# 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/
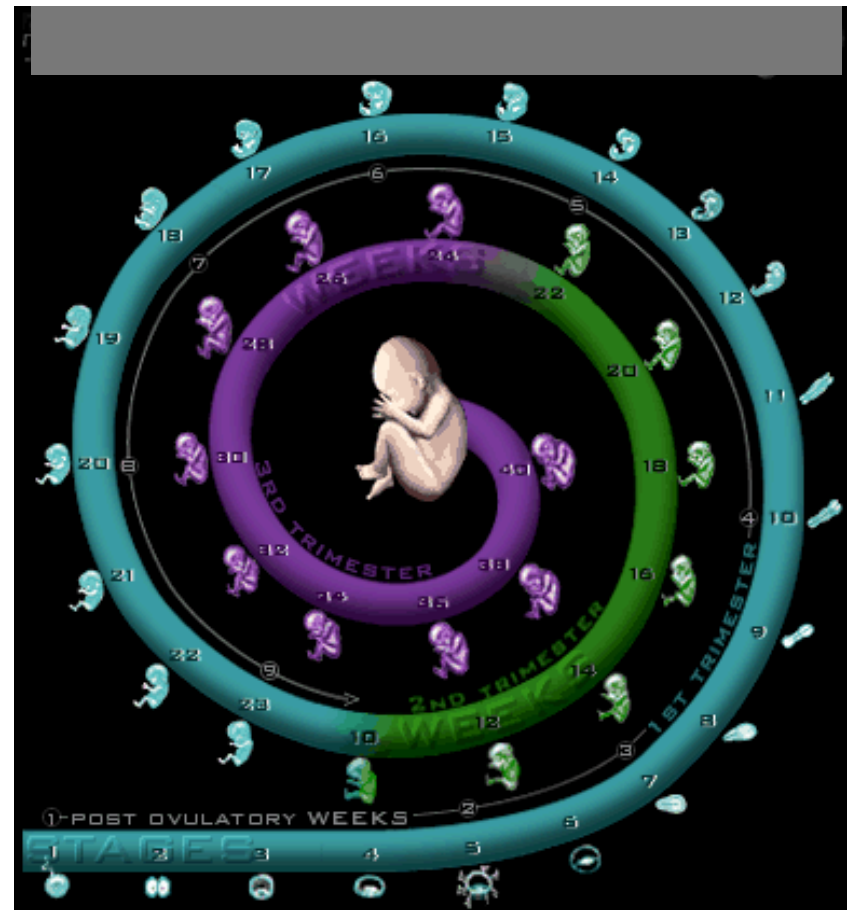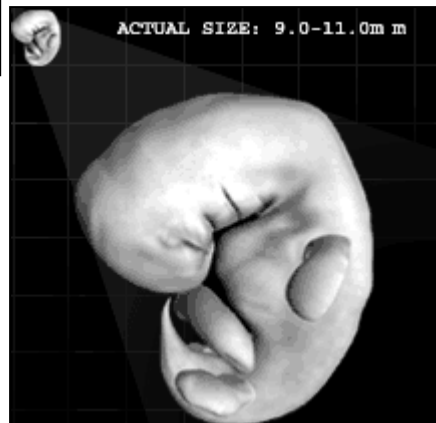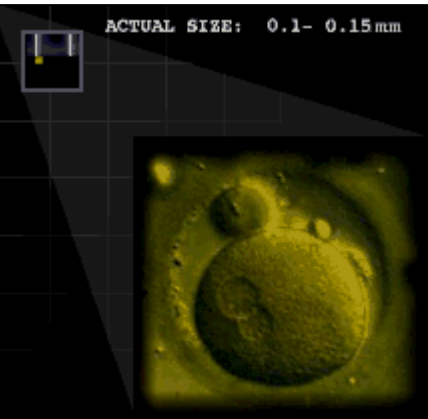
- Birbeck course on protein structure.

# Journals :

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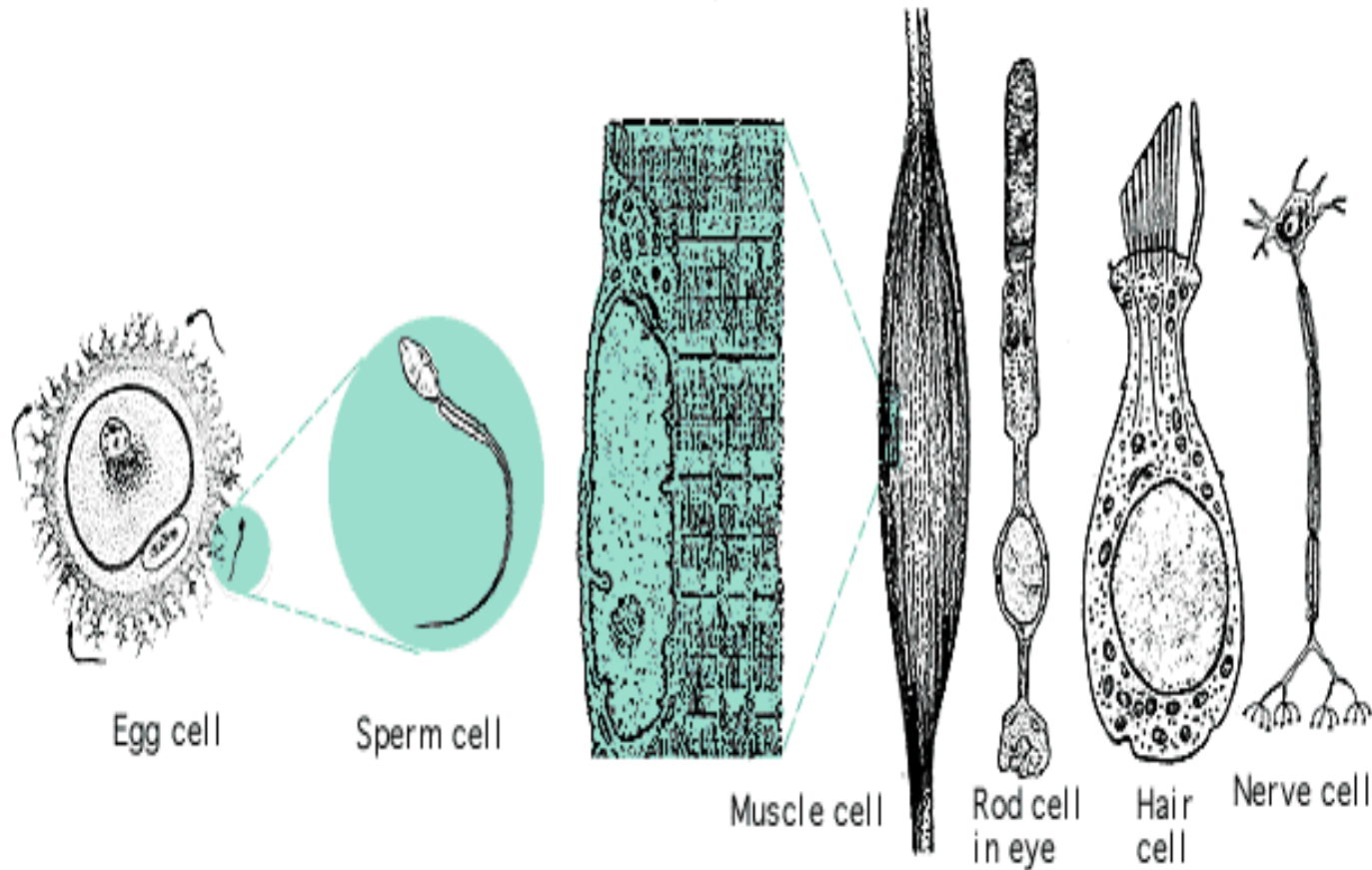- Journal of Molecular Graphics and Modelling.
- Protein Engineering.

# Computational Biology Conferences:

- **<u>ISMB</u>** - International Conference on Intelligent Systems in Molecular Biology.

- **<u>RECOMB</u>** - Int. Conference of Computational Molecular Biology.

- **<u>ECCB</u>** - European Conference on Computational Bio.

- **<u>WABI</u>** - Workshop of Algorithms in Bioinformatics .
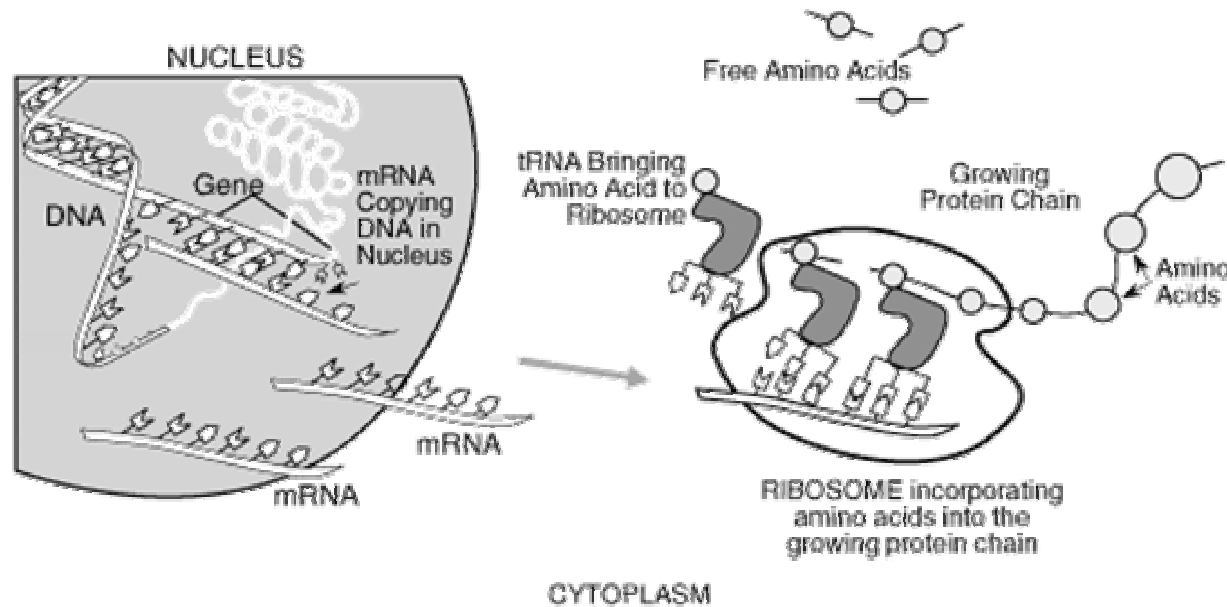
# Cell- the basic life unit



ACTUAL SIZE: 0.1- 0.15 mm

ACTUAL SIZE: 0.1- 0.2 mm

ACTUAL SIZE: 9.0-11.0m m

# Different cell types

Egg cell    Sperm cell    Muscle cell    Rod cell in eye    Hair cell    Nerve cell

# Size of protein molecules (diameter)

- cell         $(1 \times 10^{-6} \, \mathbf{m})$ $\mu$   microns

- ribosome     $(1 \times 10^{-9} \, \mathbf{m})$   nanometers

- protein      $(1 \times 10^{-10} \, \mathbf{m})$ angstroms

# The central dogma

- DNA     --->     RNA     --->     Protein

- {A,C,G,T}     {A,C,G,U}     {A,D,..Y}

- *4 letter alphabets*     *20 letter alphabet*

- Sequence of nucleic acids     seq of amino acids

NUCLEUS

Free Amino Acids

DNA

Gene

mRNA
Copying
DNA in
Nucleus

tRNA Bringing
Amino Acid to
Ribosome

Growing
Protein Chain

Amino
Acids

mRNA

mRNA

RIBOSOME incorporating
amino acids into the
growing protein chain

CYTOPLASM

**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 ) bases 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).
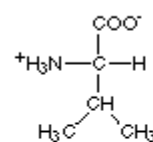- Structural proteins (hair, skin, nails).
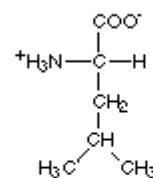
# Amino acids and the peptide bond



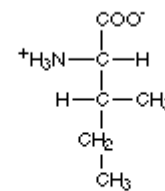$C_\beta$ – first side chain carbon (except for glycine).
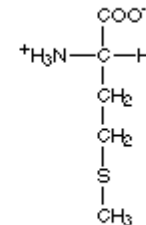


Amino acids with hydrophobic side groups
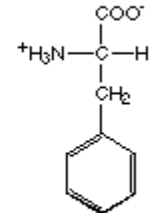
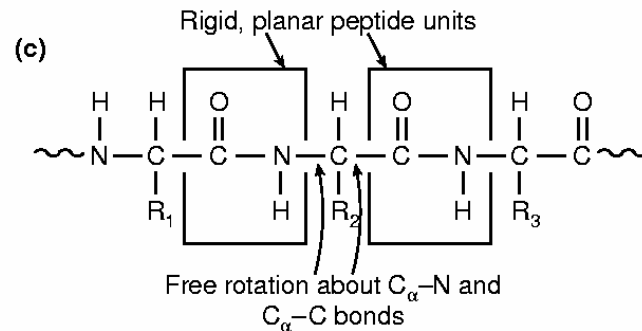Valine
(val)

Leucine
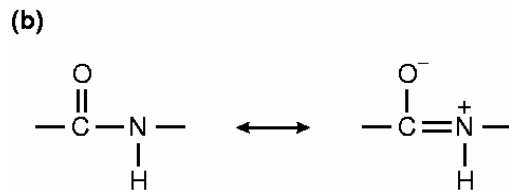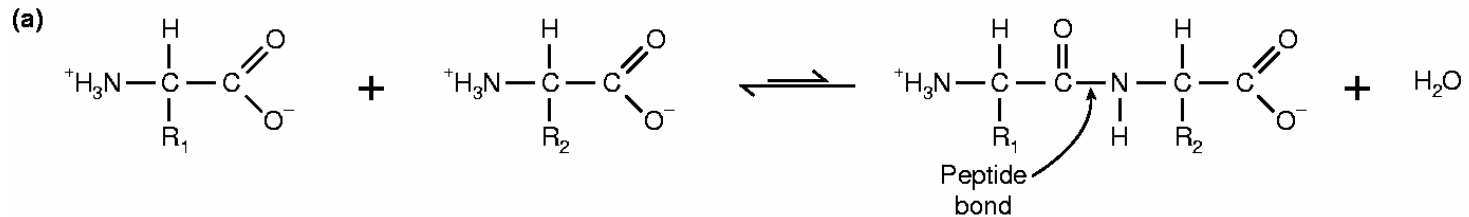(leu)

Isoleucine
(ile)

Methionine
(met)

Phenylalanine
(phe)

# Primary through Quaternary structure

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

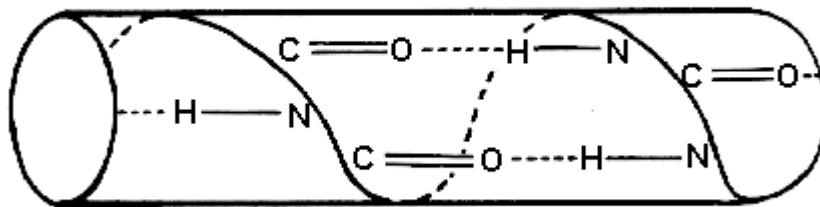- AASGDXSLVEVHXXVFIVPPXIL.....
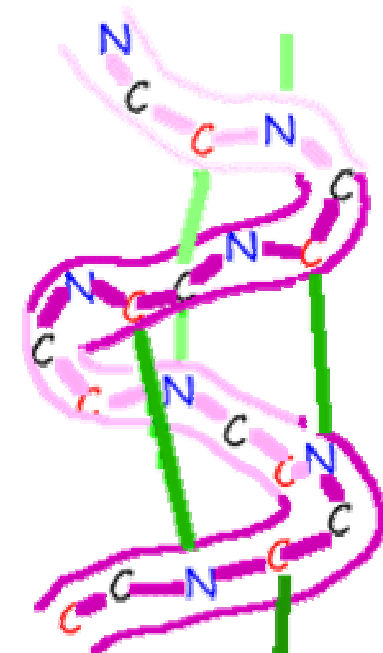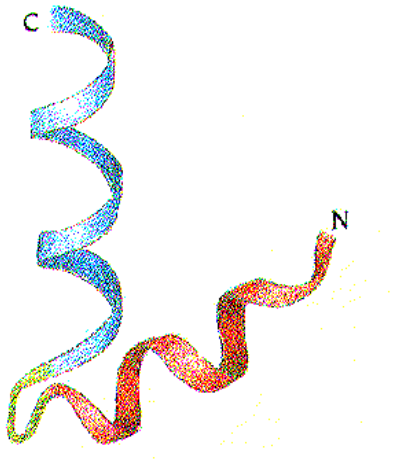
# Folding of the Protein Backbone

# 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
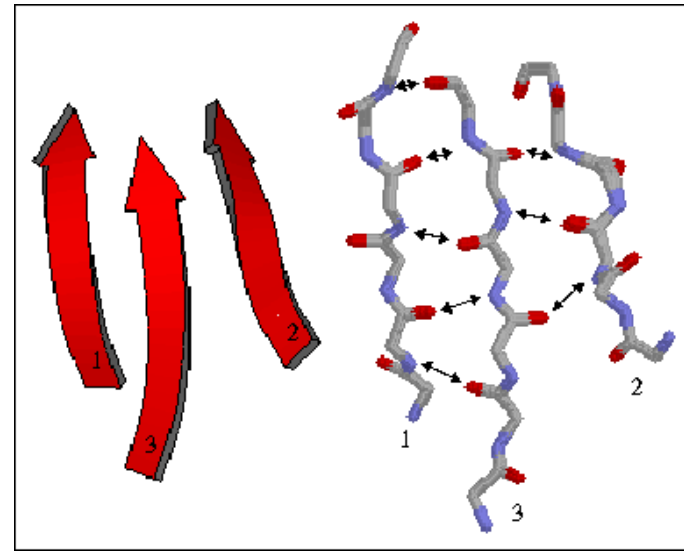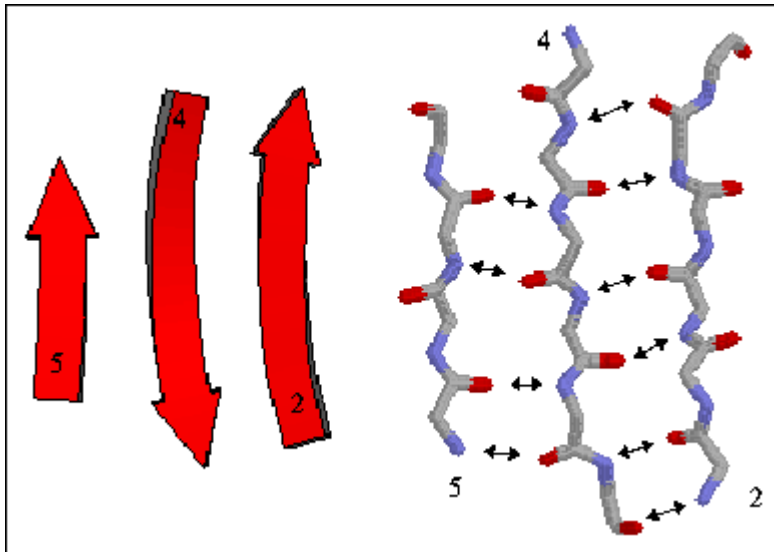
3.6 residues/turn (5.4 A dist.)
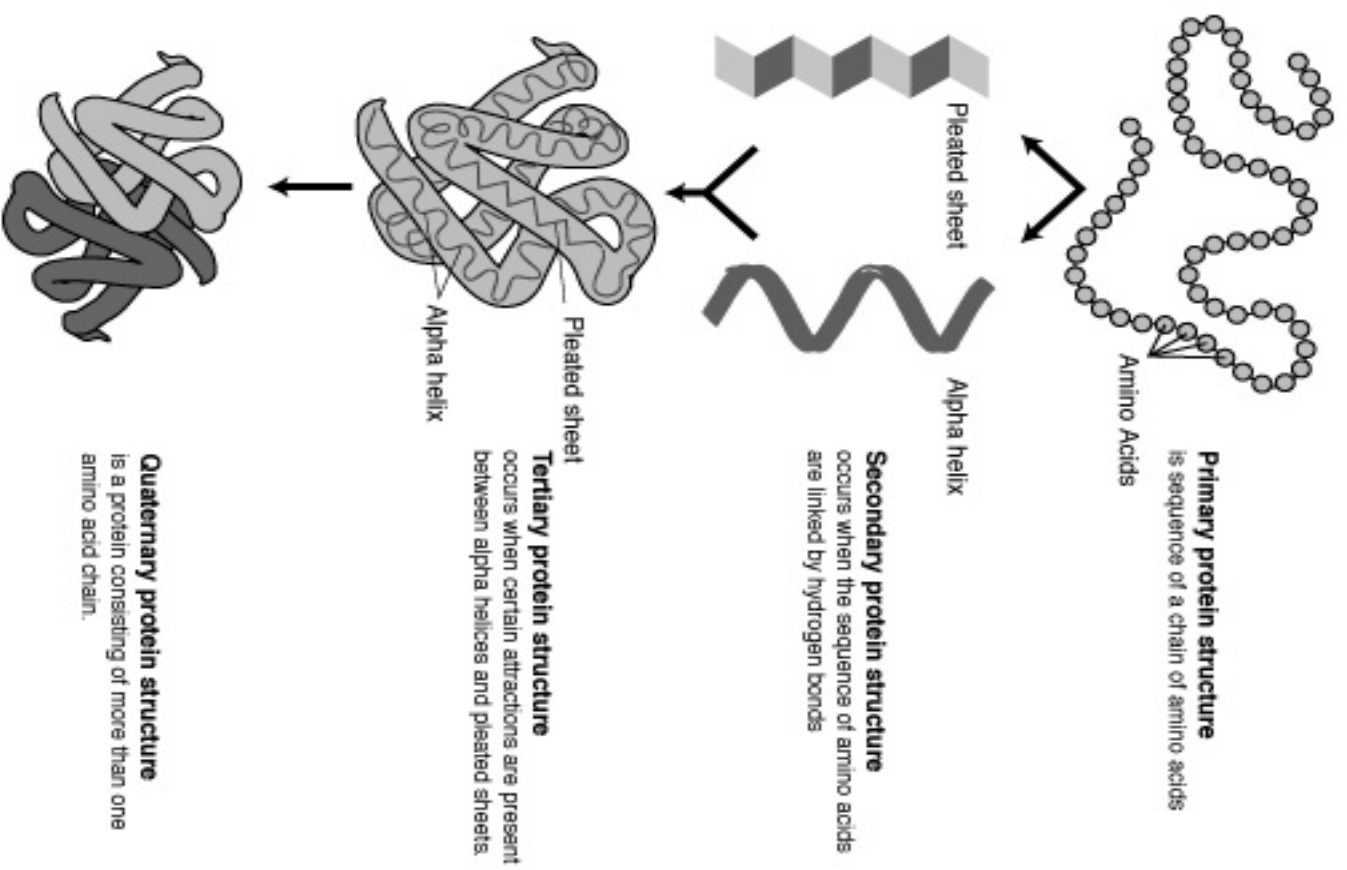
2: Backbone:

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

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

Alpha helix

Pleated sheet

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

Pleated sheet

Alpha helix

**Quaternary protein structure**
is a protein consisting of more than one
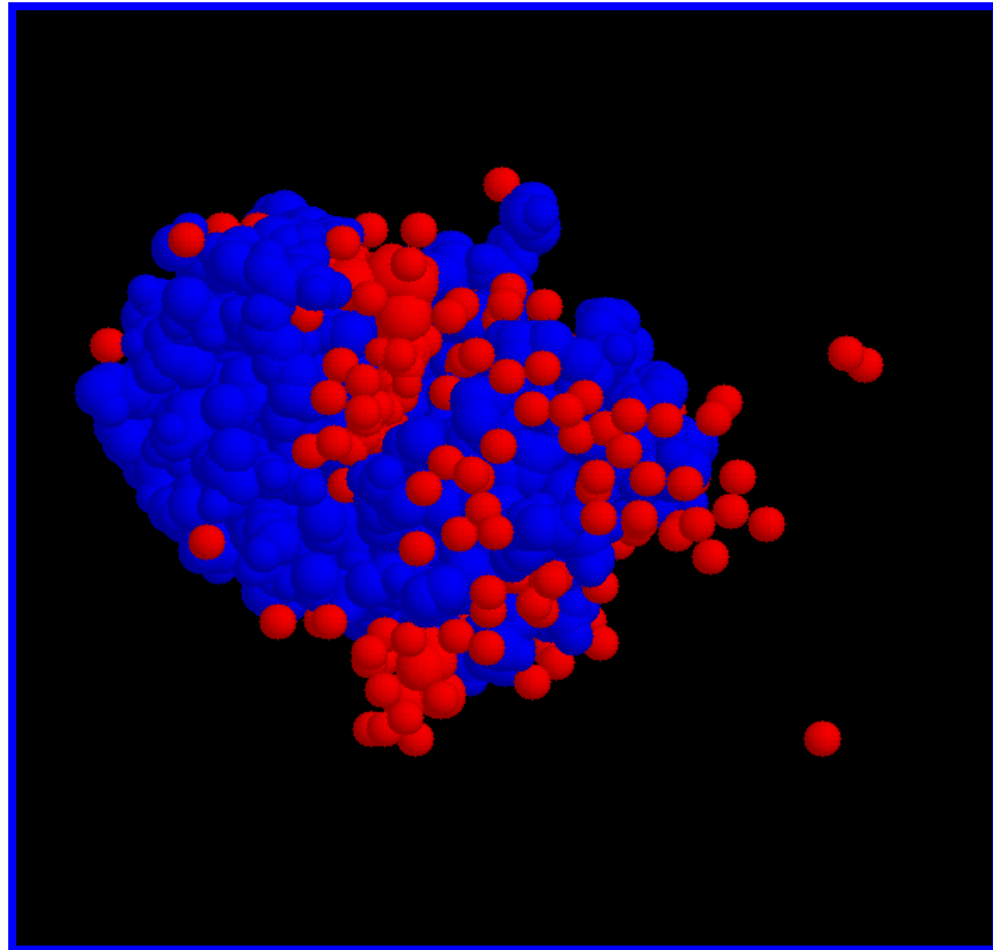amino acid chain.
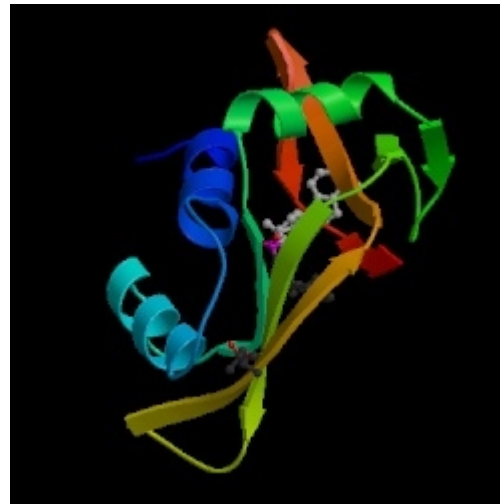
# Wire-frame or ribbons display

# Space-fill display

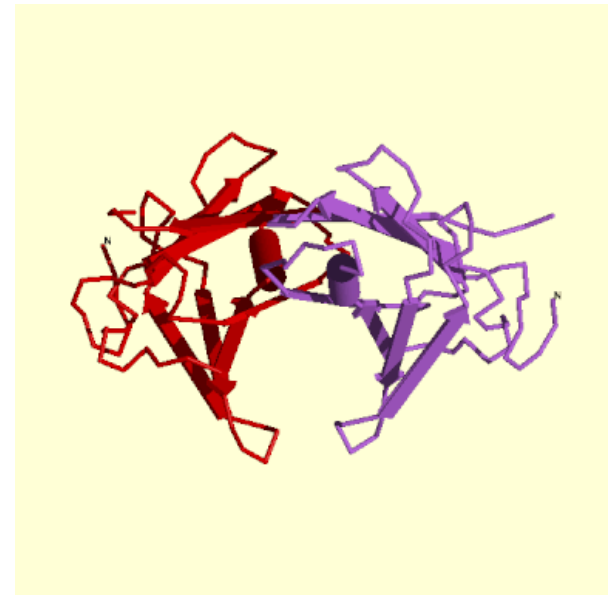# Tertiary structure: full 3D folded structure of the polypeptide chain

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)

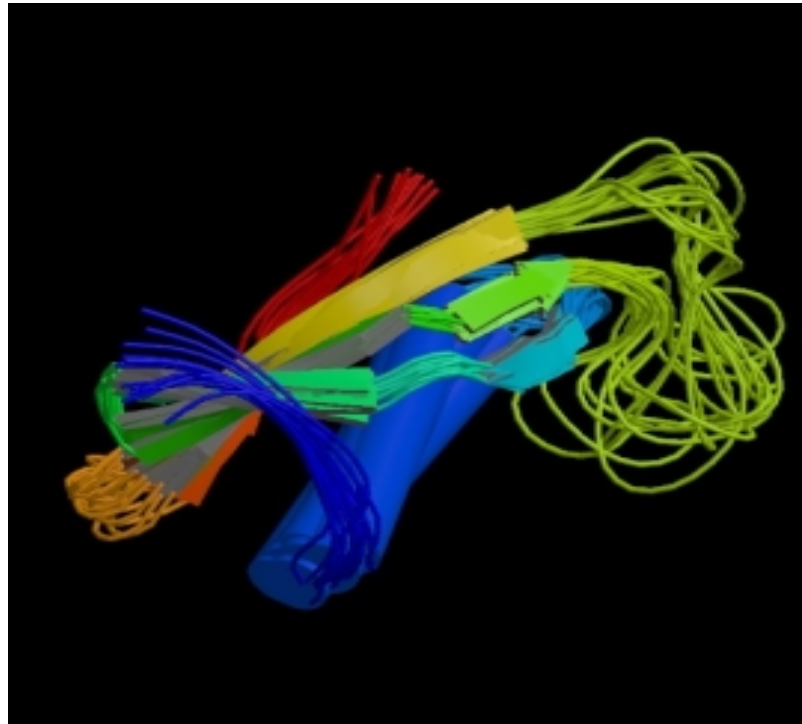- Nano – sensors (?)

# 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.

**P D B** ™

**P R O T E I N   D A T A   B A N K**

Research Collaboratory for Structural Bioinformatics

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**

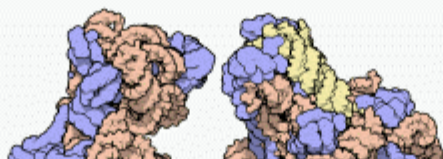*DOWNLOAD* **Retrieve structure files (FTP)**

*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.

# Structure Explorer - 1IRS

## Summary Information

**Summary Information**

View Structure

Download/Display File

Structural Neighbors

Geometry

Other Sources

Sequence Details

Explore [ ]

SearchLite  SearchFields

*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. Feisk**

*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** | $C_9H_{12}N_1O_6P_1$ |

Prof. Haim J. Wolfson   34

# 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

## 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

Prof. Haim J. Wolfson                                    37

# CATH - Protein Structure Classification
## http://www.biochem.ucl.ac.uk/bsm/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.

- **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

```
---GAATTC---
|||||||||||
---CTTAAG---
```



```
---G      AATTC---
||||      |||||
---CTTAA      G---
```

**Normal Hemoglobin**

Glutamate 6   Glutamate 6

**Sickle Hemoglobin**

Valine 6   Valine 6

Valine 6

Valine 6

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

44

# The Structural Genomics Pipeline
# (X-ray Crystallography)

Basic Steps

*Target Selection*

*Crystallomics*
- Isolation,
- Expression,
- Purification,
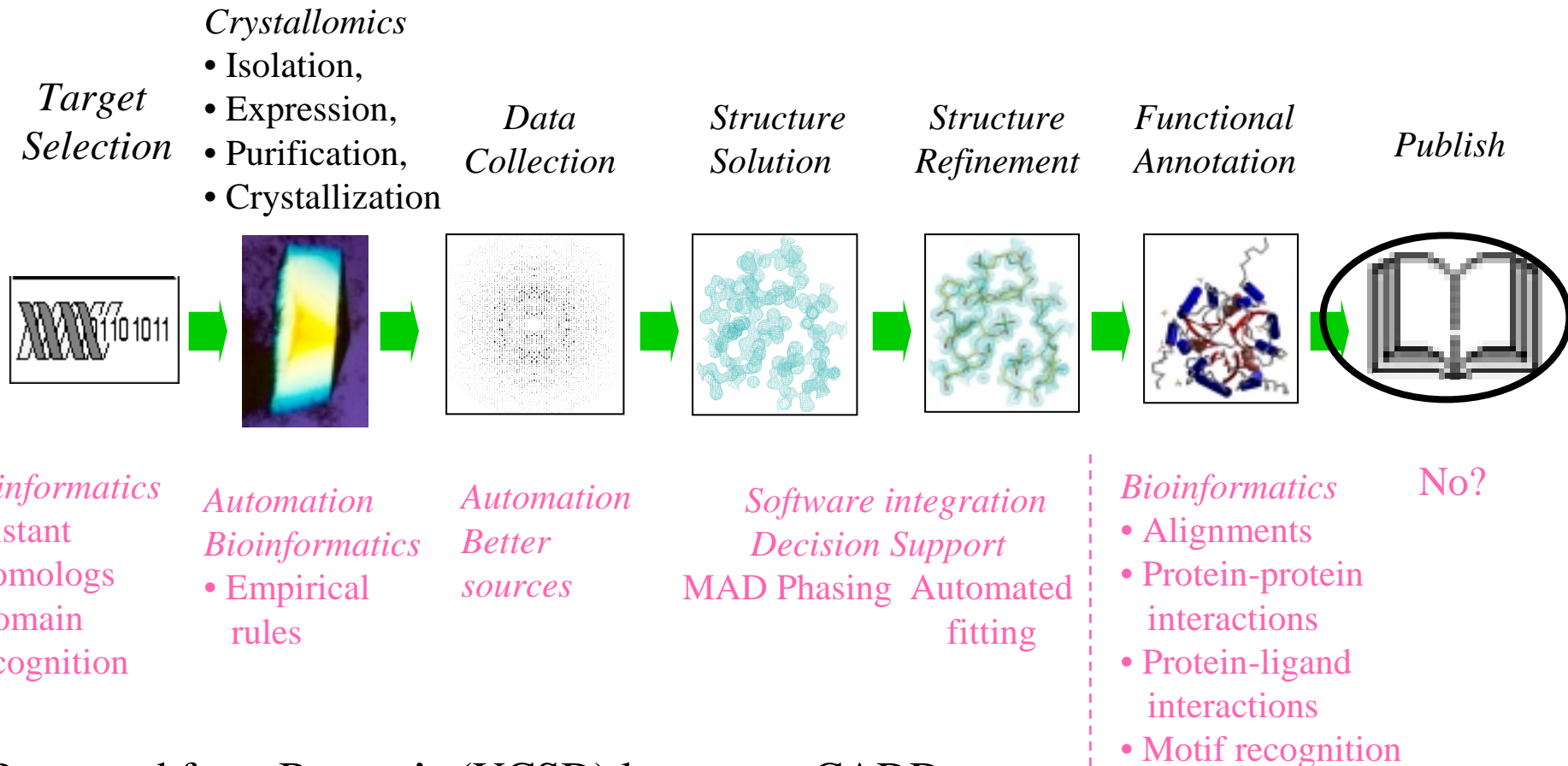- Crystallization

*Data Collection*

*Structure Solution*

*Structure Refinement*

*Functional Annotation*

*Publish*



*Bioinformatics*
- Distant homologs
- Domain recognition

*Automation Bioinformatics*
- Empirical rules

*Automation Better sources*

*Software integration Decision Support*
MAD Phasing   Automated fitting

*Bioinformatics*
- Alignments
- Protein-protein interactions
- Protein-ligand interactions
- Motif recognition

No?

Borrowed from Bourne's (UCSD) lecture on CADD

TAU Structural Bioinformatics Lab
(Wolfson-CS, Nussinov – MB)

**Human Genome Project**
DNA&Protein Sequences

**X-ray cryst. NMR, EM**

**PROTEIN STRUCTURE**

Computer Assisted Drug Design

Biological Function

*OUR EFFORTS*

# _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.

- Pharmaceutics.

- 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.
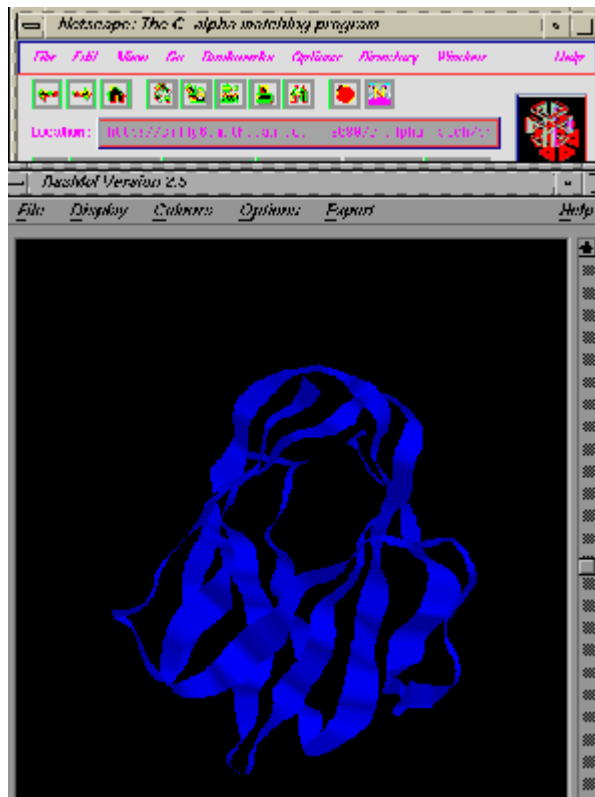
- Graph - theoretic  methods.

# *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_.**

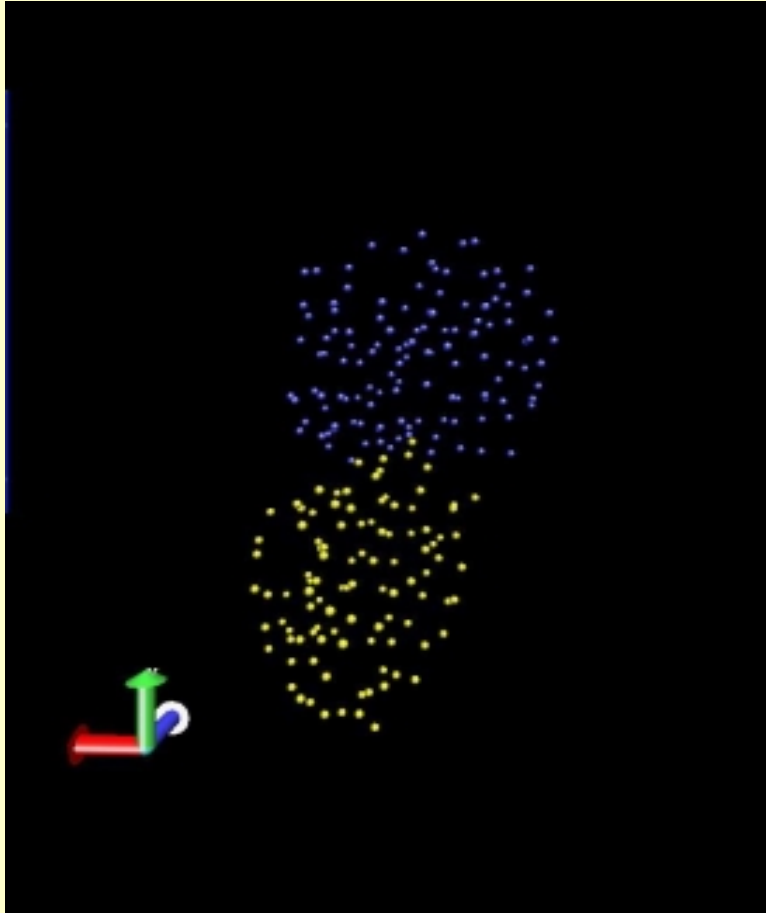# *Algorithmic Solution*



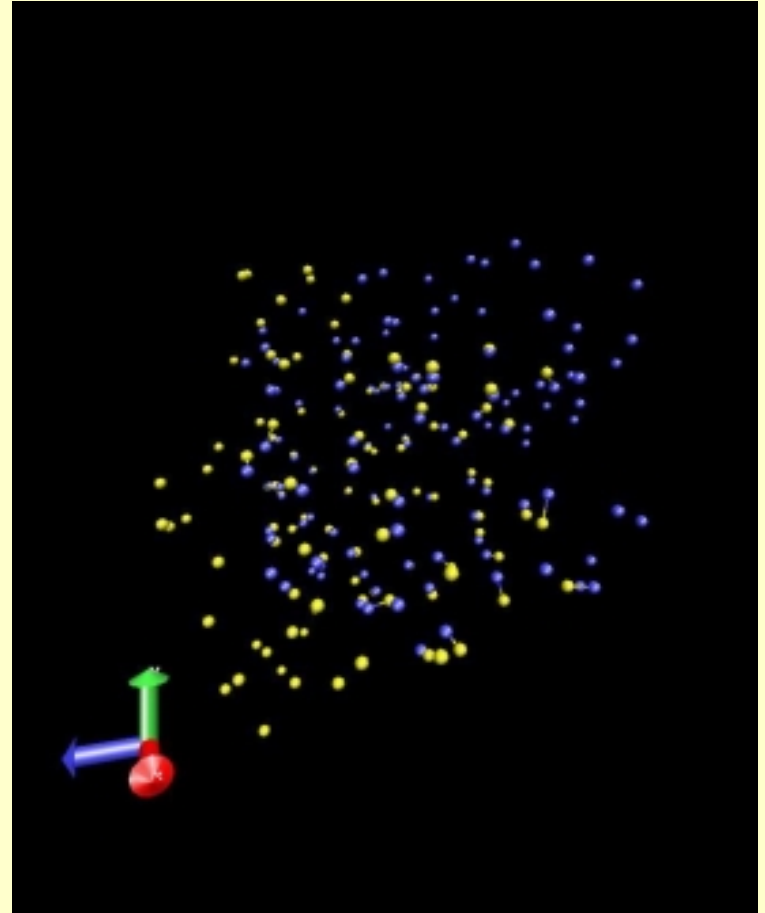**About 1 sec. Fischer, Nussinov, Wolfson ~ 1990.**

# Applications

- Classification of protein databases by structure.

- Search of <span style="color:red">partial and disconnected</span> structural patterns in large databases.

- Detection of structural pharmacophores in an ensemble of drugs.

- Comparison and detection of drug receptor *active sites*.

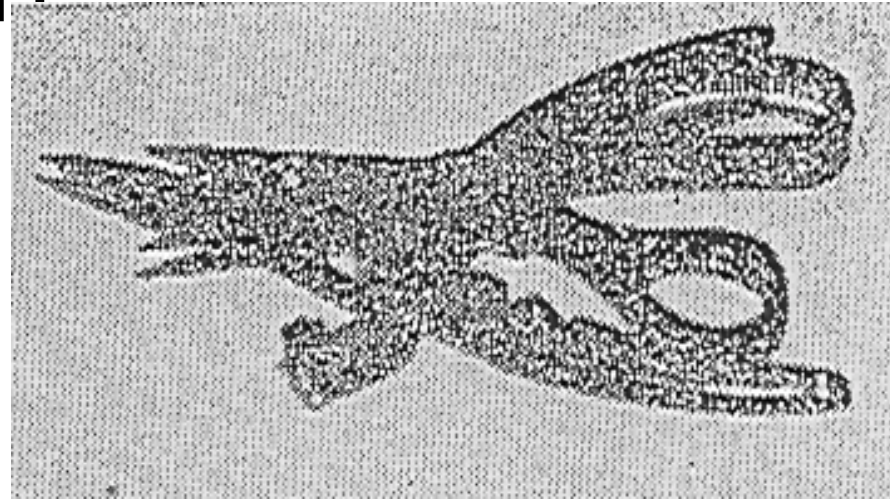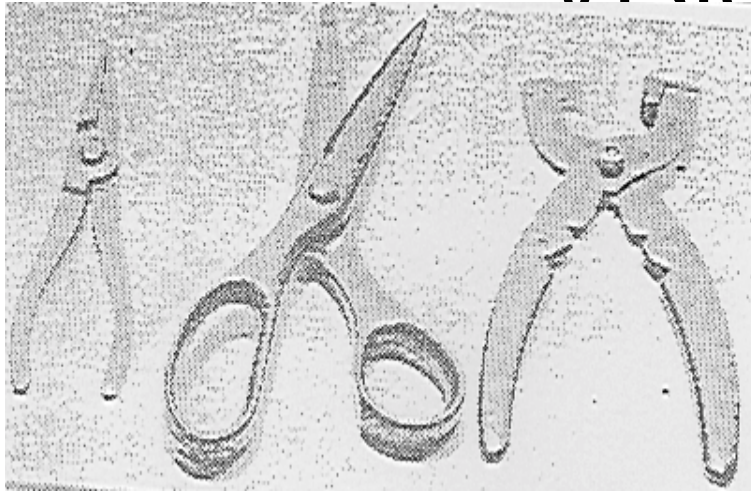# Geometric Matching task = Geometric Pattern Discovery



$C_\alpha$ 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 <span style="color:red">biologically relevant</span> computational techniques.

# Conclusions 2

- The <span style="color:red">b.c.</span> (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.