

In the Name of God, the Merciful, the Compassionate

Introduction to Bioinformatics

12 - Protein Structure Visualization, Comparison, and Classification

Instructor: Hossein Zeinali

Amirkabir University of Technology



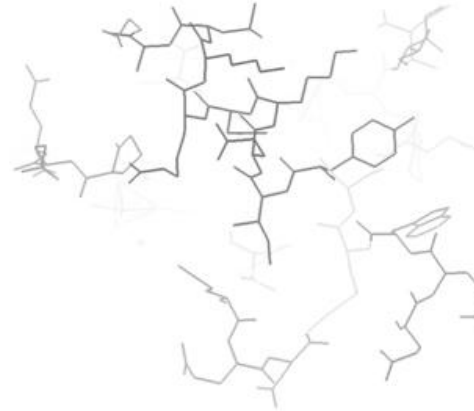
Protein Structure Visualization

- RASMOL & decendents: PyMol, MolMol
 - <http://www.umass.edu/microbio/rasmol/index2.htm>
- Cn3D - esp. good for structural alignments
 - <http://www.biosino.org/mirror/www.ncbi.nlm.nih.gov/Structure/cn3d/>
- CHIME (Protein Explorer)
 - <http://www.umass.edu/microbio/chime/getchime.htm>
- MolviZ.Org
 - <http://www.umass.edu/microbio/chime>
- Deep View = Swiss-PDB Viewer
 - <http://www.expasy.org/spdbv>

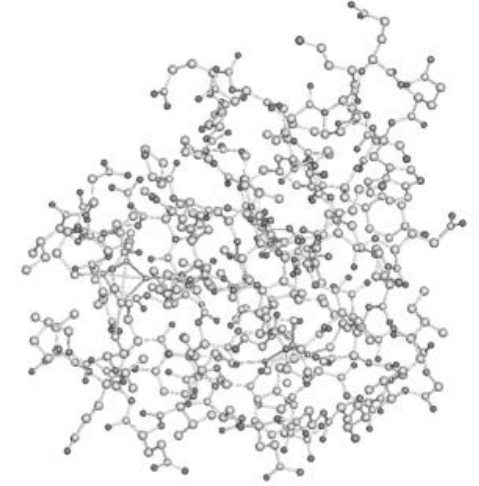
Molecular Structure Visualization Forms

- **(A) Wireframes**
 - A line drawing representing bonds between atoms
- **(B) Balls and sticks**
 - Solid spheres and rods, representing atoms and bonds
- **(C) Space-filling spheres**
 - each atom is described using large solid spheres
- **(D) Ribbons**
 - use cylinders or spiral ribbons to represent α -*helices* and broad, flat arrows to represent β -*strands*

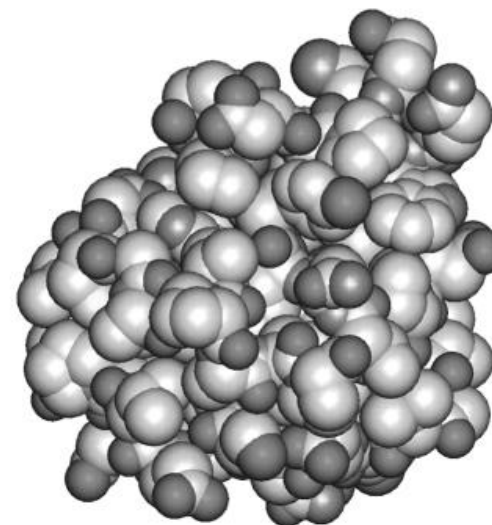
A



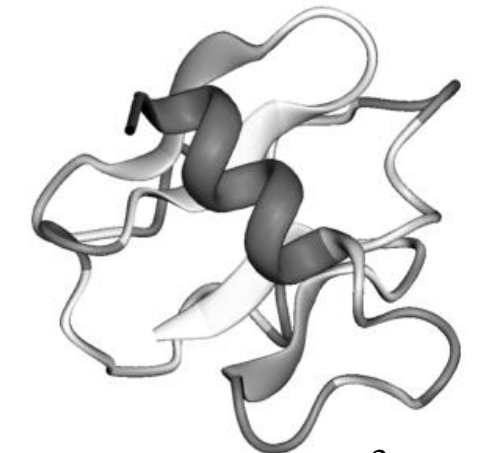
B



C



D



PyMol

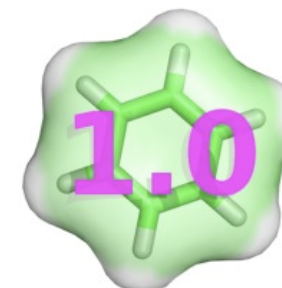
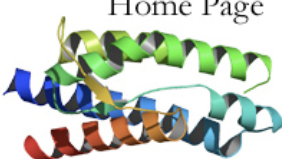
PyMOL Home Page

http://pymol.sourceforge.net/

GDCB Webmail ISU BCB Dobbs Lab BCB 444/544 Schedule of Classes Entrez PSB 2008 Conferences Seminars FadiWiki PPIDB Google Image


Download Documentation Mailing List Community Wiki F.A.Q. Plugins Links Project

PyMOL
Home Page

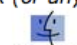


PyMOL v1.0 Released!


PyMOL Runs On:



Linux (or any Unix)



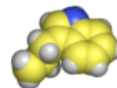
Mac OS
and



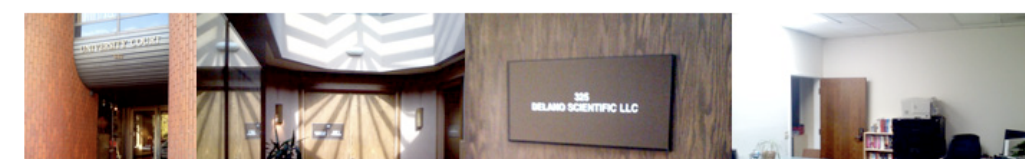
Windows

Click to Sponsor

PyMOL is a USER-SPONSORED molecular visualization system on an OPEN-SOURCE foundation. Please support our mission to create open, effective, and affordable tools for research by **purchasing a subscription** to maintenance and/or support. Thanks!

 **Latest News**

- October 4th, 2007: DeLano Scientific has moved!** We have relocated from temporary space into our new office at 540 University Ave., Suite 325, Palo Alto, CA 94301. Although this should be a long-term home, for certain delivery, please continue to send mail to P.O. Box 1118, Palo Alto, CA, 94302-1118, USA.



Cn3D

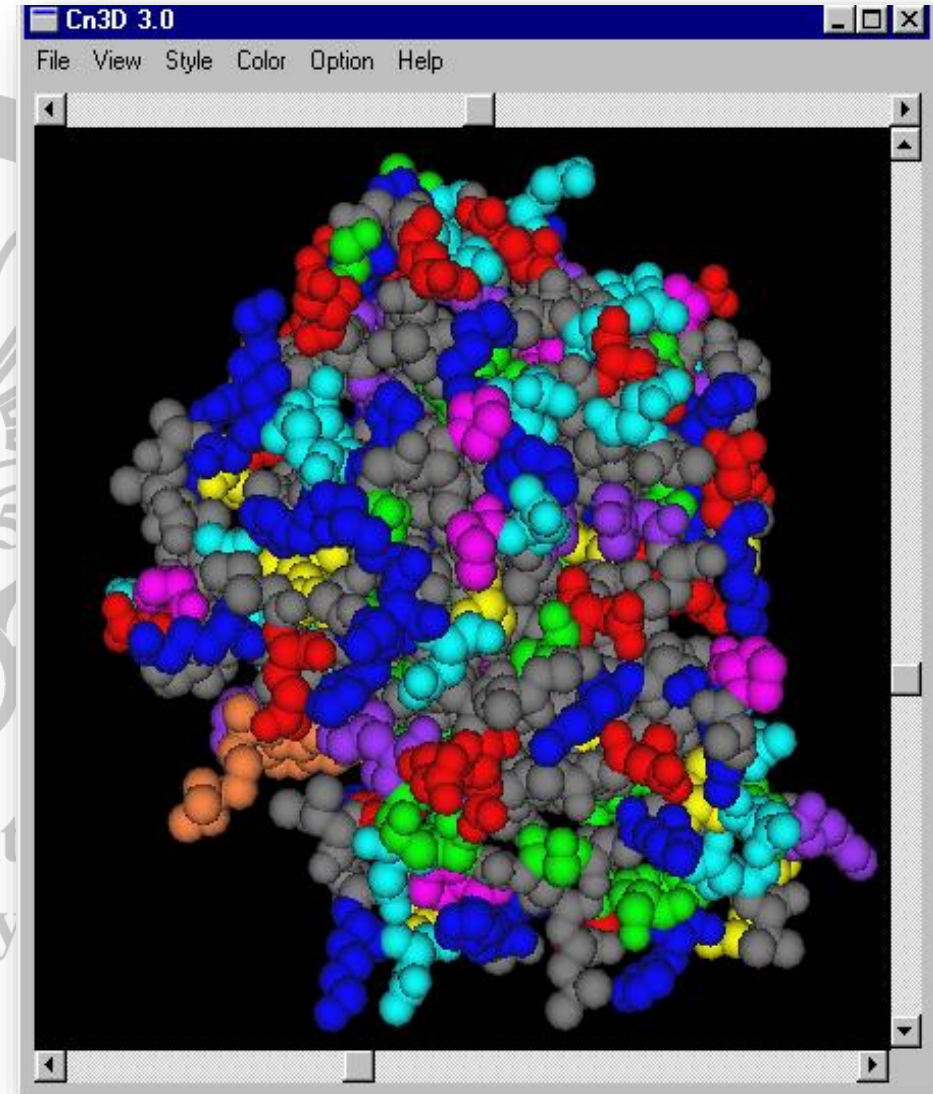
The screenshot shows a Netscape browser window titled "Cn3D Home Page" with the URL <http://www.ncbi.nlm.nih.gov/Structure/CN3D/cn3d.shtml>. The page features the NCBI logo and a navigation bar with links to PubMed, Entrez, BLAST, OMIM, Books, TaxBrowser, and Entrez Structure. A search bar is present with the text "Search Entrez" and a dropdown menu set to "Structure".

On the left side, there is a vertical menu with the following links: **Cn3D 4.1 Homepage**, **Cn3D Tutorial** (with a sub-link "Cn3D feature highlights"), **Cn3D FAQ** (with a sub-link "Frequently Asked Questions"), **Cn3D Install** (with a sub-link "Installation and Configuration"), and **MMDB** (with a sub-link "NCBI's structure database").

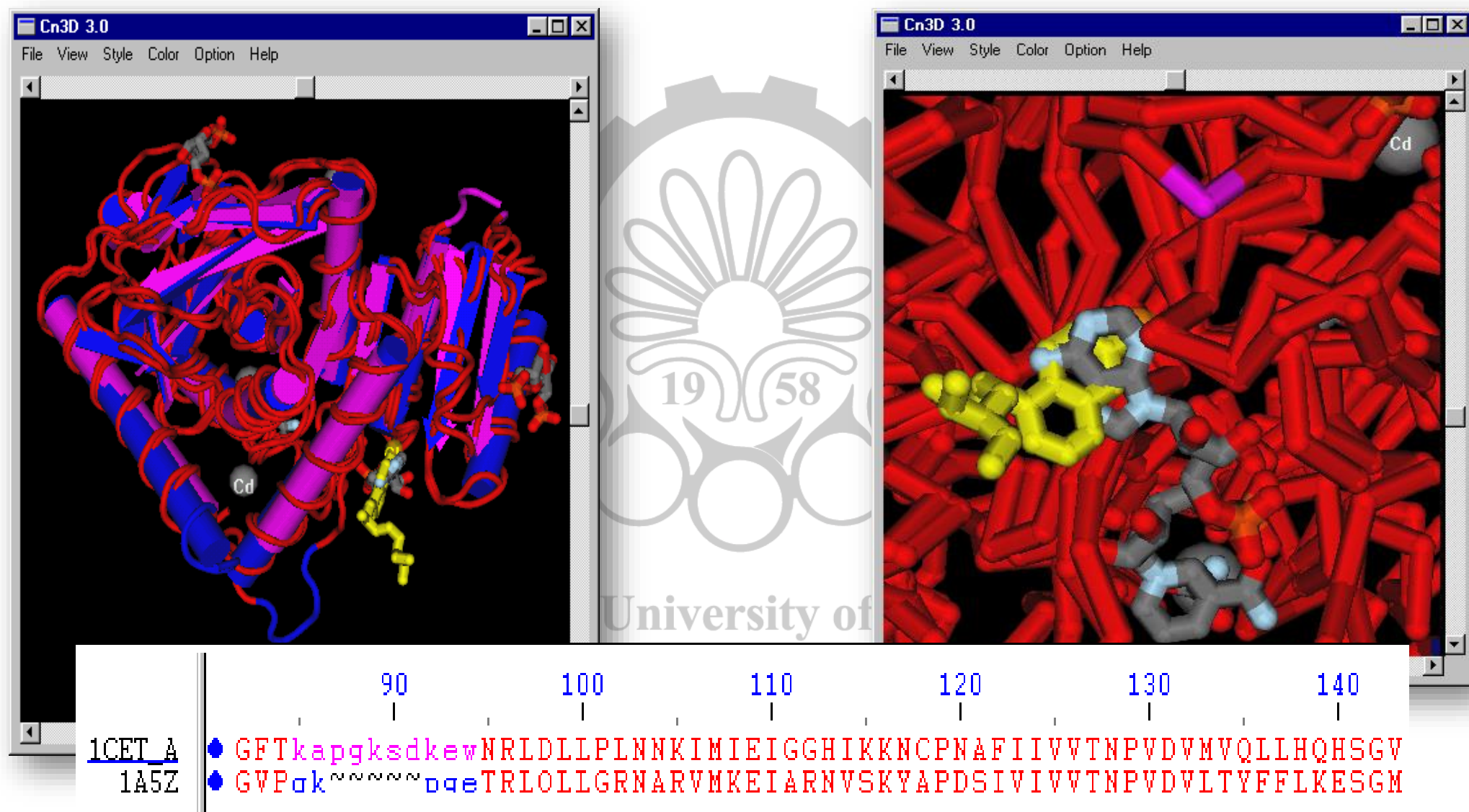
The main content area includes a link to **Download Cn3D 4.1 for PC, Mac and Unix**. Below this, a paragraph describes Cn3D as a helper application for viewing 3-dimensional structures from NCBI's Entrez retrieval service, noting its compatibility with Windows, Macintosh, and Unix, and its features for displaying structure, sequence, and alignment, as well as annotation and alignment editing.

A sample of Cn3D's capabilities is shown in a window titled "WD40 - Cn3D 4.1". This window displays a 3D ribbon diagram of a protein structure. A panel on the right, titled "CDD Descriptive Items", shows the name "WD40" and a description: "WD40 domain, found in a number of e a wide variety of functions including ac in signal transduction, pre-mRNA proc assembly; typically contains a GH dipe its N-terminus and the WD dipeptide a".

Cn3D : Displaying 3' Structures



Cn3D: Structural Alignments




Protein Explorer (Chime)

Protein Explorer FrontDoor

http://www.umass.edu/microbio/chime/pe_beta/pe/protexpl/frntdoor.htm

GDCB Webmail ISU BCB Dobbs Lab BCB 444/544 Schedule of Classes Entrez PSB 2008 Conferences Seminars FadiWiki PPIDB

 **proteinexplorer.org**
FrontDoor to Protein Explorer 2.80
from the University of Massachusetts -- also available from the [San Diego mirror](#)
Check Out MolSlides!
[Compatible browsers: Firefox \(recommended\), Internet Explorer, etc.](#)

Award: 2003 MERLOT CLASSICS
"Protein Explorer has revolutionized the teaching of biology at a molecular level."

Copyright © 2007 by Eric Martz
Thanks to [Timothy Driscoll](#) who implemented the [command script recorder](#).
Bookmark this page. Bookmarking subsequent pages in this site may not work!

[Protein Explorer en español](#)
Versión 2.25, compatible con Microsoft Internet Explorer.

[PE en Français](#)

Mac PPC OSX users: Please see IMPORTANT PPC OSX HELP.

Beginners start here:
[What Is Protein Explorer?](#)

- [Watch the Protein Explorer Demo.](#)**
Relax while an interactive Flash movie shows you how to get the most out of Protein Explorer.
- [Quick-Start Protein Explorer](#)**
to explore DNA complexed with a protein (yeast Gal4 transcriptional regulator, [PDB identification code](#) 1d66).
- [Do the 1-Hour Tour](#)**
for an introduction to Protein Explorer.

Can't get PE to work? Here is [Troubleshooting help!](#)

Experienced users:
Enter [PDB code](#):

[Start Empty Explorer](#) to load a [saved](#) PDB file from your disk.

[What's new?](#)
[MolSlides!](#)

[Help/Index/Glossary](#)

[About Protein Structure.](#)
[Short Courses.](#)

Protein Structure Comparison

- The comparative analysis involves the *direct alignment* and *superimposition of structures* in a three-dimensional space.
- Structure comparison is one of the fundamental techniques in protein structure analysis.
 - The comparative approach is important in finding remote protein homologs.
- Proteins can share common structures even without sequence similarity
 - Structures have a much higher degree of conservation than the sequences
- Structure comparison can often reveal distant evolutionary relationships between proteins
- Protein structure comparison is a prerequisite for protein structural classification

(Tehran Polytechnic)

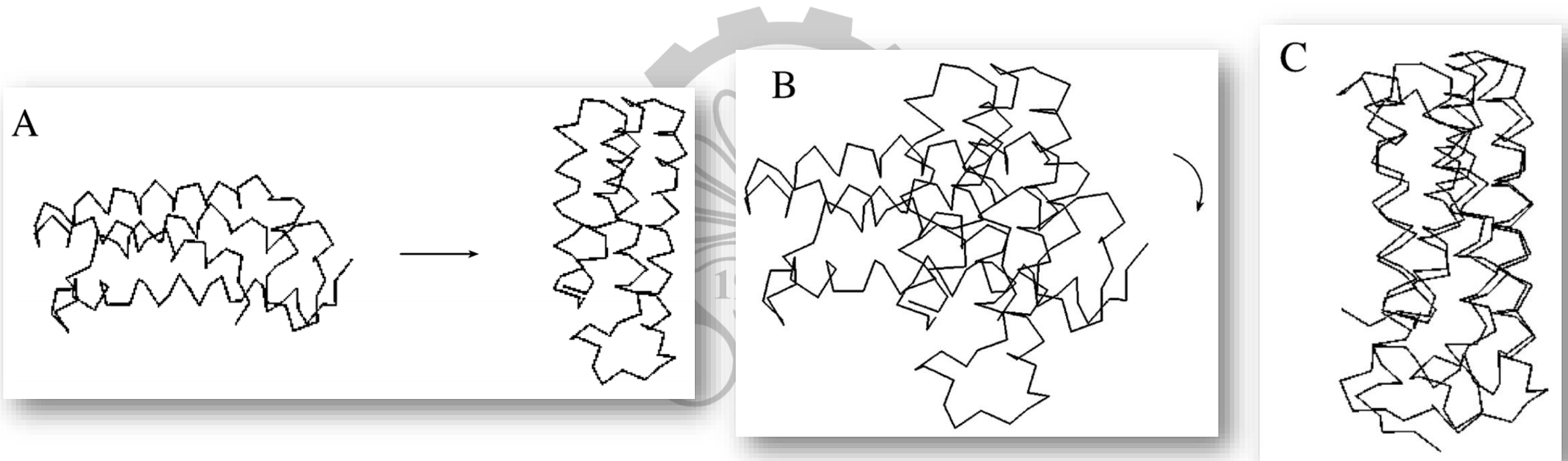
Protein Structure Comparison Methods

- Approaches to comparing protein geometric properties can be divided into three categories
 - Minimizing intermolecular distances
 - Measuring intramolecular distances
 - Combine both intermolecular and intramolecular
- DALI/FSSP (most commonly used)
 - Fully automated structure alignments using intramolecular distance method
 - DALI server: <http://www.ebi.ac.uk/dali/index.html>
 - DALI Database (fold classification):
<http://ekhidna.biocenter.helsinki.fi/dali/start>

Intermolecular Method

- Normally applied to relatively similar structures
- One of the structures has to be moved with respect to the other in such a way that the two structures have a maximum overlap.
- Procedure:
 - Identifying equivalent residues or atoms by sequence based alignment
 - *Translation* : one of the structures is moved laterally and vertically toward the other structure
 - The structures are further rotated relative to each other around the three-dimensional axes
 - The distances between equivalent positions are constantly measured
 - The rotation continues until the shortest intermolecular distance is reached: *optimal superimposition*
 - Equivalent residue pairs can be identified

Intermolecular Method



Amirkabir University of Technology
(Tehran Polytechnic)

Root Mean Square Deviation (RMSD)

- An important measurement of the structure fit during superposition is the distance between equivalent positions on the protein structures.
- *Root mean square deviation:*

$$RMSD = \sqrt{\frac{1}{N} \sum_{i=1}^N D_i^2}$$

where D is the *distance between coordinate data points* and N is the total number of corresponding residue pairs.

- In practice, only the distances between C α carbons of corresponding residues are measured.

Amirkabir University of Technology
(Tehran Polytechnic)

Root Mean Square Deviation (RMSD)

- The goal of structural comparison is to achieve a minimum RMSD:
 - The problem with RMSD is that it depends on the size of the proteins being compared
 - For the same degree of sequence identity, large proteins tend to have higher RMSD values
- Correct this size-dependency by logarithmic factor:

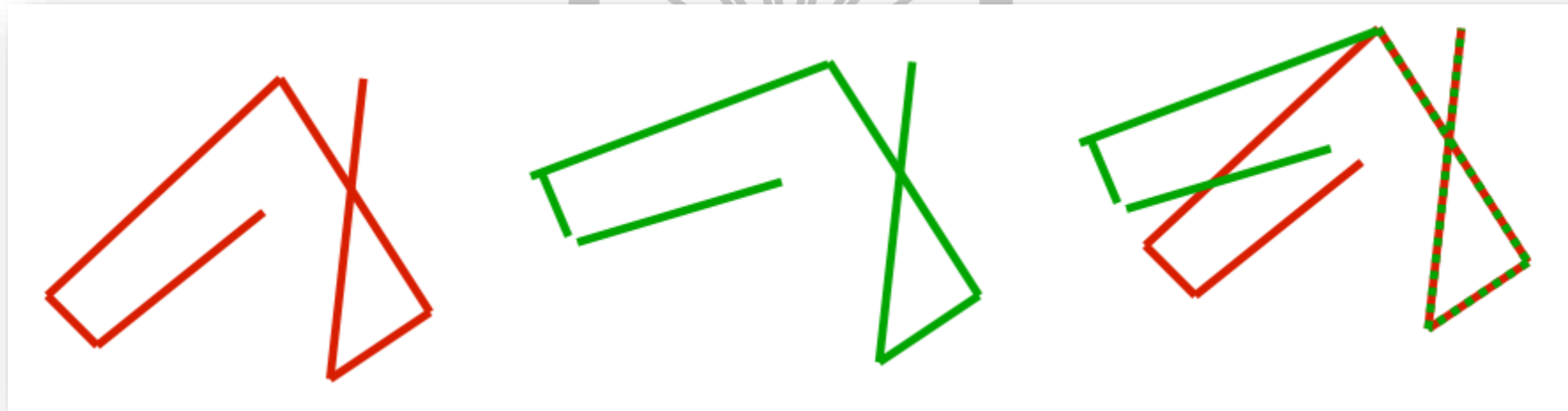
$$RMSD_{100} = \frac{RMSD}{-1.3 + 0.5 \ln(N)}$$

Intermolecular Method (Cont.)

- The most challenging part of using the intermolecular method is to identify equivalent residues in the first place.
 - Depends on sequence alignment methods
- Solutions to compare more distantly related structures:
 - Delete sequence variable regions outside secondary structure elements to reduce the search time
 - Divide the proteins into small, match similar regions is then done fragment by fragment. Finally, a joint superposition for the entire structure is performed.
 - Using *iterative optimization*

Problems with RMSD

- A small local alignment error can propagate and the quality of alignment may be underestimated



Amirkabir University of Technology
(Tehran Polytechnic)

Intramolecular Method

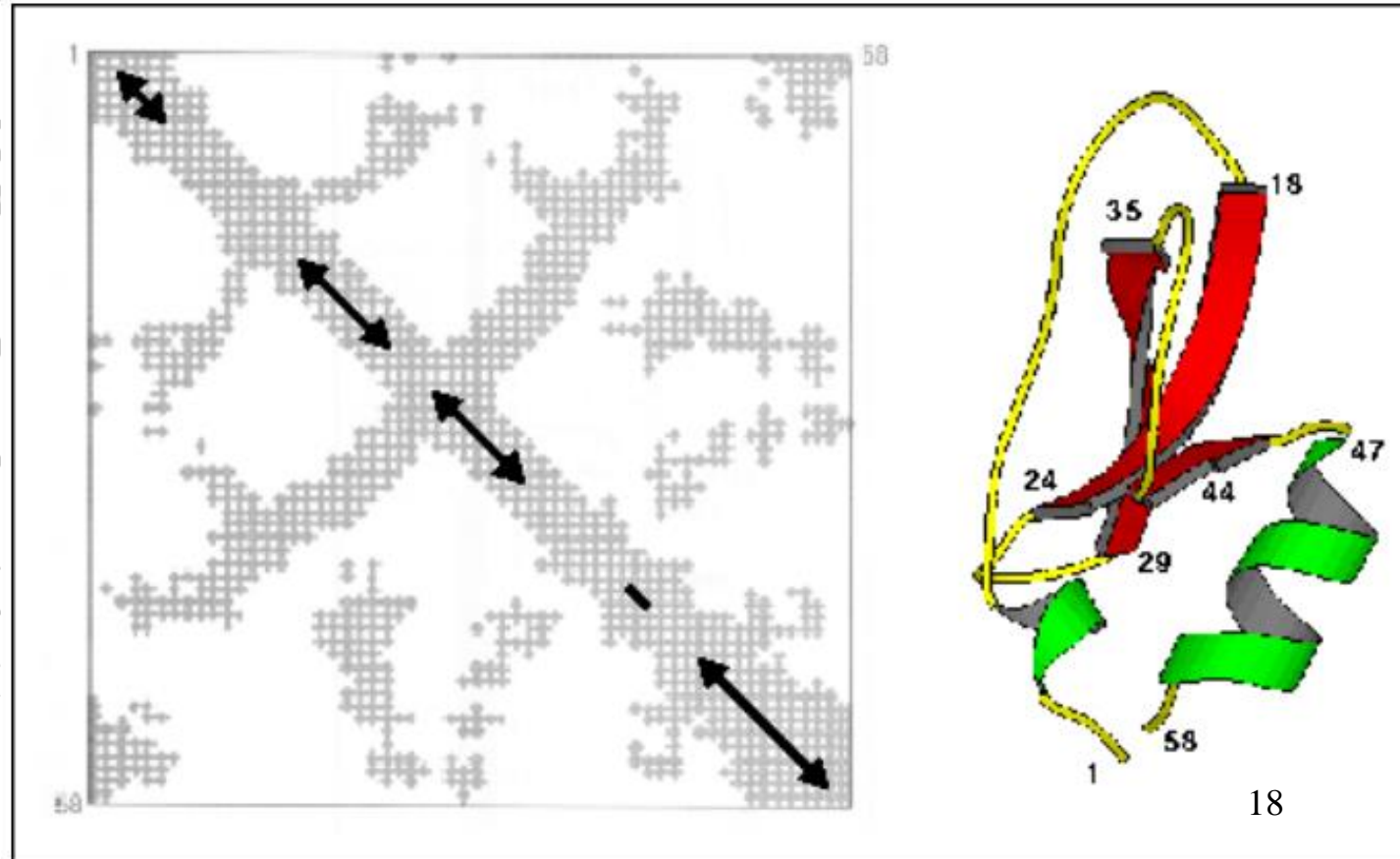
- Relies on structural internal statistics
 - Does not depend on sequence similarity between the proteins.
- Does not generate a physical superposition of structures
 - Provides a quantitative evaluation of the structural similarity between corresponding residue pairs.
- The method works by generating a distance matrix between residues of the same protein.
 - The distance matrices from the two structures are moved relative to each other to achieve maximum overlaps.
 - Similar intramolecular distance patterns representing similar structure folding regions can be identified.
- For the ease of comparison, each matrix is decomposed into smaller submatrices.

(Tehran Polytechnic)

Contact matrix: the Used Distance Matrix

- Contact matrix $n \times n$ matrix where $n = \text{\#residues}$
 $d(i, j) = \text{distance}(C\alpha_i, C\alpha_j)$
- Example: pairs with $d(i, j)$ below a certain threshold are gray and the rest is white
- **Idea:** Similar structures have similar contact matrices

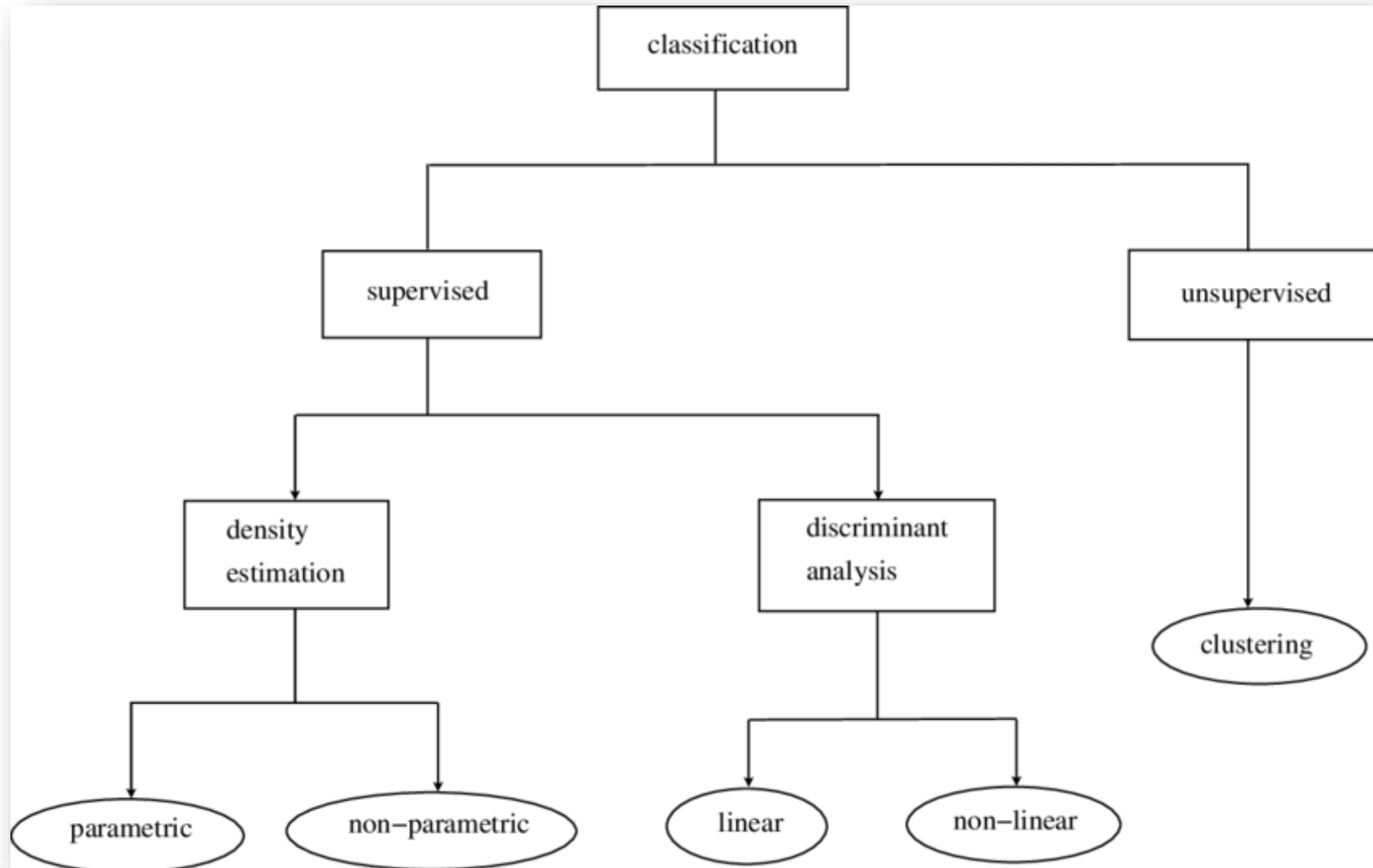
Amirkabir U
(Teh



Comparison Tools

- Combinatorial Extension (CE)
 - <http://cl.sdsc.edu/ce.html>
 - Uses the intramolecular distance
- Vector Alignment Search Tool (VAST)
 - www.ncbi.nlm.nih.gov/80/Structure/VAST/vast.shtml
 - Uses both the inter- and intramolecular approaches
- SSAP
 - www.biochem.ucl.ac.uk/cgi-bin/cath/GetSsapRasmol.pl
 - Uses an intramolecular distance-based method
- STAMP
 - www.compbio.dundee.ac.uk/Software/Stamp/stamp.html
 - Uses the intermolecular approach by *iterative alignment*

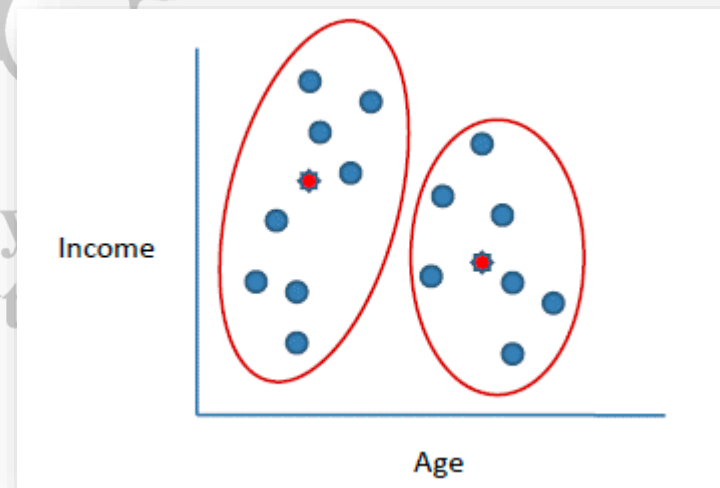
Classification Hierarchy



K-Means Clustering

1. Clusters the data into k groups where k is predefined.
2. Select k points at random as cluster centers.
3. Assign objects to their closest cluster center according to the *Euclidean distance* function.
4. Calculate the centroid or mean of all objects in each cluster.
5. Repeat steps 2, 3 and 4 until the same points are assigned to each cluster in consecutive rounds.

Amirkabir University
(Tehran Poly)



Protein Structure Classification

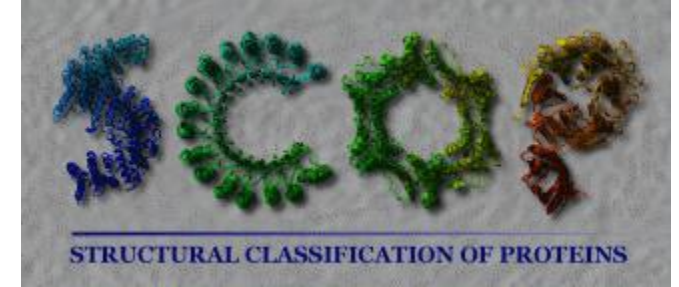
- One of the applications of protein structure comparison is structural classification.
- The ability to compare protein structures allows classification of the structure data and identification of relationships among structures.
- Usage: establish hierarchical relationships among protein structures and to provide a comprehensive and evolutionary view of known structures.
- Method & Databases:
 - Structural Classification of Proteins (SCOP)
 - Class, Architecture, Topology and Homologous (CATH)

Protein Structure Classification (Cont.)

- The first step in structure classification is to remove redundancy from databases.
 - The redundancy can be removed by selecting representatives through a sequence alignment–based approach.
- The second step is to separate structurally distinct domains within a structure.
 - Proteins with multiple domains must be subdivided before a sensible structural comparison can be carried out.
 - Can be done either manually or based on special algorithms for domain recognition.
 - Structure comparison can be conducted at the domain level, either through manual inspection, or automated, or a combination of both.
- The last step involves grouping proteins/domains of similar structures and clustering them based on different levels.

Structural Classification of Proteins (SCOP)

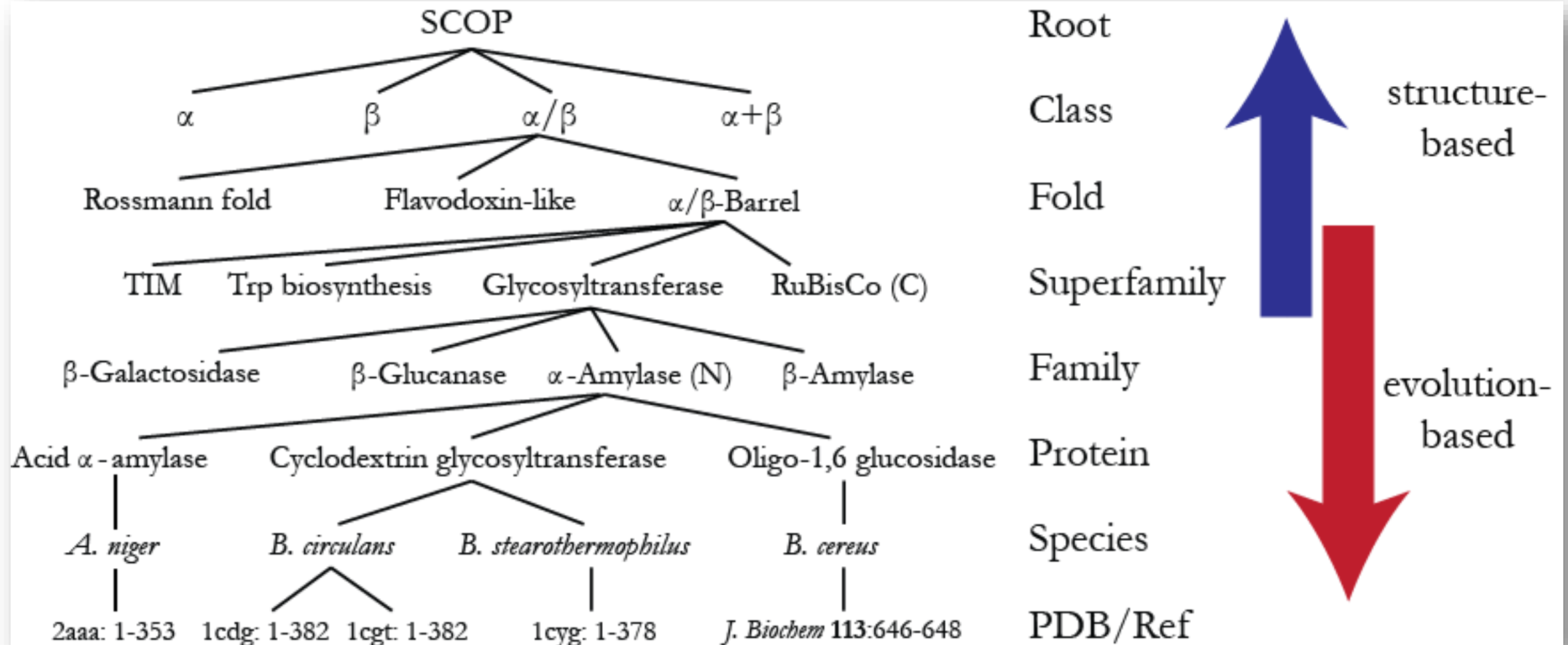
- SCOP is a database for comparing and classifying protein structures.
 - <http://scop.mrc-lmb.cam.ac.uk/scop/>
 - Is constructed based on manual examination of protein structures.
 - The proteins are grouped into hierarchies of classes, folds, super-families, and families.
 - Includes 44,769 non-redundant domains representing 540,282 protein structures.
- The SCOP families consist of proteins having high sequence identity (>30%).



Structural Classification of Proteins (SCOP)

- Super-families consist of families with similar structures, but weak sequence similarity.
- Folds consist of super-families with a common core structure, which is determined manually.
 - Fold's members do not always have evolutionary relationships.
- Classes consist of folds with similar core structures.
 - The highest level of the hierarchy
 - Distinguishes groups of proteins by secondary structure compositions

Structural Classification of Proteins (SCOP)



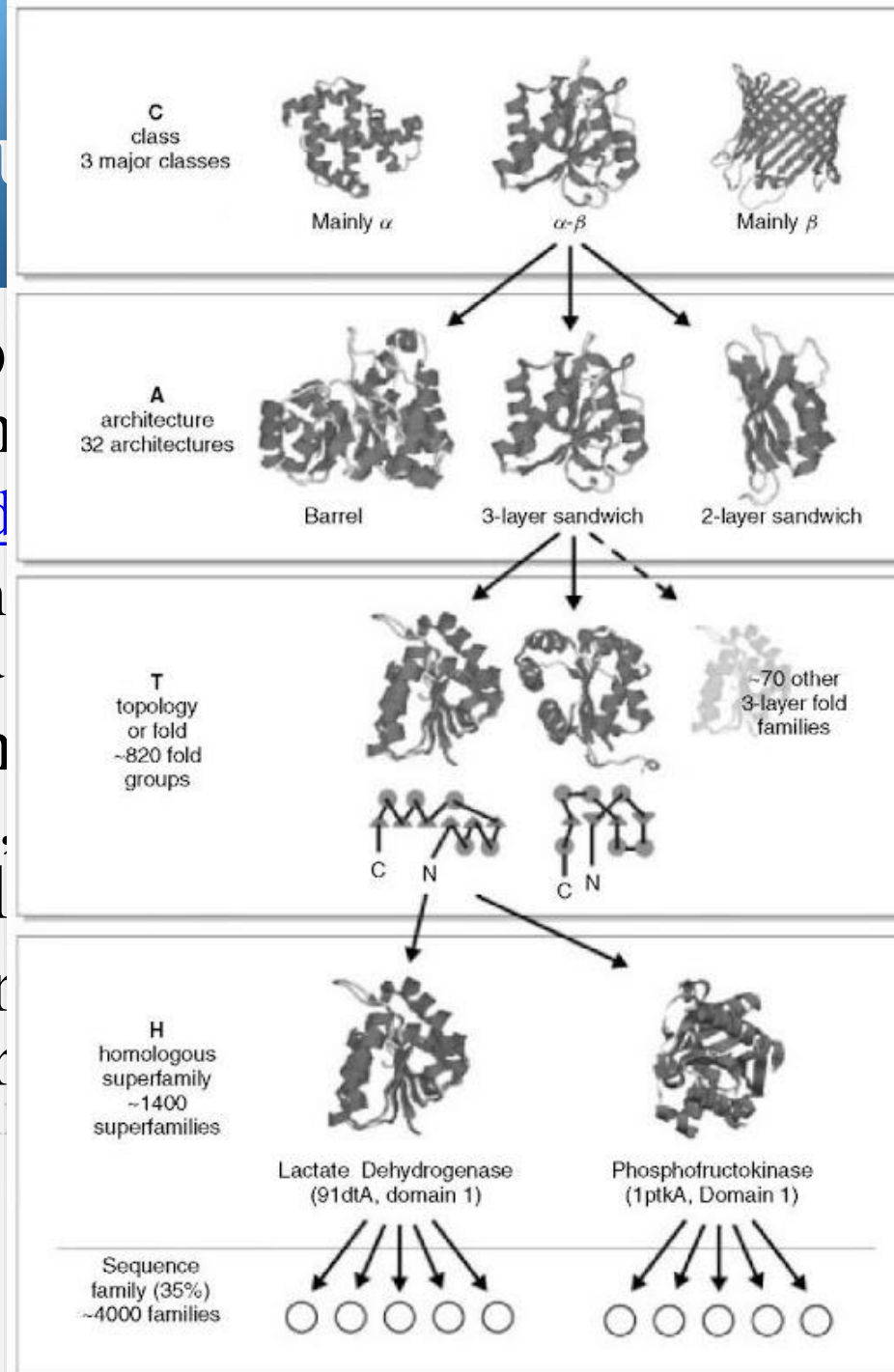
Class, Architecture, Topology and Homologous (CATH)

- CATH classifies proteins based on the automatic structural alignment program SSAP as well as manual comparison.
 - <https://www.cathdb.info/>
 - Structural domain separation is carried out as a combined effort of a human expert and computer programs.
- Individual domain structures are classified at five major levels: class, architecture, fold/topology, homologous superfamily, and homologous family.
- Architecture describes the overall packing and arrangement of secondary structures independent of connectivity between the elements.



Class, Architecture

- CATH classifies protein structural domains using a hierarchical alignment program
 - <https://www.cathdb.org/>
 - Structural domain classification is based on the combined effort of a human expert and a computer
- Individual domain classification is based on the major levels: superfamily, architecture, and topology
- Architecture describes the arrangement of secondary structural elements.



ologous (CATH)

structural comparison.



combined effort of a

major levels: superfamily, and

arrangement of y between the

Comparison of SCOP and CATH

- SCOP is almost entirely based on manual comparison of structures by human experts
 - CATH is a combination of manual curation and automated procedure
- The classification results from both systems are quite similar.
 - The results from the two systems converge at about 80% of the time.
 - Only about 20% of the structure fold assignments are different.

Amirkabir University of Technology
(Tehran Polytechnic)

References

- Mostly used:
 - Essential bioinformatics, Chapter 13 (Protein Structure Visualization, Comparison, and Classification)
- IP notice: some slides were selected from Drena Dobbs' slides.



Amirkabir University of Technology
(Tehran Polytechnic)

Thanks for your attention

