

## BIOINFORMATICS III

### „Structural Bioinformatics and Genome Analysis“

Dipl.-Ing. Noura Chelbat

Biologist: Molecular Biologist

Phone: +43-732-2468-8898

Room: T725

Consulting hours: e-mail/phone

[chelbat@bioinf.jku.at](mailto:chelbat@bioinf.jku.at)

# Lecture Dates

7 weeks 4 Units if on Tuesdays



## Stundenplan Bioinformatik, SS2009

	MONTAG	DIENSTAG	MITTWOCH	DONNERSTAG	FREITAG
8:30-9:15				365.030 Bioinformatik IV Kusel (UE)	
9:15-10:00		365.029 Bioinformatik III Chelbat (KV)			
10:15-11:00	365.033 Masterarbeitsseminar T 111  4. Sem.	entweder ...			
11:00-11:45					
12:00-12:45					
12:45-13:30					
13:45-14:30		327021 Mathematische Modellierung und wissenschaftliches Rechnen in Biowissenschaften I Beuchler (VO)	365.035 Bioinformatik II Schwarzbauer (UE)	326.007 Algebraische und diskrete Methoden der Biologie RISC	
14:30-15:15					
15:30-16:15	365.034 Bioinformatik II Hochreiter (VO) HS 14	365.029 Bioinformatik III Chelbat (KV)		365.034 Bioinformatik II Hochreiter (VO) T 111	
16:15-17:00		... oder			
17:15-18:00					
18:00-18:45					

## Total Units 28 Possible dates



March :	Tu. 3	2 U. (15.30-17.00)	
	Tu. 10	2 U. (9.15-11.00)	
	Tu. 10	4 U. (15.30-18.45)	
	Total	8 U	
May :	Tu. 5	4U. (15.30-18.45)	
	Tu. 12	4U.(15.30-18.45)	
	Tu. 19	4U.(15.30-18.45)	
	Tu. 26	4U.(15.30-18.45)	
	Total	16 U	
June :	Tu. 9	4U .(15.30-18.45)	Ready!
	Tu.16	4U .(15.30-18.45)	
	Tu.23	4U. (15.30-18.45)	Ready!
	Total	12 U	

## Brief Remind



- Part of curriculum of the master of sciences in Bioinformatics
- Included in the Compulsory modules
- Combined Courses (KV) with mainly theoretical part
- Background : Bridge modules from M1-M5
  - M1 Basics of molecular biology
  - M2 Basics of biochemistry
  - M3 Basics of algorithms and data structure
  - M4 Basics of information systems
  - M5 Basics of mathematics

DNA, RNA, Transcription, Translation, Genetic Code, Promoter, Protein folding, Gene regulation

Purification, Molecular forces, Secondary / Tertiary /quaternary structure, Folding, Molecular dynamics, instrumental analytics



## Molecular and Cell Biology

- Lodish, Berk, Matsudaira, Kaiser, Krieger, Scott, Zipursky & Darnell - Molecular Cell Biology. Fifth edition. W.H. Freeman and Company, New York, USA, 2004.
- Alberts, Johnson, Lewis, Raff, Roberts, Walter –Molecular Biology of the Cell. Fourth edition. Garland Science, Taylor and Francis Group, New York, USA, 2002.
- Mathew, Van Holde and Ahern –Biochemistry. Third edition. Benjamin/Cummings an imprint of Addison Wesley Longman, 1301 Sansome street, San Francisco, CA 94111

## General Bioinformatics

- David W. Mount. Bioinformatics – Sequence and Genome Analysis. Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York, USA, 2004



## Structural Bioinformatics

- Philip E. Bourne and Helge Weissig. Structural Bioinformatics. Wiley-Liss, Hoboken, New Jersey, USA, 2003.
- Michael J. E. Sternberg. Protein Structure Prediction. Oxford University Press, 1996.
- Arthur M. Lesk. Introduction to protein Architecture. Oxford University Press 2003
- Richard A. Friesner. Computational Methods for Protein Folding. Advances in Chemical Physics Volume 120. A John Wiley & Sons, INC. Publication. 2002



## Genome Analysis

- Steen Knudsen. Guide to Analysis of DNA Microarray Data. John Wiley & Sohns, Hoboken, New Jersey, USA, 2004.
- Ernst Wit and John McClure. Statistics for Microarrays. John Wiley & Sohns Ltd., England, 2004.
- Pierre Baldi and G. Wesley Hatfield. DNA Microarrays and Gene Expression From Experiments to Data Analysis and Modeling. Cambridge University Press, United Kingdom, 2002.
- Geoffry J. McLachlan, Kim-Anh Do, and Christophe Ambroise. Analyzing Microarray Gene Expression Data. John Wiley & Sohns Inc., Hoboken, New Jersey, USA, 2004.
- Jerome K. Percus. Mathematics of Genome Analysis. Cambridge University Press, United Kingdom, 2002.

1. Structural bioinformatics: Chapters 1-6
2. Genome analysis: Chapters 7-10

## Goals:

- Main methods in structural bioinformatics and gene analysis: from where we get them and how to use them
- How to choose the proper method from a given pool of approaches
- Adaptation of standard algorithms to the final purpose: combining the information of certain algorithms and biology to build up practical solutions
- How we can use this information to perform searches for the optimal 3D prediction, motifs, expression profiles, pattern regulation ..
- Exercises: SSEs, SCOP classes recognition, DEGs, CNVs, arrays, expression patterns...



## Structural Bioinformatics

Motivation:

From Genome sequencing to amino acids/nucleotides primary structure.  
From amino acids/nucleotides primary structure to 3D Structure Prediction.

2008 In PDB data base 49192 Structures structures

[Feb 24, 2009](#) \_ 56066 Structures

2008 SWISS PROT 356 194 entries sequence

10-Feb-2009 UniProtKB/Swiss-Prot Release 56.8 of : 410 518 entries

Ratio of 1 structure to 7 sequences

Increasing number of methods to predict 3D structures beside sequencing ones  
New approaches based on Machine learning, SVM, NNs, Dynamic programming  
and Distance matrixes.

# Part I: Structural Bioinformatics



## 1D

Linear arrangement of amino acids: chain assembled on the ribosome using the codon sequence on mRNA as a template

## 2D

Secondary structures elements: core elements for protein architecture

- $\alpha$  Helix
- $\beta$  Sheet
- Loops
- Coil coiled
- Turns

## 3D

Functional activity:  
Folding and Post-translational modifications  
Interactions among amino acids side groups  
Chaperones

## Molecular representation and viewers

- Difficulties in transforming all of the important 3D structural information about a molecule into an understandable two-dimensional representation
- A variety of molecular representation formats have been developed each of one is designed to show a particular aspect of a molecule's structure
- To visualize the three-dimensional structure of the molecule and understand the relationship between the structural features and its function
- RasMol, Pymol, Chime, .etc

# Part I: Structural Bioinformatics



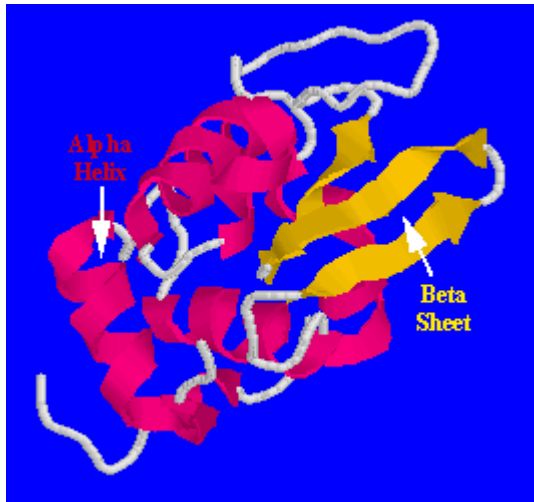
Goals at the end of this part:

- Special amino acids (proline, cysteine,.. ) and special bonds in 2D configuration
- C-C, disulfide bond (-CH<sub>2</sub>-S-S-CH<sub>2</sub>-)
- Recognition of the main types of 2D configurations  $\alpha$  helix,  $\beta$  strands, loops, turns
- Recognition of motifs
- Coil coiled, Zn Fingers, Leucine Zippers...
- Viewers and displayers:
  - Sticks, balls, sticks & balls, ribbon, . .etc
  - Atoms, bonds, angles,..
- Structural comparison and Alignment Methods, Protein Secondary structure prediction
- Molecular Dynamics

# Part I: Structural Bioinformatics

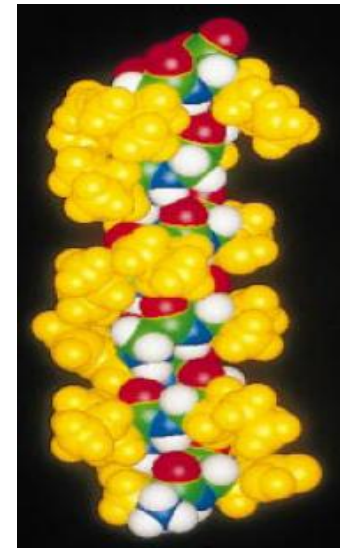
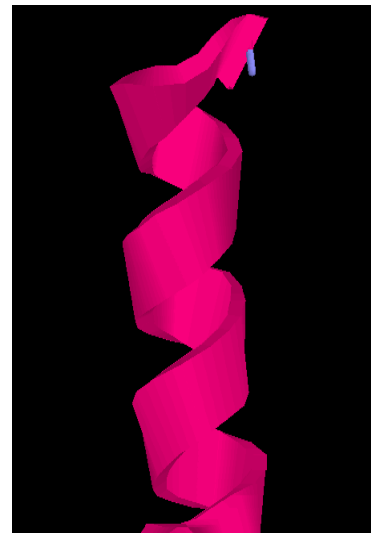


Each picture tells us something different about the structure of the molecule



Lysozyme

- To catch the main SSEs on a subunit
- To see the relative sizes of the atoms in an  $\alpha$  helix by balls representation

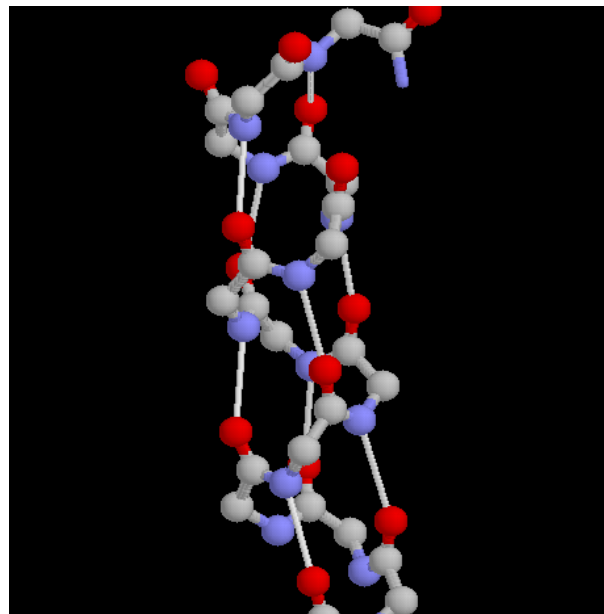
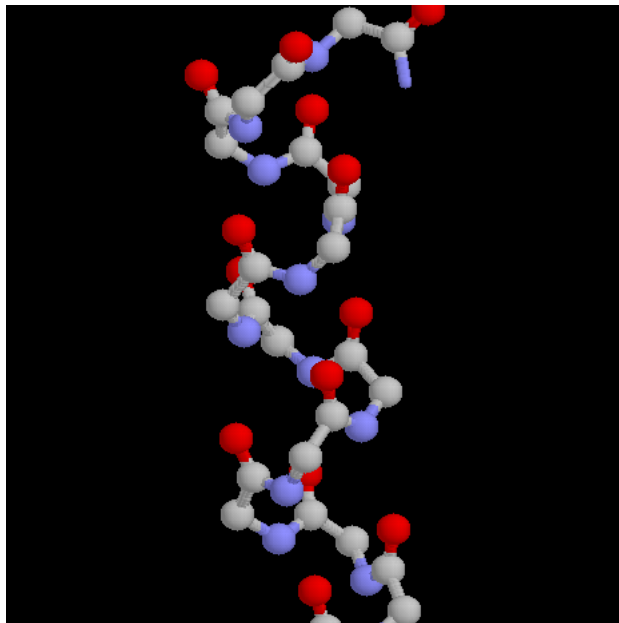


[http://project.bio.iastate.edu/Courses/BIOL202/Proteins/secondary\\_structure.htm](http://project.bio.iastate.edu/Courses/BIOL202/Proteins/secondary_structure.htm)

# Part I: Structural Bioinformatics



## $\alpha$ Helix Ball and Stick View of Lysozyme

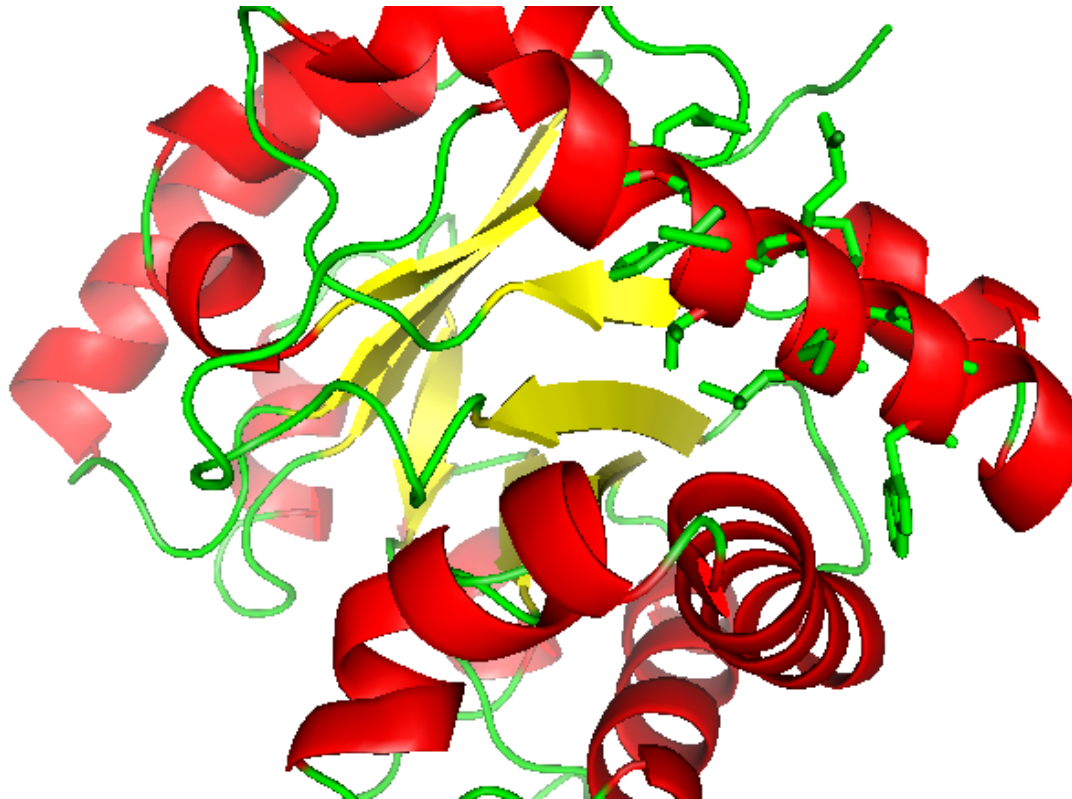


Carbon: Grey  
Oxygen: Red  
Hydrogen: White  
Nitrogen : Blue

To know how the atoms in an  $\alpha$  helix are connected to one another by sticks representation  
Hydrogen bonds location

[http://project.bio.iastate.edu/Courses/BIOL202/Proteins/secondary\\_structure.htm](http://project.bio.iastate.edu/Courses/BIOL202/Proteins/secondary_structure.htm)

# Part I: Structural Bioinformatics



$\alpha$  helix: Red 3.6 amino acid per turn and H bonds between every 4th residue

$\beta$  sheet : Yellow H bonds between a median of 5-10 amino acid

Turns: Green

Loop: Turn between  $\alpha$  helix and  $\beta$  sheet

Ribbon representation: traces the protein backbone does not include the atoms in the side chains of the residues. Bundles of helices.

# Part I: Structural Bioinformatics



For similarity and 3D structure detection

Methods from Bioinformatics I allow for homology and comparative modelling where it is assumed that similar sequences have the same 3D structure

## Troubles

Different sequences from different proteins can fold into similar three-dimensional configurations

- i. No more use of PAM or BLOSSUM matrixes to predict 3D structure on the basis of amino acids substitution because of their standardization
- ii. No more use of methods in which both the core regions and loops are equally represented
- iii. Gaps should be confined to regions not in the core when multiple alignment are used

# Part I: Structural Bioinformatics



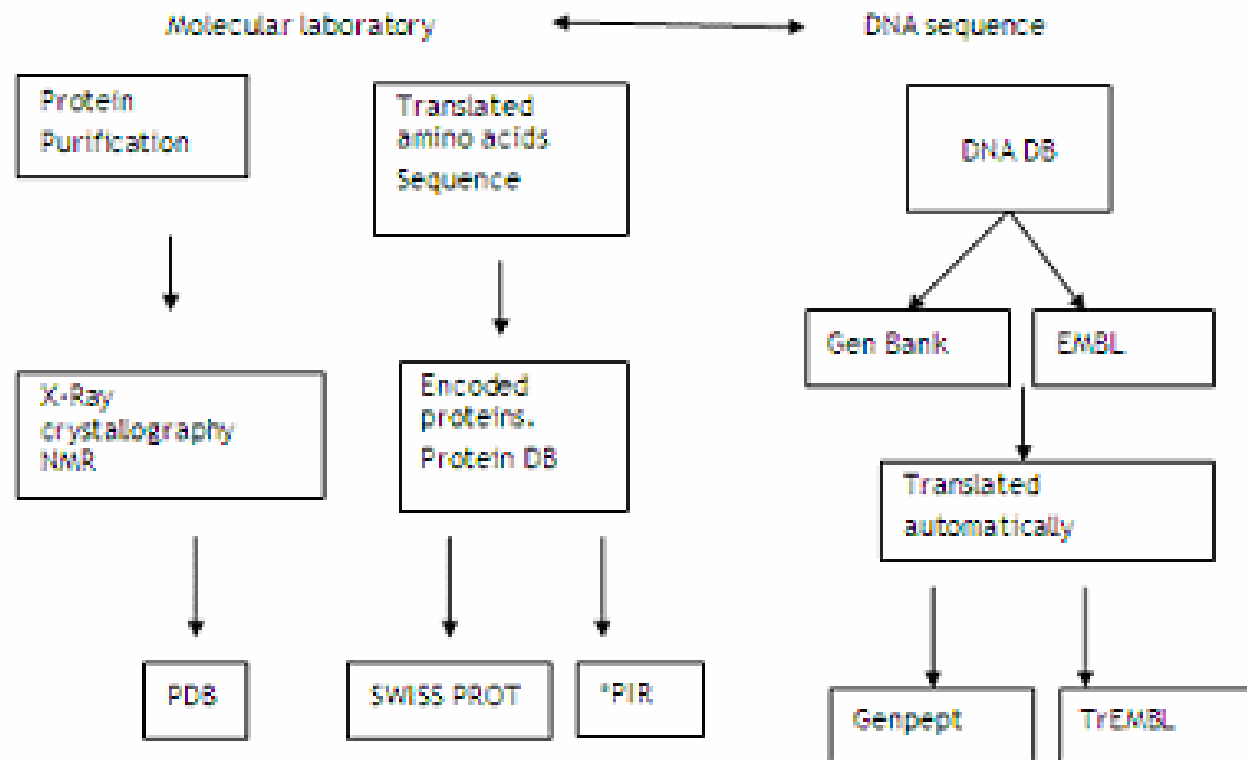
Four steps can be addressed when attempting to get information about an unknown protein structure

- **1st Structure alignment:** based on 3D known structures to find equivalent amino acids residues
- **2nd Structure comparison:** based on shared similarities of two or more proteins when comparing their 3D known structures
- **3rd Structure superposition:** based on preliminary knowledge of positive match of some residue in proteins 1 and 2. The alignment is assumed and the main goal is to search for the best solution to find what amino acids are equivalents to each other
- **4th Structure classification:** based on structural alignment beside other methods to hierarchically assign classes of proteins

# Part I: Structural Bioinformatics



How to go from DNA sequence to DB of structural elements?



type of search from the amino acids chain

# Part I: Structural Bioinformatics



What could be used??

- Comparative Modeling: Sequence to sequence, Sequence to structure (Psi-Blast, SVM, Fisher Kernels..)
- Scoring matrices
- Distance matrices
- HMMs
- Monte Carlo Optimization and Dynamic programming

## Solutions

Direct link between sequence and structure. In all a sequence representation of a known 3D structure is compared with any other sequences up to match the structure predicted by the model

Accuracy of methods to predict  $\alpha$  helix,  $\beta$  strands, coiled coil, turns and loops has an overage of 64-75 % being the highest accuracy for  $\alpha$  helix

# Part I: Structural Bioinformatics



Methods like CE, DALI, SSAP, VAST, SARF2 and COMPARER



## Spatial Arrangement of Backbone Fragments

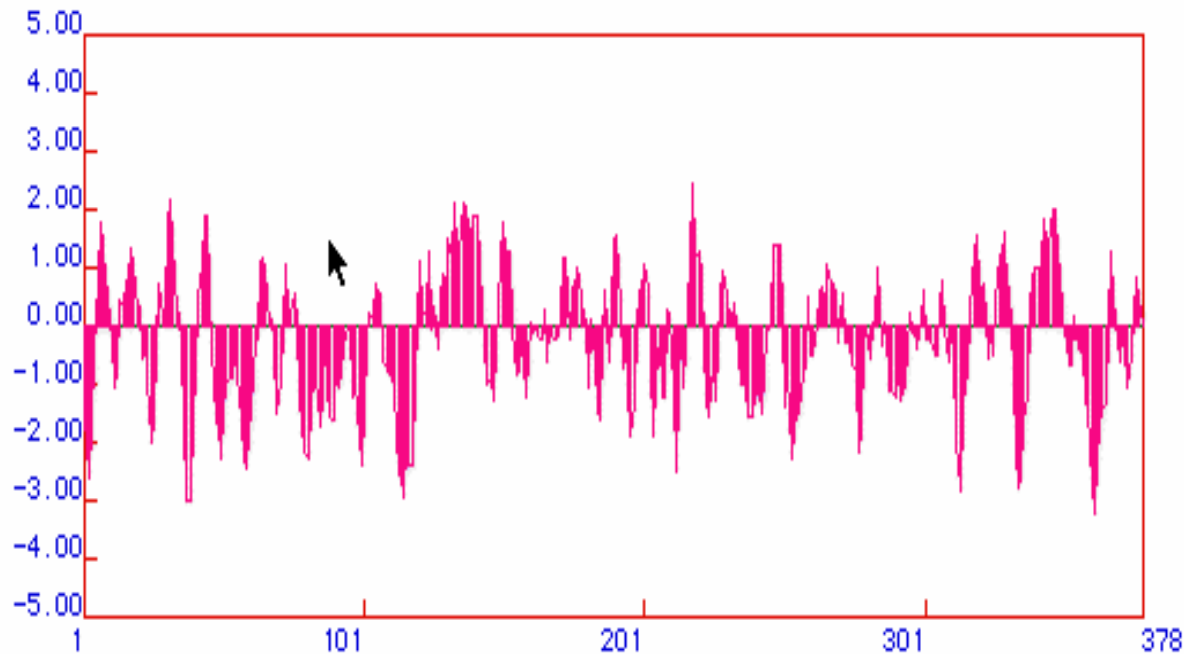
Method based in the comparison of the  $C\alpha$  of each residue in the Secondary Structure Elements (SSEs)

The procedure is design to find out these SSEs which could form similar spatial arrangements but with different topological connections

Manose represented by the SARF2 software. Pectate, lyase and agglutinin

<http://123d.ncifcrf.gov/sarfex.html>

# Part I: Structural Bioinformatics



Hydrophobicity plot for the human actin in which peaks above 2.00 Suggest hydrophobic chains

Pattern of hydrophobicity as approximation to predict transmembrane  $\alpha$  helix of proteins

# Part I: Structural Bioinformatics



## Protein 2D structure

- GOR
- Chou-Fasman
- Lim's
- Neural Network
- SVMs approximations

The ability also depends on predicting types of SSEs and defining classes of protein structures and patterns

- PHD (Profile Network from Heidelberg) for  $\alpha$  helices
- DSSP (Dictionary of Secondary Structure of Proteins)
- STRIDE (STRuctural IDentification)

# Part I: Structural Bioinformatics



	SEQUENCE ALIGNMENT	STRUCTURAL COMPARISONS
HOW TO	Sequences of proteins written one above the other so the similar amino acids are placed in the same columns and gaps are included	Proteins domains are superimposed fitting together the atoms as closely as possible so that the average deviation between them is the minimum
EVOLUTIONARY SIGNIFICANCE	Sequence similarity = evolutionary relationship	When structural similarity is common evolutionary relationship and convergence phenomena. When no common similarities then divergence phenomena but possible temporary folds

# Part I: Structural Bioinformatics



## 3D homology structure

There are available more than 356 000 known protein sequences but just 56 000 known structures

- New sequence has an homolog with about the same structure
- No homologues do exist and new structures also must be predicted
  - If two proteins share significant sequence similarity they should have also similar 3D structure
  - When the global alignment is performed and the identity shared between the proteins is 25-45 % then the two structures are likely to be similar
  - When approximately 45% , then the amino acids could be superimposed in the 3D structure

Some methods like

- SVMs (when remote homology search)
- PSI-BLAST (Position specific iterative BLAST)
- FPS (Family Pairwise Search)

# Part I: Structural Bioinformatics



## Threading

How well a sequence fits to a given 3D structure

Sequence comparisons can be made on structural level by computing the sequences-to-structure-fitness

1. The target sequence is threaded through the backbone structures of a collection of template proteins
2. Fold library or dictionary of resolved structures for sequence-to-structure alignment
3. "Godness of fit" score calculated in terms of empirical energy function based on statistics derived from known protein structures

Share some of the characteristics of both comparative modelling methods (the sequence alignment aspect) and *ab initio* prediction methods

# Part I: Structural Bioinformatics



*Ab initio* and Molecular Dynamics : Insights into protein folding and stability

*Ab initio*:

Method using only the amino acid sequence to find the 3D structure

Applicable to proteins with novel structure so that threading methods would fail

Rosetta: as the most important *ab initio* method

Protein function details and docking behavior are often analyzed based on force fields

Molecular Dynamics:

Physical laws can explain folds in a protein

Basic form of computer simulation where atoms and molecules are allowed to interact for a period of time under known laws of physics

# Genome Analysis

## Motivation

High-throughput techniques developed

Major source of information about the processes performed within a cell and evolved to one of the major topics in Bioinformatics

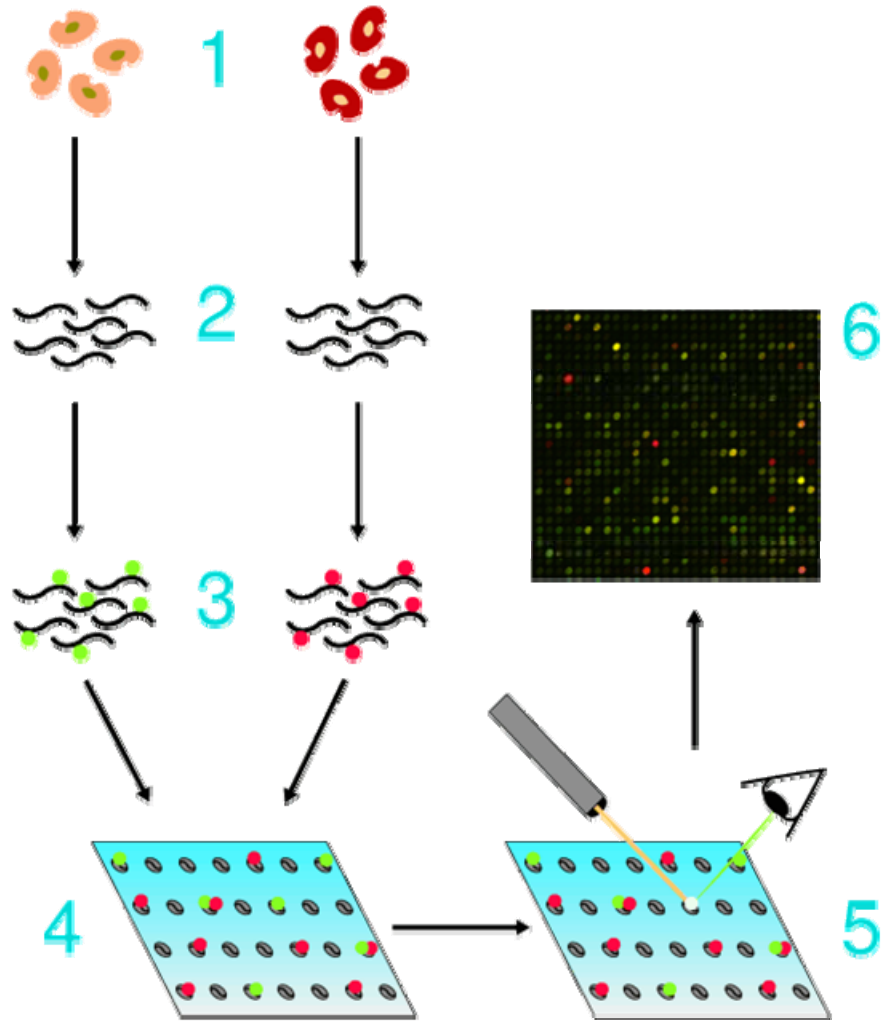
Provide means of measuring tens of thousands of genes simultaneously by measure at once cellular concentrations of thousands of mRNA: gene expression profile

Detection of genes that are differentially expressed (DEGs) in tissue samples

Basis for the functional genome analysis, molecular diagnostics, systems biology

Important applications in pharmaceutical and clinical research

# Part II: Genome Analysis



mRNA concentration ~ activity of a gene

Activity of a gene = expression level

The proportionality between the measured intensities and the number of copies of mRNA in the cell can vary in different arrays

# Part II: Genome Analysis



## 1. DNA Microarray

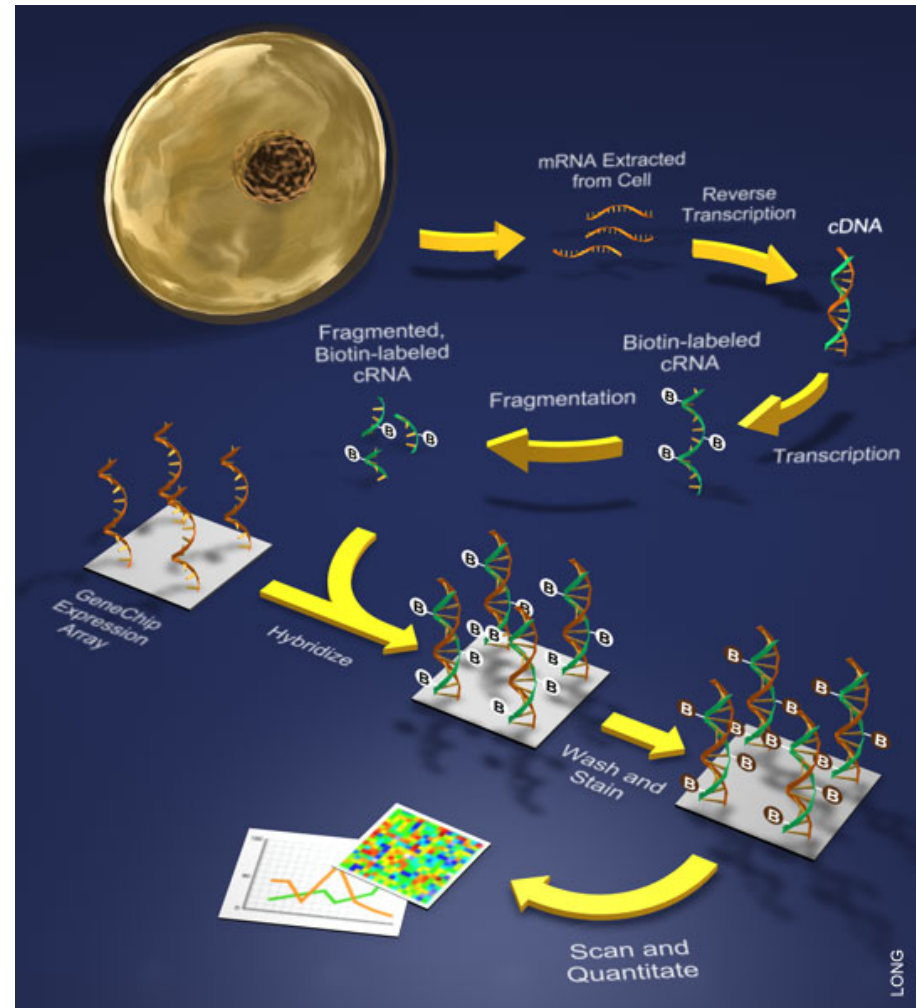
Techniques and Image analysis  
Background correction  
Normalization  
PM correction  
Summarization  
ML applications (Gene selection, clustering,...)

## 2. DNA analysis

Genome anatomy  
Genome individuality  
SNPs

## 3. Alternative splicing

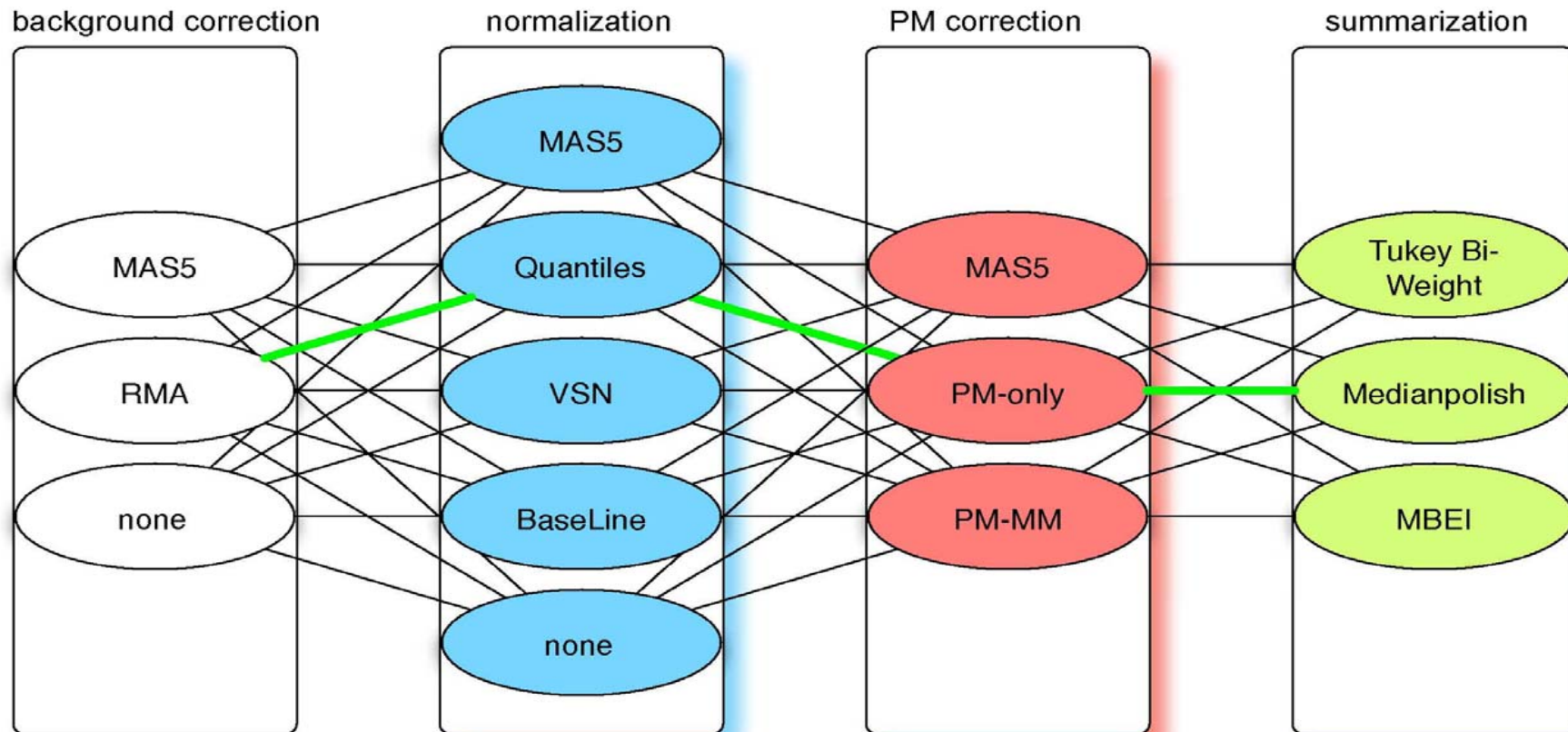
## 4. Modelling



# Part II: Genome Analysis



Many possible combinations for Microarray Summarization



## Part II: Genome Analysis



### 5. Next generation sequencing techniques:

Research community of genomics and transcriptomics as an alternative to array based methods: Illumina's Solexa, Roche's 454, or Applied Biosystems' SOLiD

→ Produces more than 50 million reads each 30 – 72 long prefix or suffix sequences of DNA fragments with length 100 to 500 base pairs

→ Reads Back-mapping to the reference genome (parallelized on multiprocessor machines or run on computer grids )

→ Analysis: to assemble a genome, to determine the transcripts and their concentrations, to detect nucleosome positions, to identify single nucleotide polymorphisms, or to estimate copy number variations