type="checkbox"/>	Dagan, T...Graur, D	MOL BIOL EVOL	2005	22	496		9	View Record
<input type="checkbox"/>	Dagan, T...Graur, D	MOL BIOL EVOL	2002	19	1022		29	View Record
<input type="checkbox"/>	Dagan, T...Graur, D	NUCLEIC ACIDS RES	2004	32	D489	10.1093/nar/gkh132	36	View Record
<input type="checkbox"/>	Elhai, Eran...Graur, Dan	BIOL DIRECT	2010	5		ARTN 10	1	View Record
<input type="checkbox"/>	Elhai, Eran...Graur, Dan	COMPUT BIOL CHEM	2008	32	147	10.1016/j.combiolchem.2007.11.003	4	View Record
<input type="checkbox"/>	Elhai, Eran...Graur, Dan	MOL BIOL EVOL	2010	27	1015	10.1093/molbev/msp307	2	View Record
<input type="checkbox"/>	Elhai, Eran...Graur, Dan	MOL BIOL EVOL	2009	26	1829	10.1093/molbev/msp100	6	View Record
<input type="checkbox"/>	Elhai, E...Graur, D	MOL BIOL EVOL	2006	23	1	10.1093/molbev/msj006	17	View Record



Databases

**Sequence
Databases**

**Bibliographic
Databases**

Clinical Databases

Integrated Databases

Structural Databases

Sequence Databases

Nucleotide Databases:

EMBL: European Molecular Biology Laboratory

Genbank

DDBJ: DNA Data Bank of Japan

Current Release: 18,324,138 entries

*Release/Up
dates*

International repository for all nucleotide sequences submitted by researchers

Accession numbers are unique to each entry.

One alphabetical character is followed by five digits, or two alphabetical characters are followed by six digits.

Sequence Databases

Nucleotide Databases

RefSeq: Reference Sequence

Current Release: 93,285
entries

NC_123456

Complete Prokaryote Genome

Complete Eukaryote
Chromosome

NG_123456

Homo sapiens Genomic Region

A database of non-redundant reference sequences standards, including genomic DNA contigs, mRNAs and proteins for known genes. Contributions are taken from the NCBI and collaborative sequencing efforts.

NM_123456

mRNA of several organisms, including *Homo sapiens*, *Mus musculus*, *Rattus*

Those accession numbers beginning with X indicate model entries produced as a result of the Genome Annotation process.

Sequence Databases

Protein Databases.

SwissProt: Swiss Protein

Current Release: 115,105 entries

Entry names are often the name of the gene followed by the species.

Accession numbers are of the following format:

[O,P,Q] [0-9] [A-Z, 0-9] [A-Z, 0-9] [A-Z, 0-9] [0-9],

e.g. P26367 (PAX6_HUMAN)

Contains translated sequences from EMBL, adaptations from PIR, extracted from the literature and directly submitted by researchers. Annotation is high quality and the data is cross-referenced to other databases.

Sequence Databases

Protein Databases.

TrEMBL: Translated EMBL

Current Release: 632,013 entries

SpTrEMBL & RemTrEMBL

Acts as a supplement to SwissProt and contains translated EMBL sequences with automatic annotation. TrEMBL entries are manually annotated before being entered into SwissProt.

Remaining TrEMBL contains entries that will never be incorporated into SwissProt. These include: immunoglobulins; T-cell receptors; small fragments; synthetic sequences; CDS not coding for real proteins; patent application sequences

SwissProt TrEMBL contains entries which will eventually be integrated into the SwissProt database. SwissProt accession numbers have been assigned.

Sequence Databases

Protein Databases.

PIR: Protein Information Resource

**Current Release: 283,175
entries**

The PIR is a computer system offering both peptide and nucleotide sequences designed to aid protein identification.

Although much of the protein information in the PIR has been integrated into SwissProt, it may contain some unique sequences.

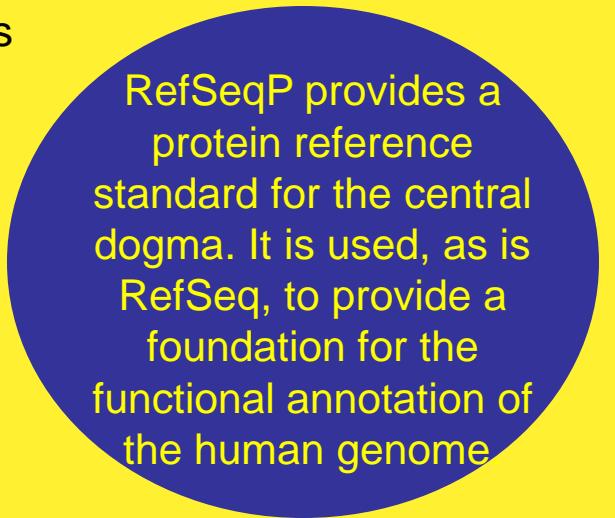
Sequence Databases

Protein Databases.

RefSeqP: Reference Sequence Proteins

Current Release: 402,006 entries

Accession numbers for all proteins are of the format: NP_123456



RefSeqP provides a protein reference standard for the central dogma. It is used, as is RefSeq, to provide a foundation for the functional annotation of the human genome

Sequence Databases



Searching for a sequence:

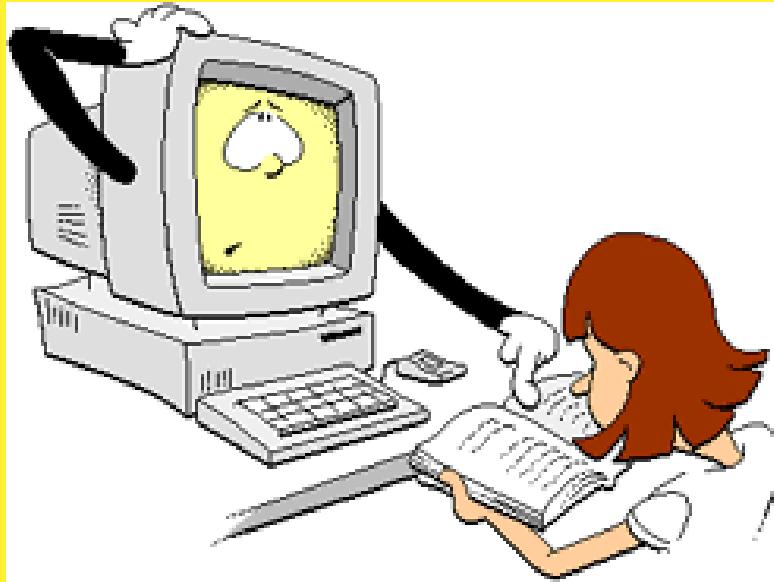
Text Search:

Use text with a boolean operator

BRCA1 & BRCA2 – searches for BRCA1 **AND** BRCA2

BRCA1 | BRCA2 – searches for one gene **OR** the other

BRCA1 ! BRCA2 – searches for BRCA1 **BUT NOT** BRCA2



Computers are
THICK!

Database entries often presented as **flatfiles**

Each piece of information is on a separate line,
distinguished by a code. Computers index this code,
so they can search for the relevant entry.

EMBL entry for a sequence fragment implicated in Human Breast Cancer

Identification	ID	AY144588 standard; DNA; HUM; 68 BP.
Accession	AC	AY144588;
Sequence Version	SV	AY144588.1
Date	DT	23-SEP-2002 (Rel. 73, Created)
	DT	23-SEP-2002 (Rel. 73, Last updated, Version 1)
Description	DE	Homo sapiens truncated breast and ovarian cancer susceptibility protein
Keyword	DE	(BRCA1) gene, partial cds.
Organism	KW	.
Source	OS	Homo sapiens (human)
Organism Classification	OC	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;
	OC	Eutheria; Primates; Catarrhini; Hominidae; Homo.

Reference Number RN [1]

Reference Position RP 1-68

Reference Author RA Rajkumar T., Soumitra N., Nirmala Nancy K., Shanta V.;

Reference Title RT "Novel 5bp deletion in BRCA1 gene in South Indian family";

Reference Location RL Unpublished.

 RN [2]

 RP 1-68

 RA Rajkumar T., Soumitra N., Nirmala Nancy K., Shanta V.;

 RT ;

 RL Submitted (27-AUG-2002) to the EMBL/GenBank/DDBJ databases.

 RL Molecular Oncology, Cancer Institute (WIA), Canal Bank Road, Adyar, RL Chennai, TN 600020, India

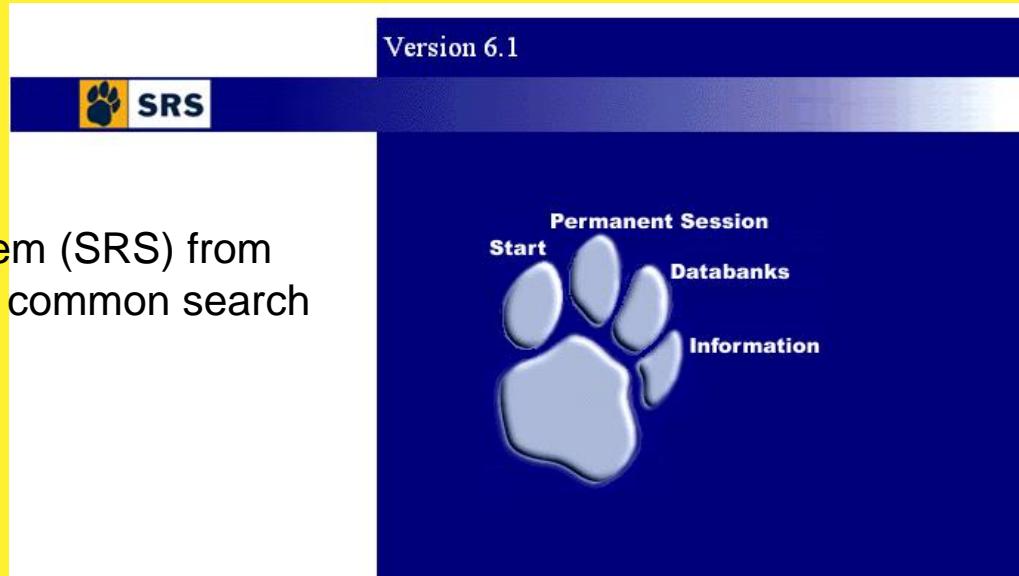
Feature Table Header	FH	Key	Location/Qualifiers
	FH		
Feature Table Data	FT	source	1..68
	FT	/country	"India: South India"
	FT	/db_xref	"taxon:9606"
	FT	/note	"identical sequence found in daughter with breast
	FT		cancer"
	FT	/sex	"female"
	FT	/organism	"Homo sapiens"
	FT	/isolation_source	"mother with breast cancer"
	FT	/dev_stage	"adult"
	FT		mRNA 68
	FT	/gene	"BRCA1"
	FT	/product	"truncated breast and ovarian cancer susceptibility protein"
	FT		

Sequence Header

FT CDS <1..68
FT /codon_start=3
FT /note="contains premature stop codon due to frameshift
FT caused by deletion"
FT /product="truncated breast and ovarian cancer susceptibility protein"
FT /protein_id="AAN10167.1"
FT /translation="EAASGCESETSVSEDCSGLSE"
FT exon 1..68
FT /number=12
FT /gene="BRCA1"
FT misc_feature 61..62
FT /note="site of deletion"
FT /gene="BRCA1"
SQ Sequence 68 BP; 19 A; 12 C; 23 G; 14 T; 0 other;
gtgaaggcagc atctgggtgt gagagtgaaa caagcgtctc tgaagactgc tcagggctat 60
cagagtga
// 68

Searching the databases with a “search engine”:

The Sequence Retrieval System (SRS) from LION Bioscience AG is a very common search tool

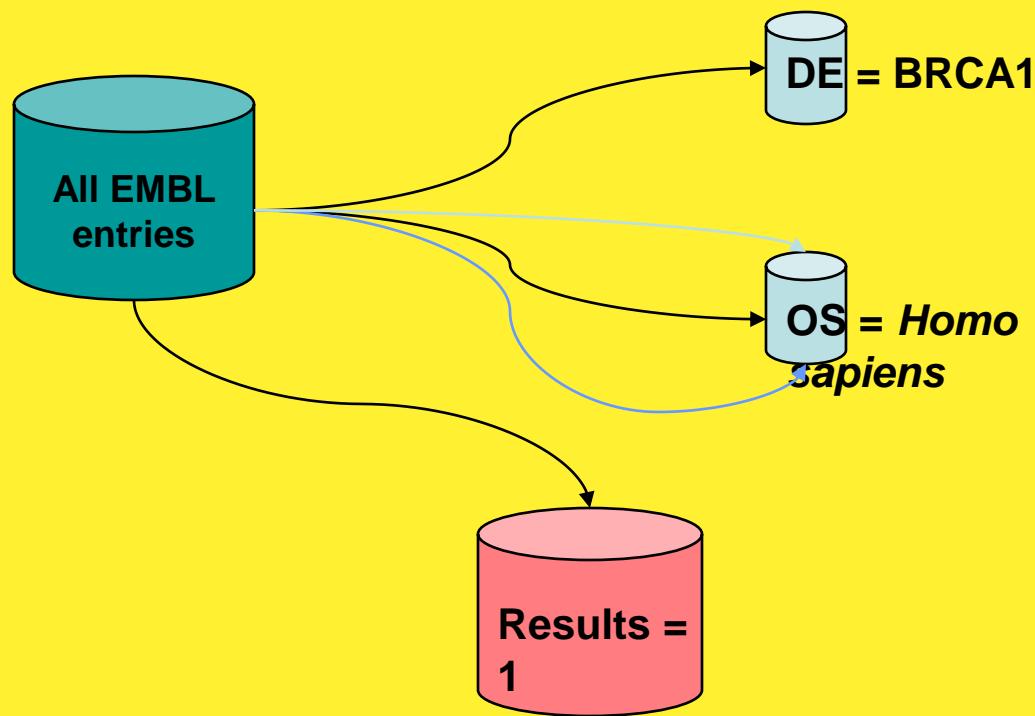


The NCBI in the USA has its own search engine called Entrez.

A screenshot of the NCBI Entrez search interface. At the top, there's a navigation bar with links for PubMed, Nucleotide, Protein, Genome, and Structure. Below that is a search bar with dropdown menus for "Search" (set to "PubMed"), "for", and "Go" and "Clear" buttons. To the left of the search bar is a sidebar with links for "About Entrez", "SITE MAP", "PubMed Help", "Protein Help", "Entrez Help", and "The Entrez Databases". The main content area has a large "Entrez" logo with the text "search and retrieval system". It also contains the text "Entrez is a retrieval system for searching several linked databases." and "It provides access to:" followed by a list of database resources: "PubMed: The biomedical literature (PubMed)", "Nucleotide sequence database (Genbank)", "Protein sequence database", "Structure: three-dimensional macromolecular structures", "Genome: complete genome assemblies", "PopSet: Population study data sets", "Taxonomy: organisms in GenBank" (marked as NEW), and "OMIM: Online Mendelian Inheritance in Man" (marked as NEW).

To search for the BRCA1 gene in Homo sapiens in the EMBL database:

BRCA1 [DE] & Human [OC]

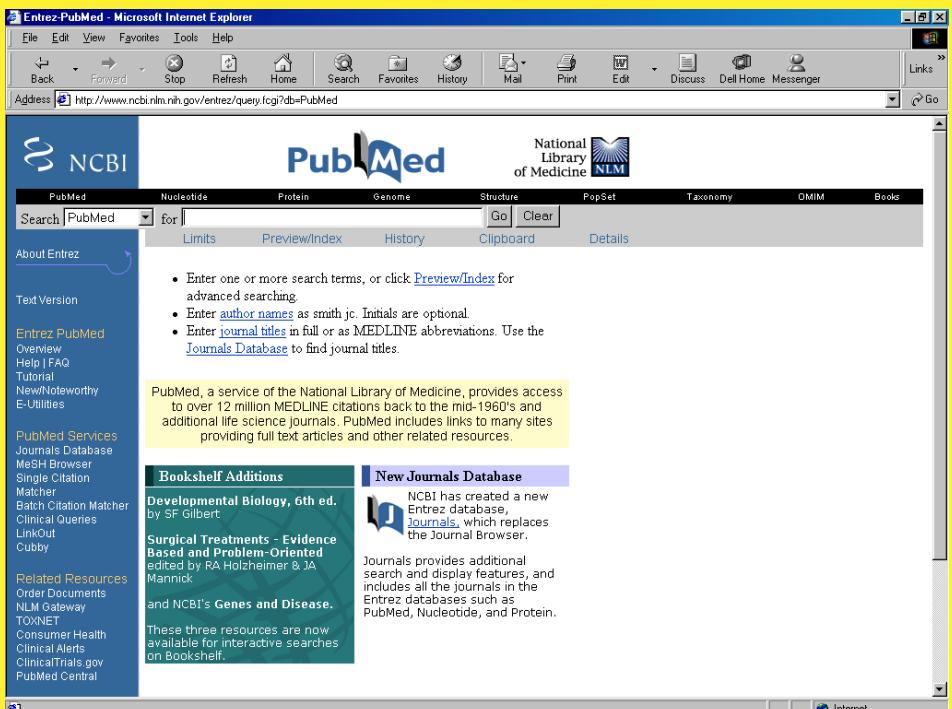


Bibliographic Databases

Used for searching for reference articles

For all (loosely) medically related papers, use *PubMed* from the NCBI

Currently holds over 12 million MEDLINE entries.



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

Bibliographic Databases

Other scientific databases may include:

Web of Science: <http://wos.mimas.ac.uk>

Free to academics, but requires username and password

PubCrawler: <http://www.pucrawler.ie>

Free to academics, will search journals and sequences daily, weekly or monthly and alert the user when results are found corresponding to their search

Clinical Databases

Generally contain information from the Human.

Human Gene Mutation Database, Cardiff, UK:

<http://www.hgmd.org>

Registers known mutations in the human genome and the diseases they cause.

dbSNP, Bethesda, USA:

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

The largest database for single nucleotide polymorphisms. Accession numbers used in dbSNP are not compatible with other SNP databases.

Integrated Databases

These contain overview information garnered from a variety of different databases, and then offer links to further information.

GeneCards: <http://bioinformatics.weizmann.ac.il/cards>

An extremely thorough overview of a particular gene, with links to various other integrated and clinical databases.

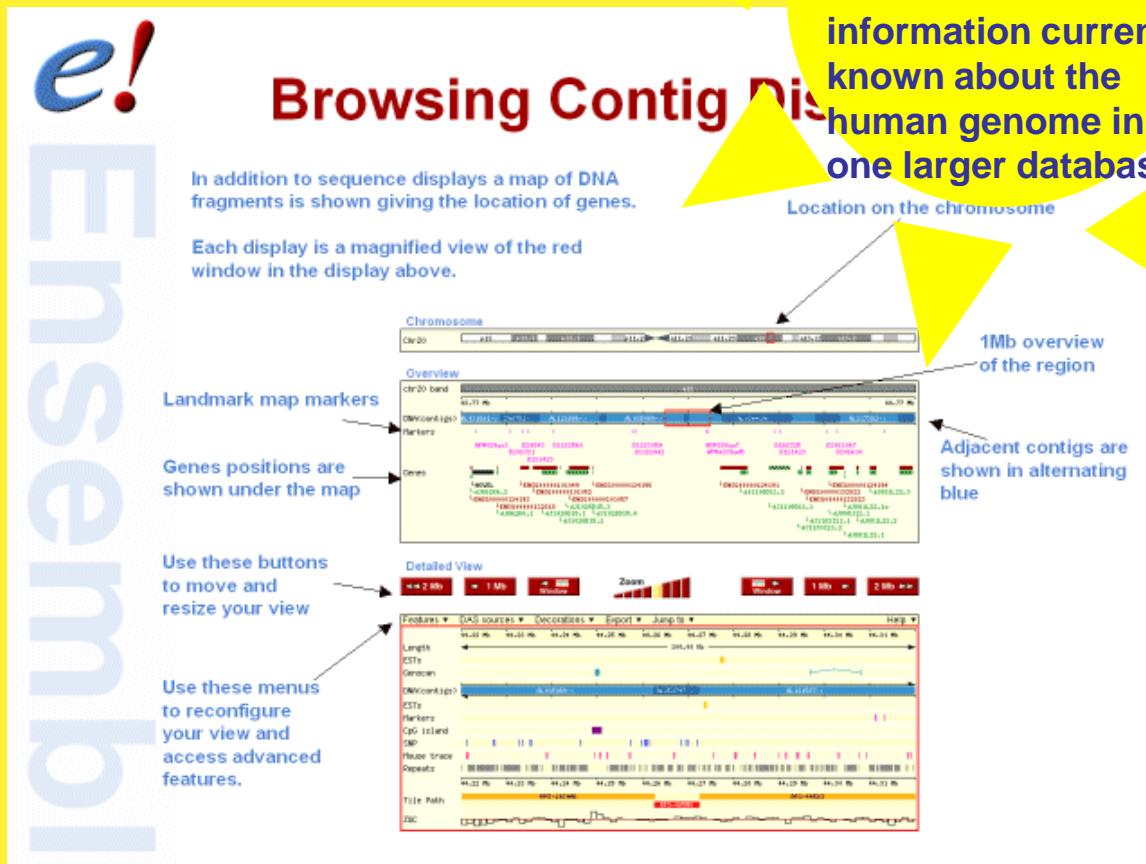
Interpro: <http://www.ebi.ac.uk/interpro>

Integration of individual protein resources PRINTS; PROSITE; SMART; ProDom; Pfam; TIGRfam into one database. A search will scan entries of each and output results.

Integrated Databases

Ensembl: <http://www.ensembl.org>

A joint project by EBI
and Sanger to
annotate all the
information currently
known about the
human genome in
one larger database



Structural Databases

Tertiary protein structure prediction is possibly the Holy Grail of bioinformatics.

PDB: Protein DataBank, New Jersey, USA

<http://www.rcsb.org/>

EMSD: EBI Macromolecular Structure Database

<http://www.ebi.ac.uk/msd/index.html>

Management and distribution of data on macromolecular structures in close collaboration with the PDB.

This houses a collection of 3D coordinates of each atom in a protein, allowing the structure to be displayed by viewing software. Protein structures are submitted by individual researchers and have been determined by x-ray diffraction, or NMR.

Structural Databases

SCOP: Structural Classification of
Proteins

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

**Current Release: 686 folds; 1073 Superfamilies; 1827
Families
representing 15,979 PDB
entries**

CATH: Classification, Architecture, Topology, Homology

http://www.biochem.ucl.ac.uk/bsm/cath_new/

**Current Release: 36,480
Domains**

Alignment III

PAM Matrices

PAM250 scoring matrix

Scoring Matrices

$S = [s_{ij}]$ gives score of aligning character i with character j for every pair i, j .

C	12			
S	0	2		
T	-2	1	3	
P	-3	1	0	6
A	-2	1	1	1 2

STPP
CTCA

$$0 + 3 + (-3) + 1$$

$$= 1$$

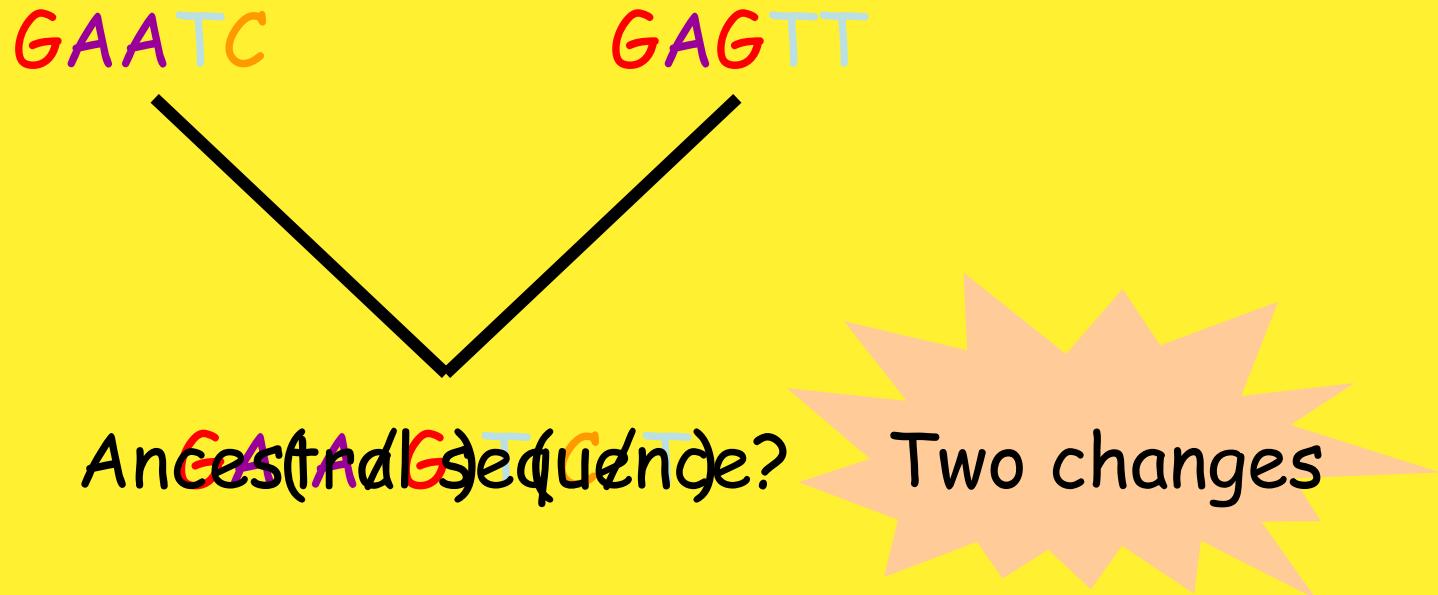
Scoring with a matrix

- Optimum alignment (global, local, end-gap free, etc.) can be found using dynamic programming
 - No new ideas are needed
- Scoring matrices can be used for any kind of sequence (DNA or amino acid)

Types of matrices

- PAM
- BLOSUM
- Gonnet
- JTT
- DNA matrices
- PAM, Gonnet, JTT, and DNA PAM matrices are based on an explicit evolutionary model;
BLOSUM matrices are based on an implicit model

PAM matrices are based on a simple evolutionary model



- Only mutations are allowed
- Sites evolve independently

Log-odds scoring

- What are the odds that this alignment is meaningful?

$$\begin{matrix} X_1 & X_2 & X_3 & \dots & X_n \\ Y_1 & Y_2 & Y_3 & \dots & Y_n \end{matrix}$$

- **Random model:** We're observing a chance event. The probability is

$$\prod_i p_{X_i} \prod_i p_{Y_i}$$

where p_X is the frequency of X

- **Alternative:** The two sequences derive from a common ancestor. The probability is

$$\prod_i q_{X_i Y_i}$$

where q_{XY} is the joint probability that X and Y evolved from the same ancestor.

Log-odds scoring

- Odds ratio:

$$\frac{\prod_i q_{X_i Y_i}}{\prod_i p_{X_i} \prod_i p_{Y_i}} = \prod_i \frac{q_{X_i Y_i}}{p_{X_i} p_{Y_i}}$$

- Log-odds ratio (score):

where

$$S = \sum_i s(X_i, Y_i)$$

$$s(X, Y) = \log\left(\frac{q_{XY}}{p_X p_Y}\right)$$

is the **score** for X, Y . The $s(X, Y)$'s define a **scoring matrix**

PAM matrices: Assumptions

- Only mutations are allowed
- Sites evolve independently
- Evolution at each site occurs according to a simple (“first-order”) Markov process
 - Next mutation depends only on current state and is independent of previous mutations
- Mutation probabilities are given by a **substitution matrix** $M = [m_{XY}]$, where $m_{xy} = \text{Prob}(X \rightarrow Y \text{ mutation}) = \text{Prob}(Y|X)$

PAM substitution matrices and PAM scoring matrices

- Recall that

$$s(X, Y) = \log\left(\frac{q_{XY}}{p_X p_Y}\right)$$

- Probability that X and Y are related by evolution:

$$q_{XY} = \text{Prob}(X) \cdot \text{Prob}(Y|X) = p_X \cdot m_{XY}$$

- Therefore:

$$s(X, Y) = \log\left(\frac{m_{XY}}{p_Y}\right)$$

Mutation probabilities depend on evolutionary distance

- Suppose M corresponds to one unit of evolutionary time.
- Let f be a frequency vector (f_i = frequency of a.a. i in sequence). Then
 - $M \cdot f$ = frequency vector after one unit of evolution.
 - If we start with just amino acid i (a probability vector with a 1 in position i and 0s in all others) column i of M is the probability vector after one unit of evolution.
 - After k units of evolution, expected frequencies are given by $M^k \cdot f$.

PAM matrices

- Percent Accepted Mutation: Unit of evolutionary change for protein sequences [Dayhoff78].
- A PAM unit is the amount of evolution that will on average change 1% of the amino acids within a protein sequence.

PAM matrices

- Let M be a PAM 1 matrix. Then,

$$\sum_i p_i (1 - M_{ii}) = 0.01$$

- **Reason:** M_{ii} 's are the probabilities that a given amino acid does not change, so $(1 - M_{ii})$ is the probability of mutating away from i .

The PAM Family

Define a *family* of substitution matrices — PAM 1, PAM 2, etc. — where PAM n is used to compare sequences at distance n PAM.

$$\text{PAM } n = (\text{PAM } 1)^n$$

Do not confuse with scoring matrices!

Scoring matrices are derived from PAM matrices to yield log-odds scores.

Generating PAM matrices

- **Idea:** Find amino acids substitution statistics by comparing evolutionarily close sequences that are highly similar
 - Easier than for distant sequences, since only few insertions and deletions took place.
- Computing PAM 1 (Dayhoff's approach):
 - Start with highly similar aligned sequences, with known evolutionary trees (71 trees total).
 - Collect substitution statistics (1572 exchanges total).
 - Let m_{ij} = observed frequency (= estimated probability) of amino acid A_i mutating into amino acid A_j during one PAM unit
 - Result: a 20×20 real matrix where columns add up to 1.

Dayhoff's PAM matrix

	<i>A</i>	<i>R</i>	<i>N</i>	<i>D</i>	<i>C</i>
<i>A</i>	9867	2	9	10	3
<i>R</i>	1	9913	1	0	1
<i>N</i>	4	1	9822	36	0
<i>D</i>	6	0	42	9859	0
<i>C</i>	1	1	0	0	9973

All entries $\times 10^4$