

CAP 5510: Introduction to Bioinformatics  
CGS 5166: Bioinformatics Tools

**Giri Narasimhan**

ECS 254; Phone: x3748

[giri@cis.fiu.edu](mailto:giri@cis.fiu.edu)

[www.cis.fiu.edu/~giri/teach/BioinfS15.html](http://www.cis.fiu.edu/~giri/teach/BioinfS15.html)

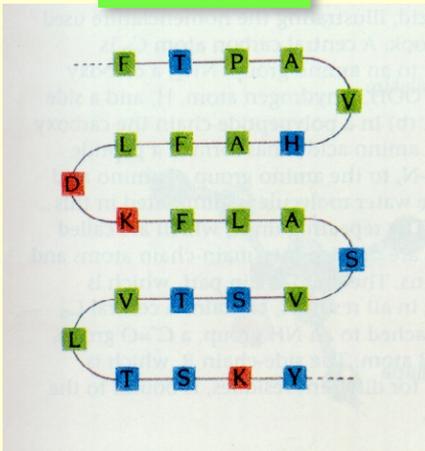
---

# Proteins and Protein Structure

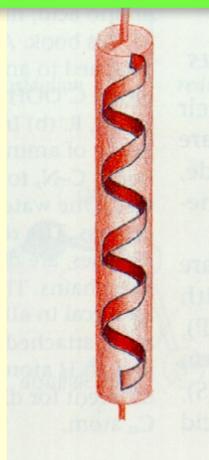
# Protein Structures

- Sequences of amino acid residues
- 20 different amino acids

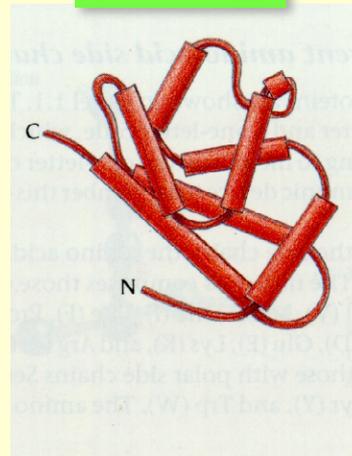
Primary



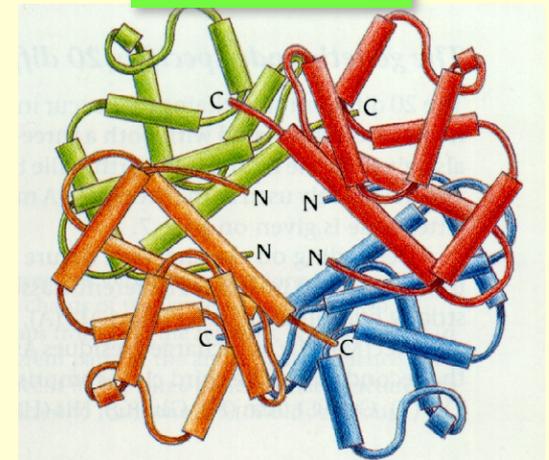
Secondary



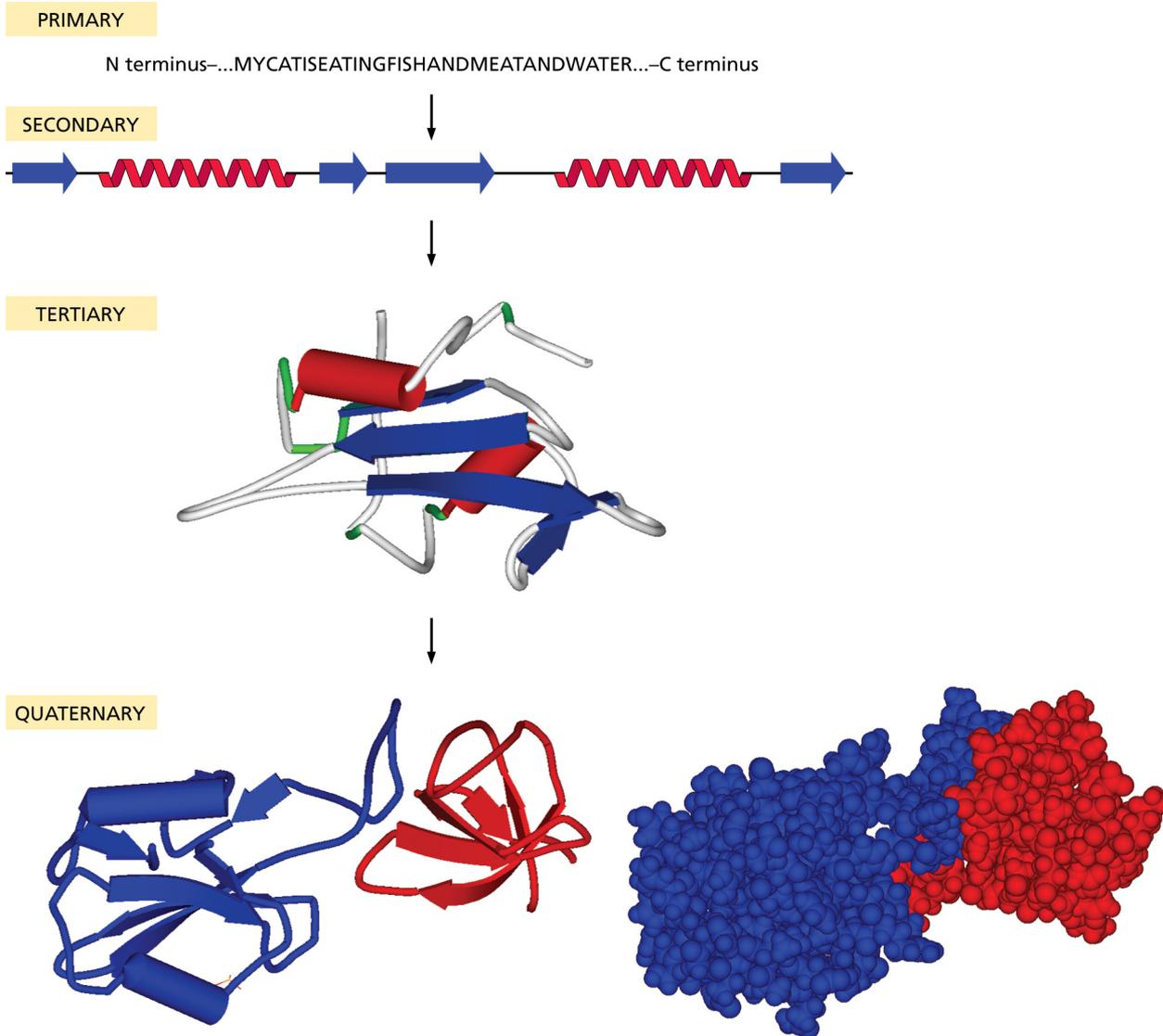
Tertiary



Quaternary



# Proteins: Levels of Description



# Proteins

- **Primary structure** is the sequence of amino acid residues of the protein, e.g.,

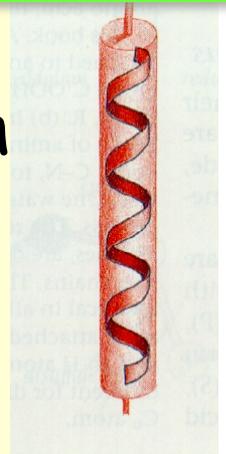
**Flavodoxin:**

AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADA...

- Different regions of the sequence form local regular **secondary structures**, such
  - **Alpha helix**, **beta strands**, etc.

AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADA...

Secondary



# More on Secondary Structures

## □ $\alpha$ -helix

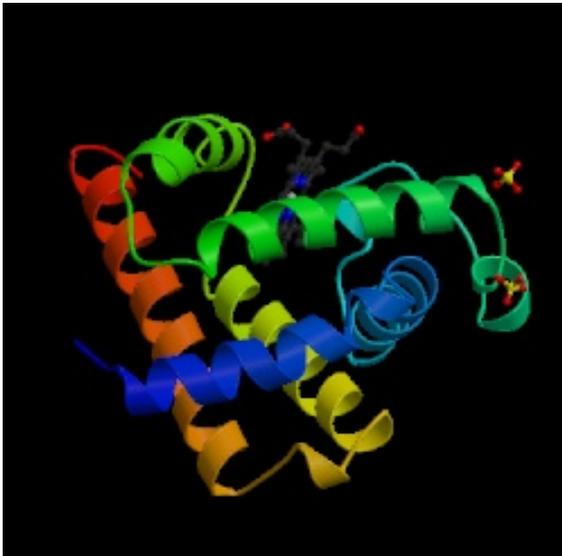
- Main chain with peptide bonds
- Side chains project outward from helix
- Stability provided by H-bonds between CO and NH groups of residues 4 locations away.

## □ $\beta$ -strand

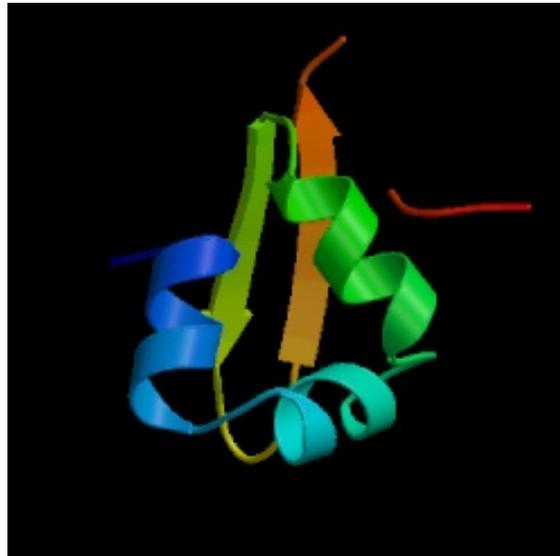
- Stability provided by H-bonds with one or more  $\beta$ -strands, forming  $\beta$ -sheets. Needs a  $\beta$ -turn.

# Proteins

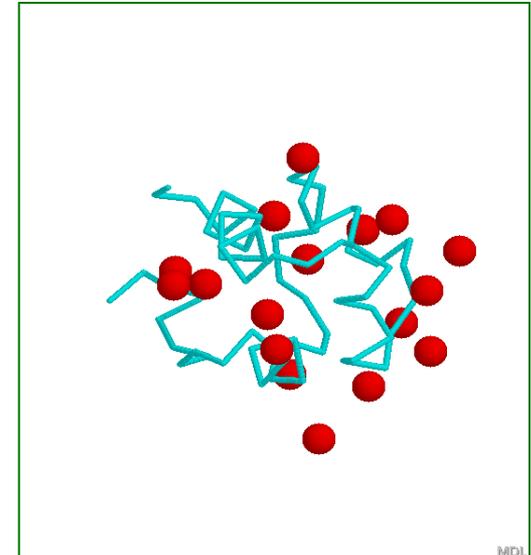
- **Tertiary structures** are formed by packing secondary structural elements into a globular structure.



Myoglobin

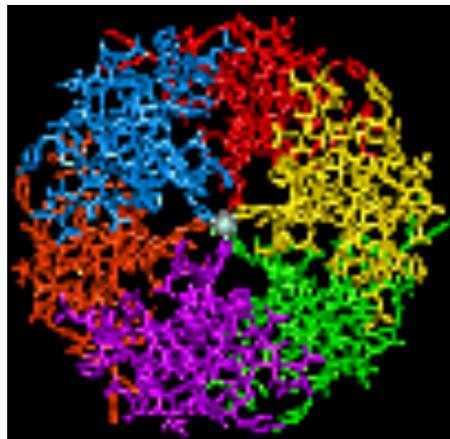
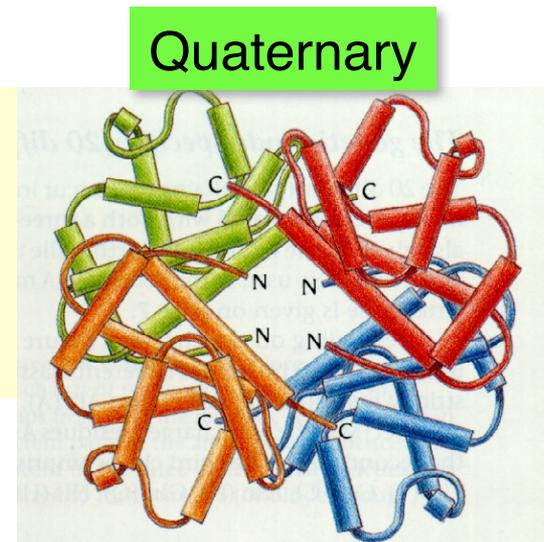


Lambda Cro



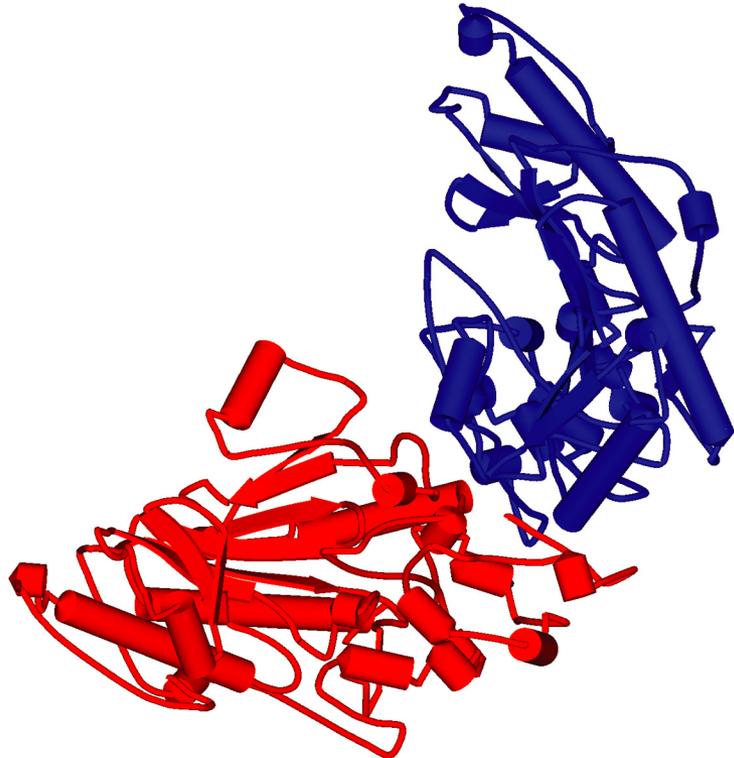
# Quaternary Structures in Proteins

- The final structure may contain more than one “chain” arranged in a **quaternary structure**.



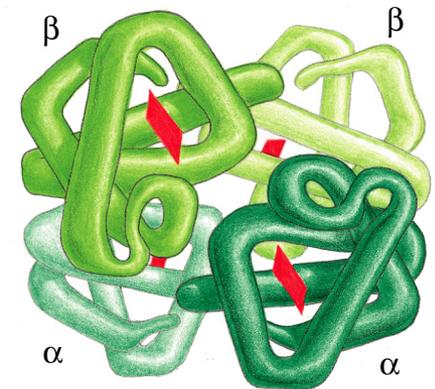
Insulin Hexamer

# More quaternary structures



Muscle creatine kinase  
(Homodimer)

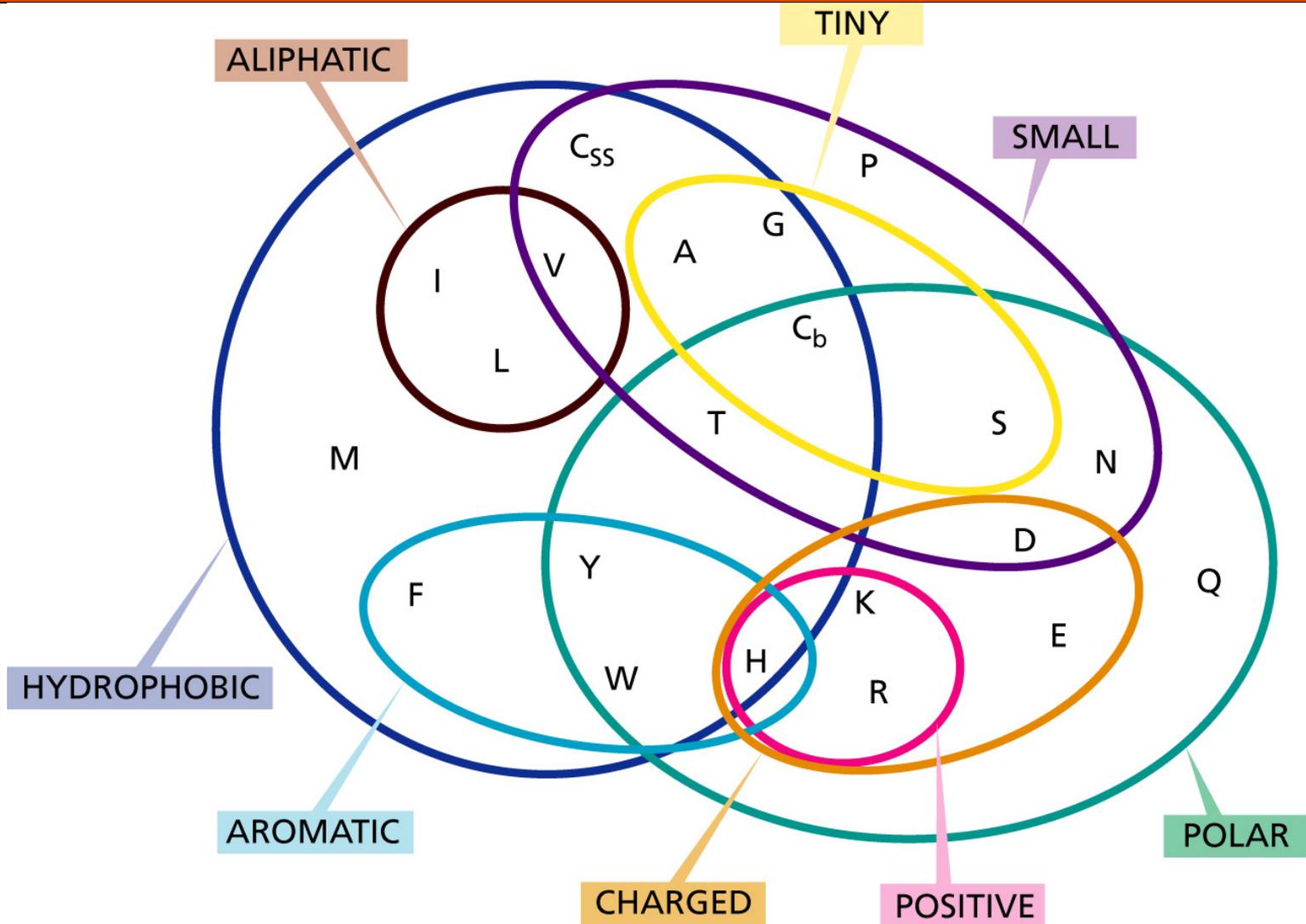
Bovine deoxyhemoglobin  
(Heterotetramer)



# Amino Acid Types

- Hydrophobic**    **I, L, M, V, A, F, P**
- Charged**
  - Basic**    **K, H, R**
  - Acidic**    **E, D**
- Polar**    **S, T, Y, H, C, N, Q, W**
- Small**    **A, S, T**
- Very Small**    **A, G**
- Aromatic**    **F, Y, W**

# Amino Acid Types



# Structure of a single amino acid

All 3 figures are cartoons of an amino acid residue.

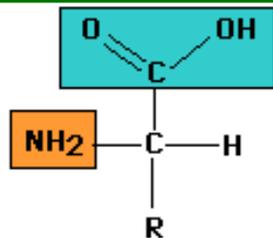
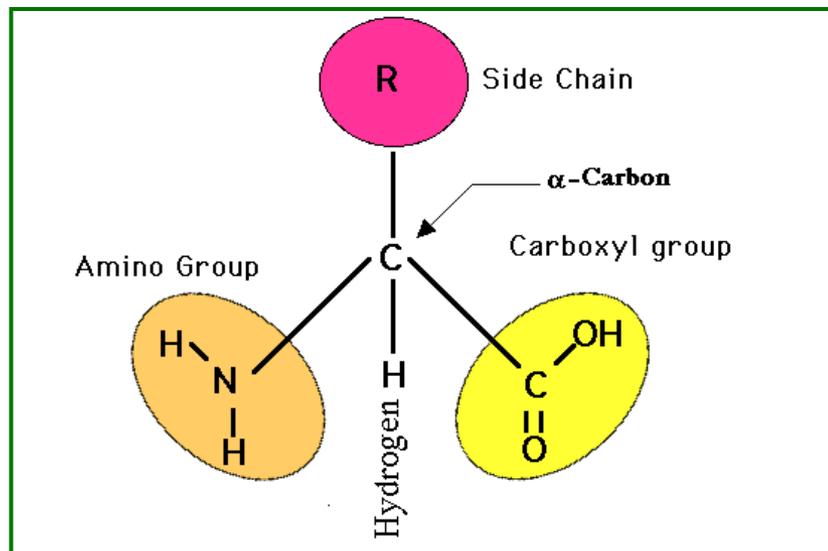
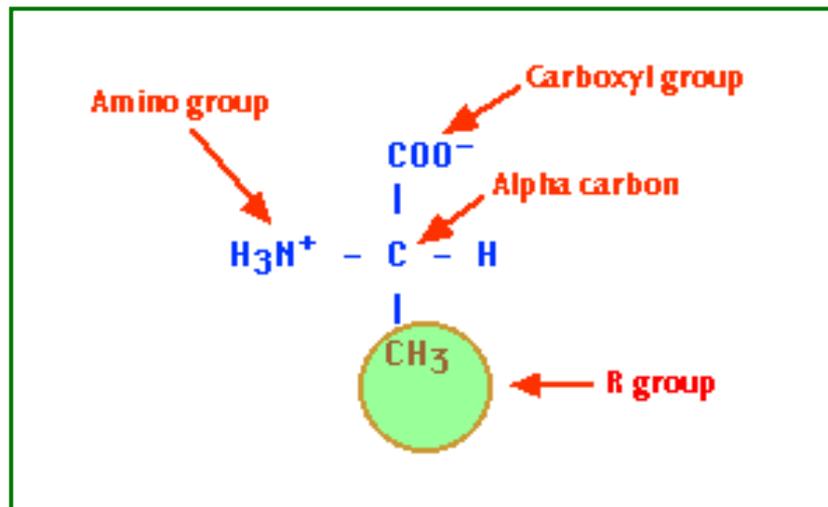
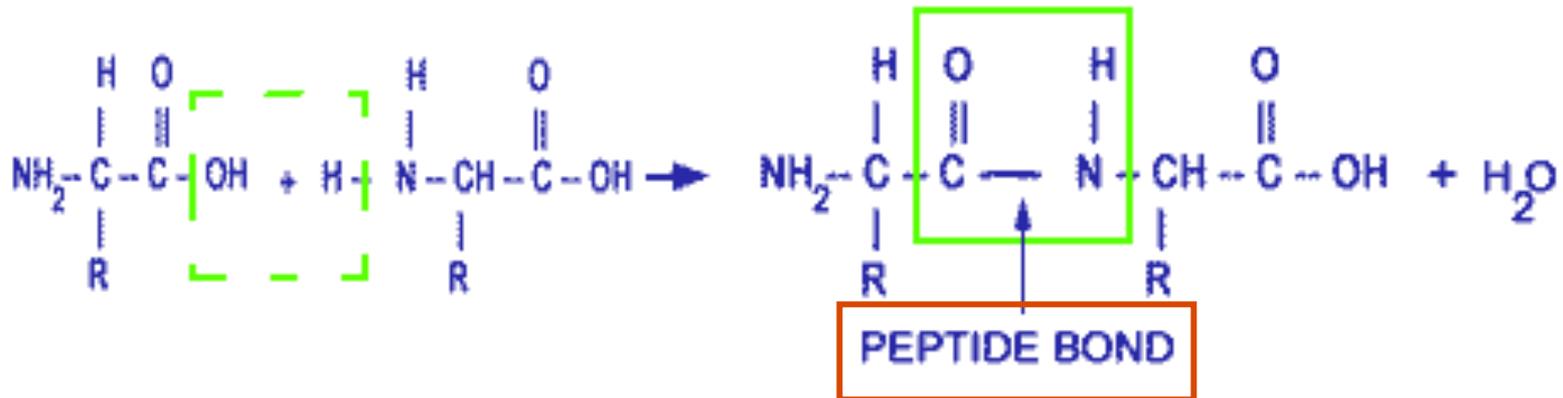


Fig. General formula for an amino acid molecule. "R" represents the variable groups that are attached to this basic molecule to make up the 20 common amino acids

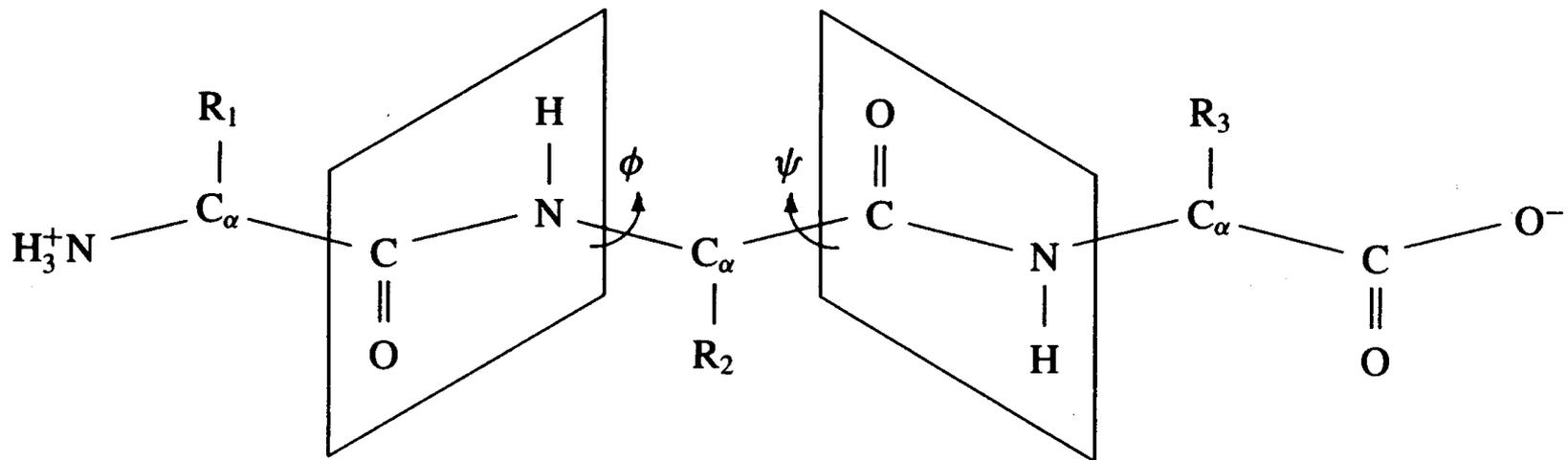


# Chains of amino acids



**Amino acids vs Amino acid residues**

# Angles $\phi$ and $\psi$ in the polypeptide chain



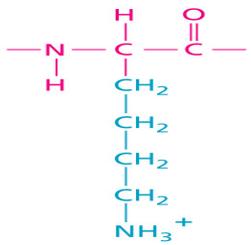
**FIGURE 1.2**

*A polypeptide chain. The  $\text{R}_i$  side chains identify the component amino acids. Atoms inside each quadrilateral are on the same plane, which can rotate according to angles  $\phi$  and  $\psi$ .*

## BASIC SIDE CHAINS

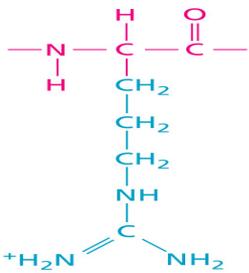
### lysine

(Lys, or K)



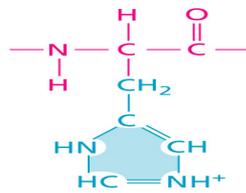
### arginine

(Arg, or R)



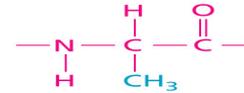
### histidine

(His, or H)



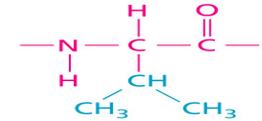
### alanine

(Ala, or A)



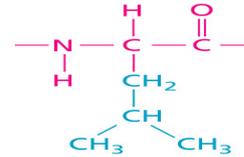
### valine

(Val, or V)



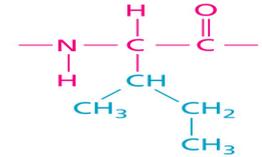
### leucine

(Leu, or L)



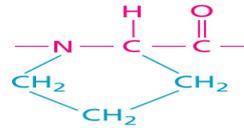
### isoleucine

(Ile, or I)



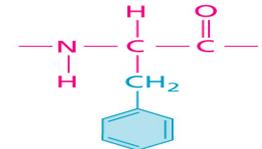
### proline

(Pro, or P)



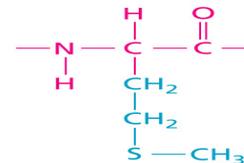
### phenylalanine

(Phe, or F)



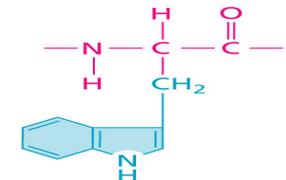
### methionine

(Met, or M)



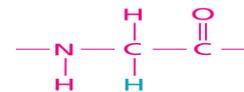
### tryptophan

(Trp, or W)



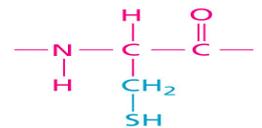
### glycine

(Gly, or G)



### cysteine

(Cys, or C)

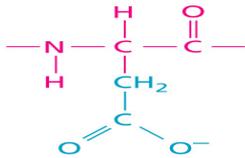


## NONPOLAR SIDE CHAINS

## ACIDIC SIDE CHAINS

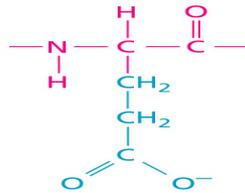
### aspartic acid

(Asp, or D)



### glutamic acid

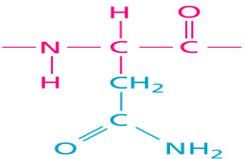
(Glu, or E)



## UNCHARGED POLAR SIDE CHAINS

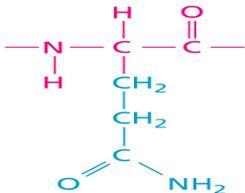
### asparagine

(Asn, or N)



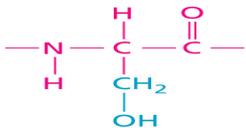
### glutamine

(Gln, or Q)



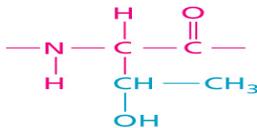
### serine

(Ser, or S)



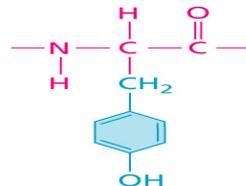
### threonine

(Thr, or T)

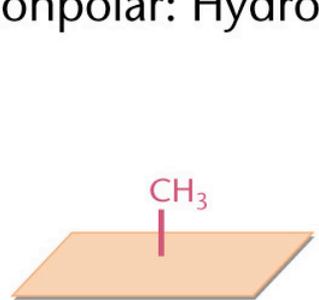


### tyrosine

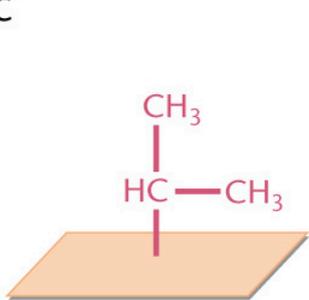
(Tyr, or Y)



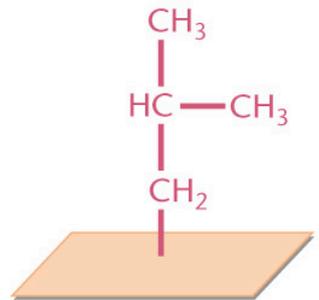
# 1. Nonpolar: Hydrophobic



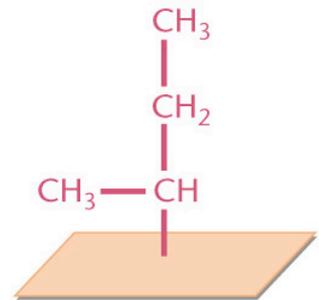
Alanine (ala-A)



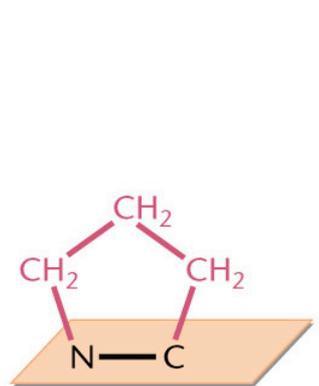
Valine (val-V)



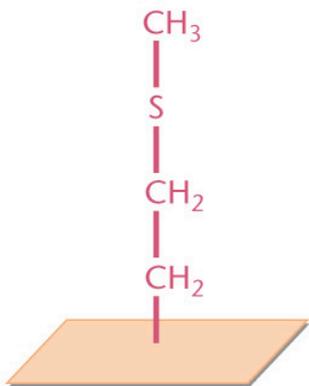
Leucine (leu-L)



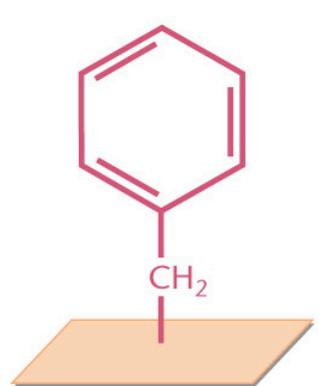
Isoleucine (ile-I)



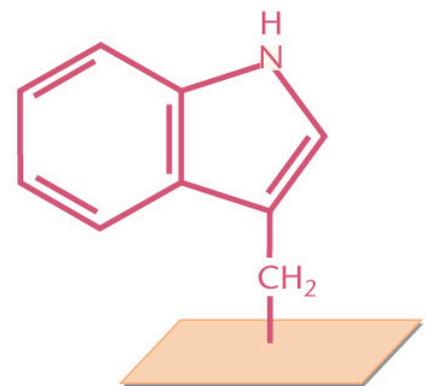
Proline (pro-P)



Methionine (met-M)



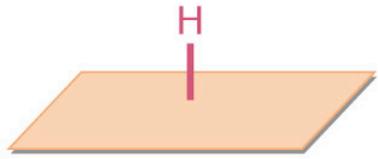
Phenylalanine (phe-F)



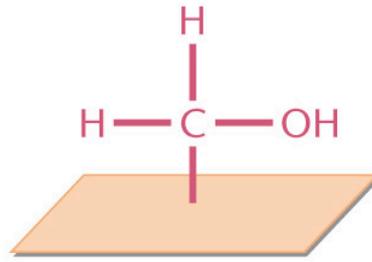
Tryptophan (trp-W)

Amino Acid Structures from Klug & Cummings

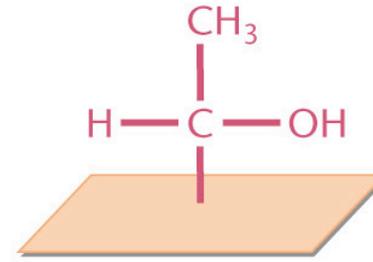
## 2. Polar: Hydrophilic



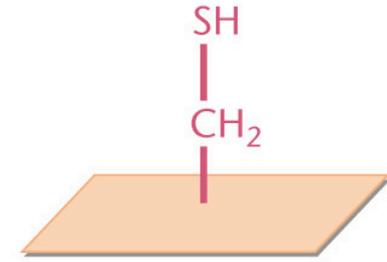
Glycine (gly-G)



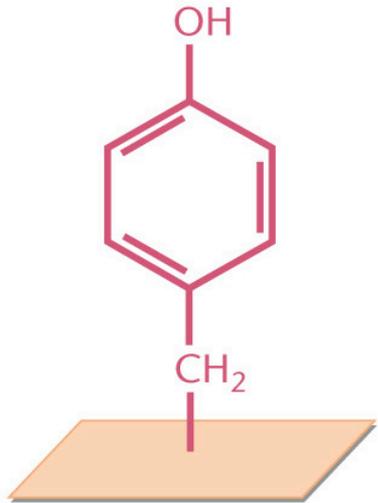
Serine (ser-S)



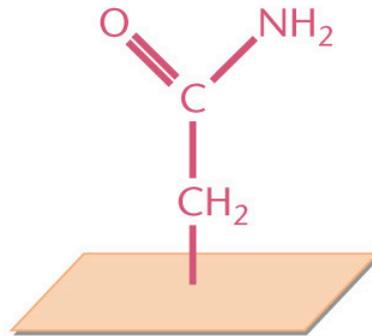
Threonine (thr-T)



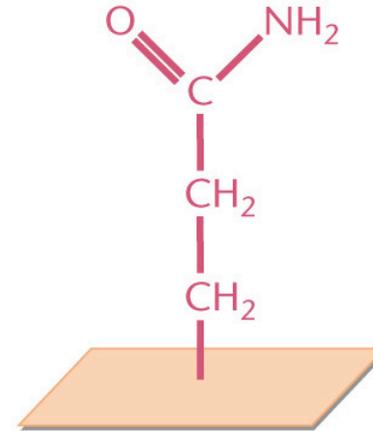
Cysteine (cys-C)



Tyrosine (tyr-Y)



Asparagine (asn-N)



Glutamine (gln-Q)

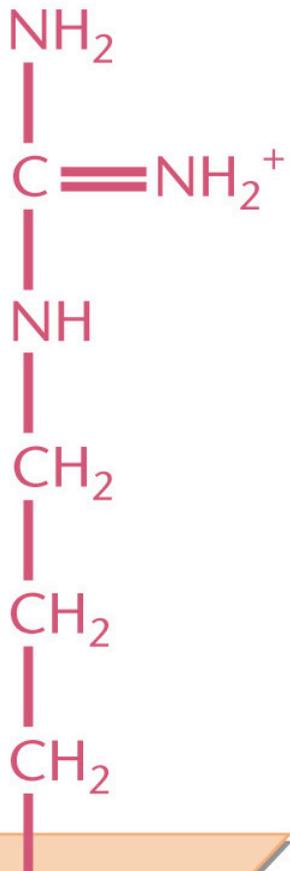
Amino Acid Structures from Klug & Cummings

### 3. Polar: positively charged (basic)

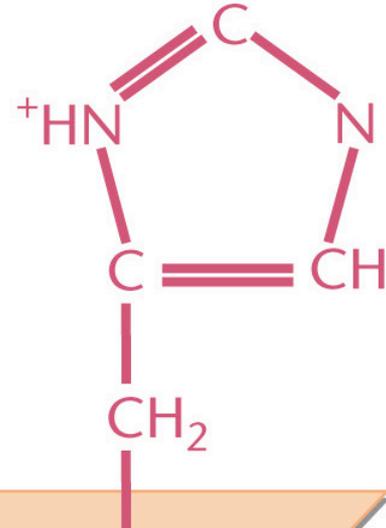
Amino Acid Structures  
from Klug & Cummings



Lysine (lys-K)

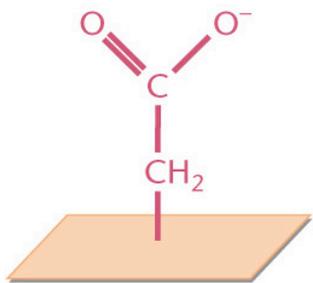


Arginine (arg-R)

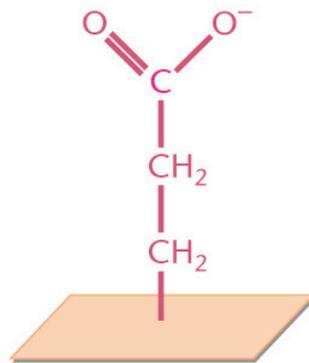


Histidine (his-H)

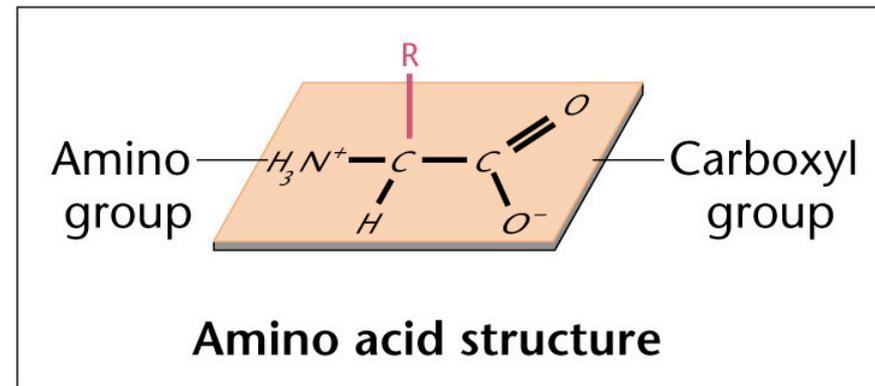
#### 4. Polar: negatively charged (acidic)



Aspartic acid (asp-D)

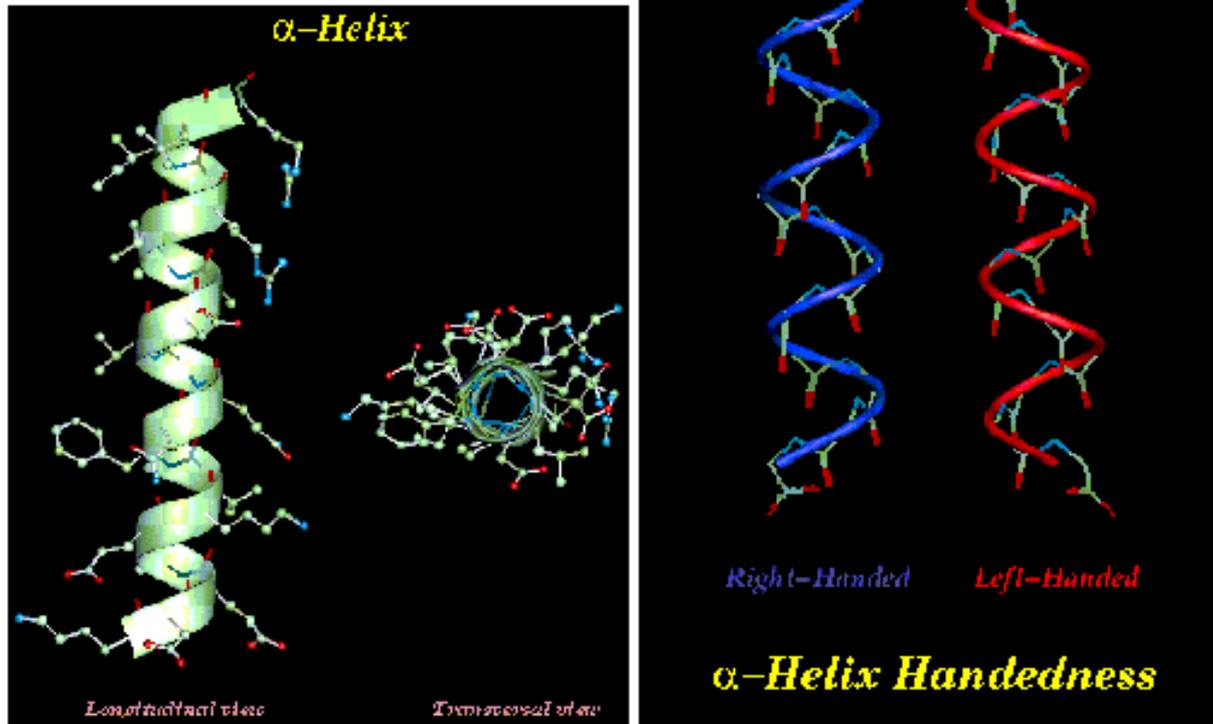


Glutamic acid (glu-E)



Amino Acid Structures from Klug & Cummings

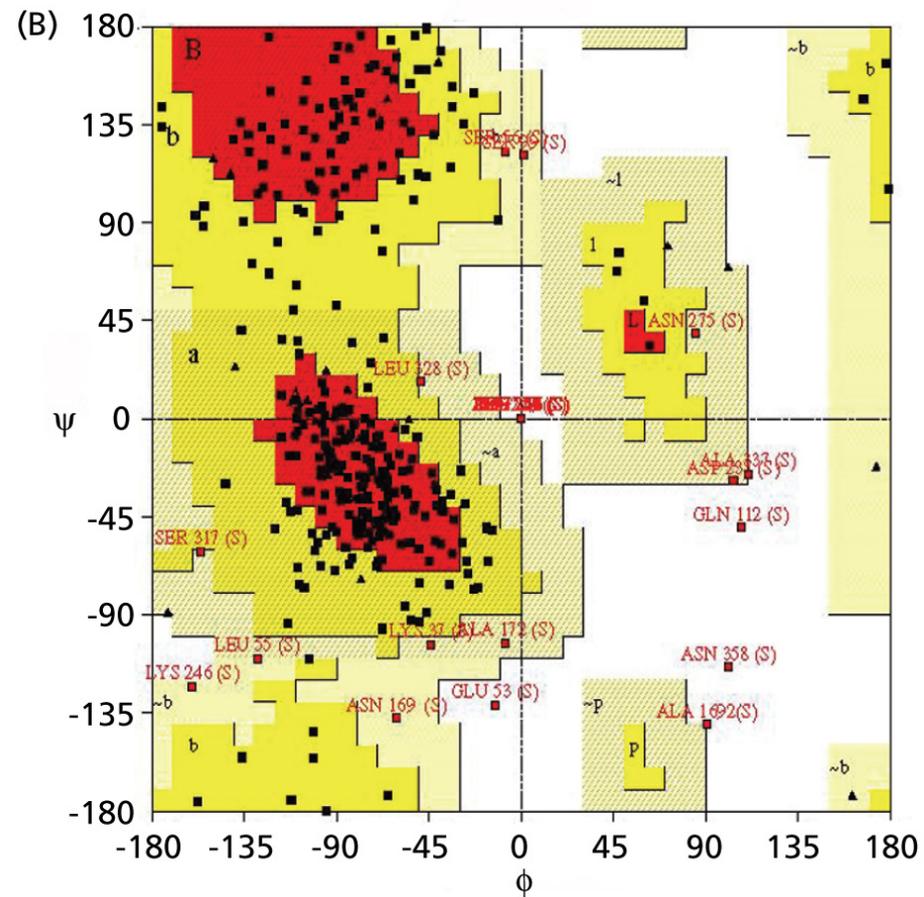
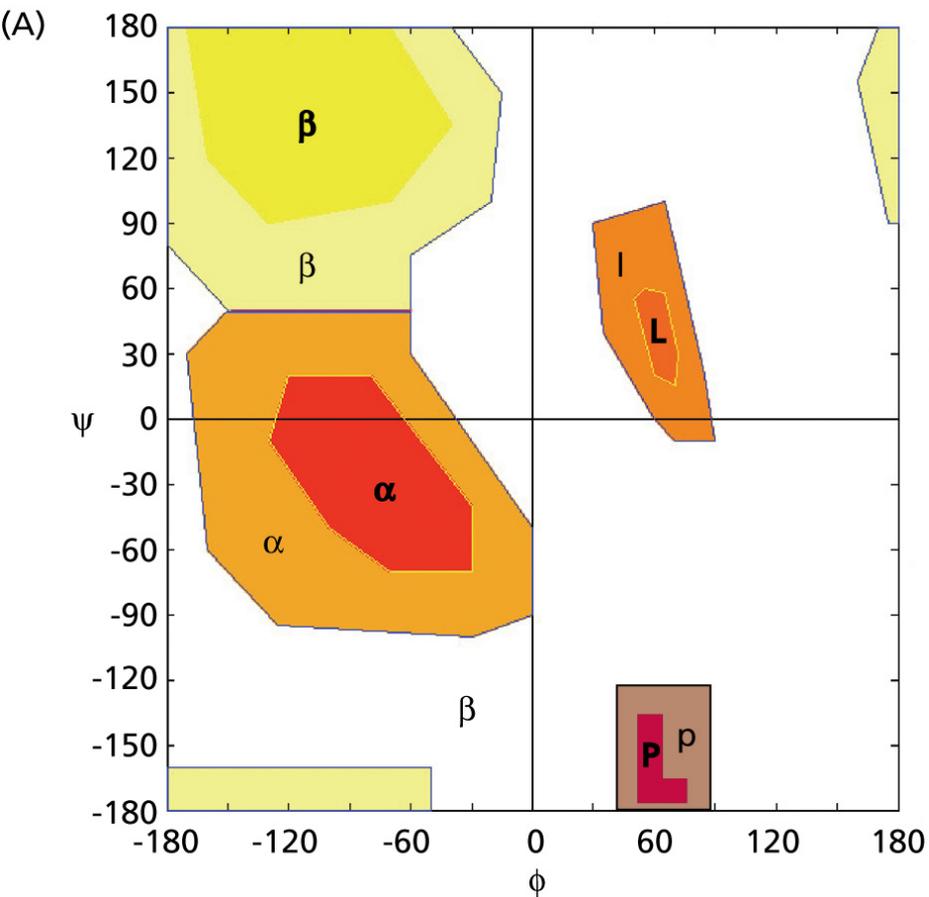
## Alpha helices

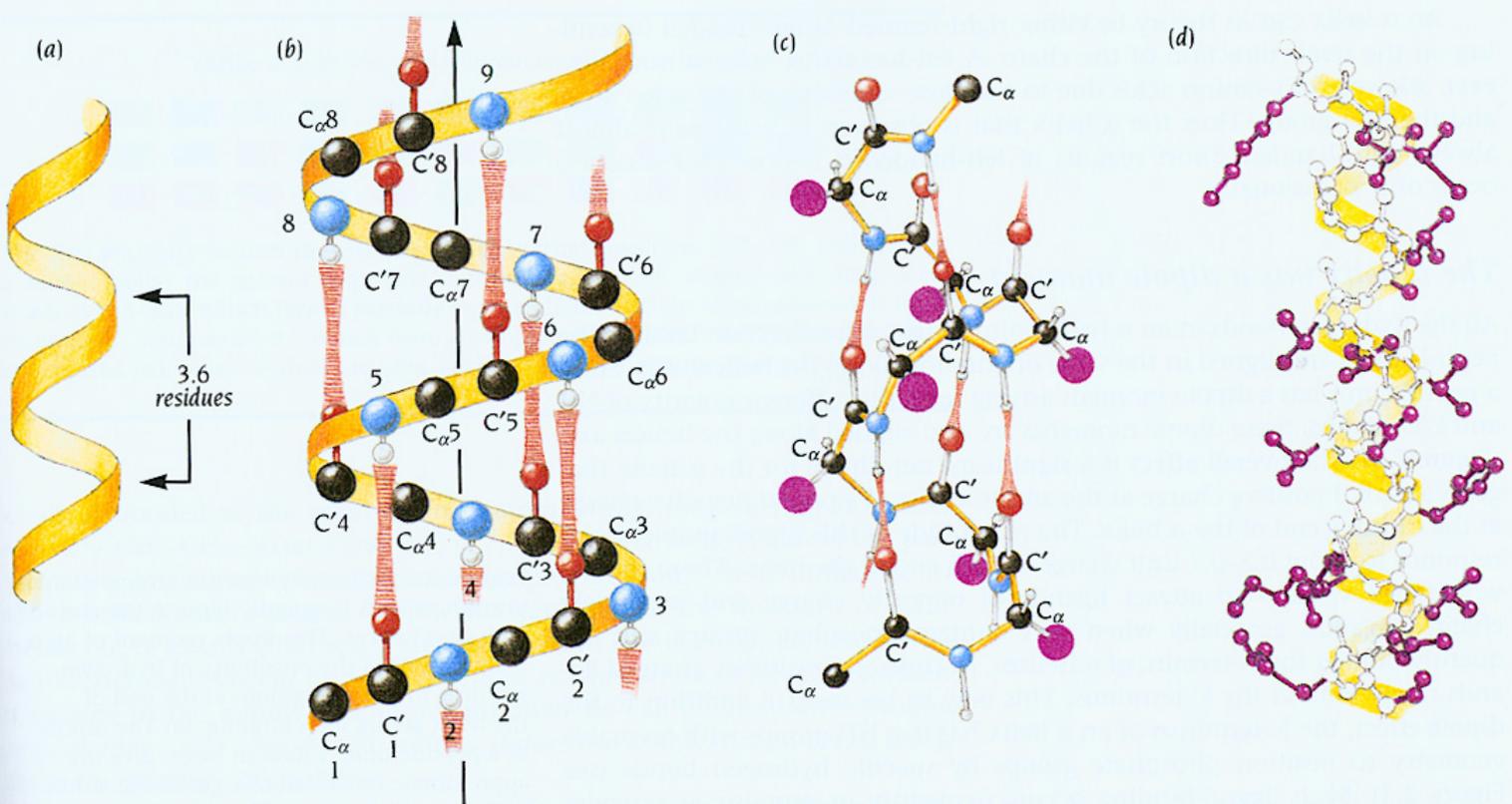


(c) David Gilbert, Aik Choon Tan, Gilleain Torrance and Mallika Veeramalai 2002

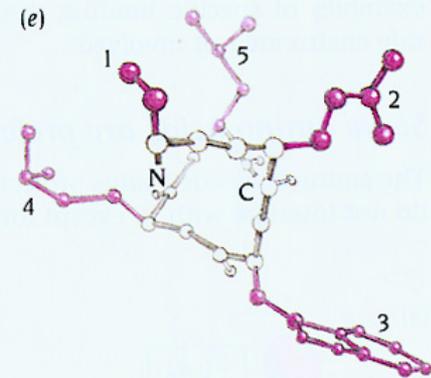
16

# Ramachandran Plot

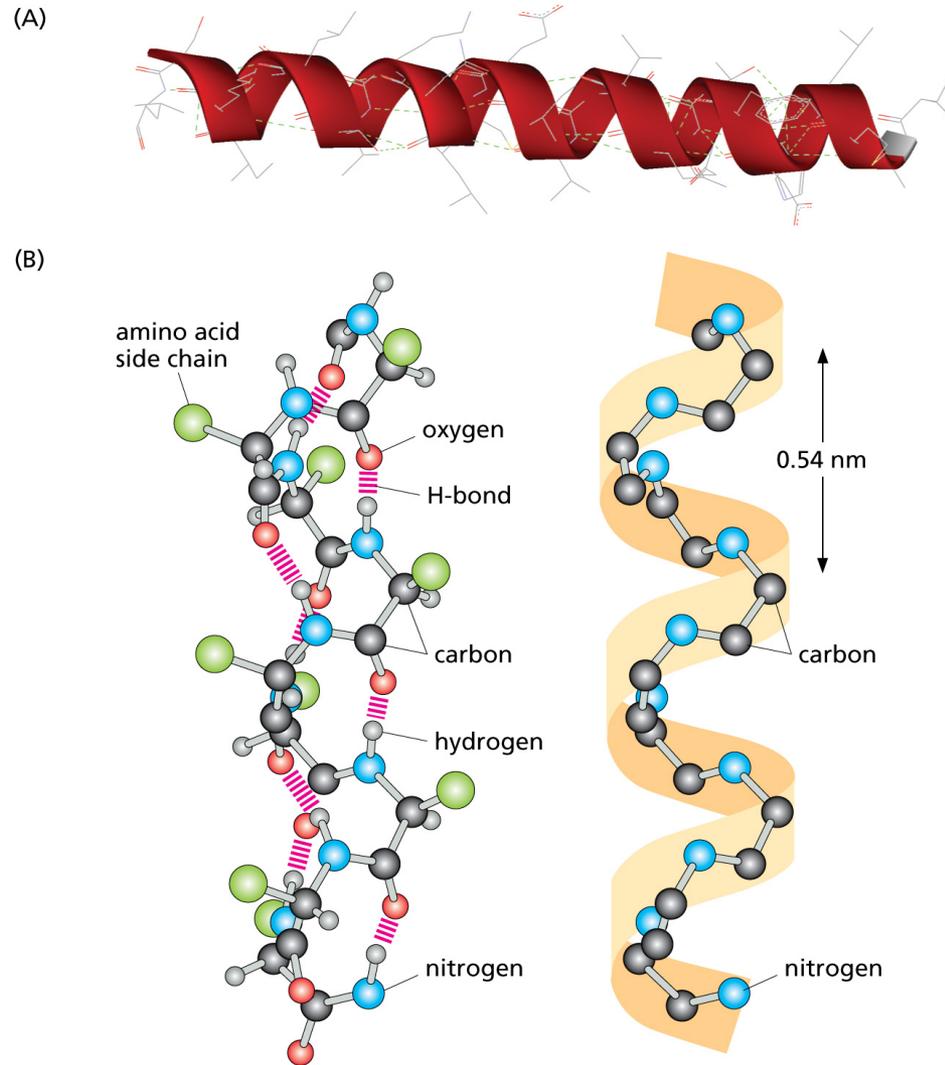




**Figure 2.2** The  $\alpha$  helix is one of the major elements of secondary structure in proteins. Main-chain N and O atoms are hydrogen-bonded to each other within  $\alpha$  helices. (a) Idealized diagram of the path of the main chain in an  $\alpha$  helix. Alpha helices are frequently illustrated in this way. There are 3.6 residues per turn in an  $\alpha$  helix, which corresponds to 5.4 Å (1.5 Å per residue). (b) The same as (a) but with approximate positions for main-chain atoms and hydrogen bonds included. The arrow denotes the direction from the N-terminus to the C-terminus. (c) Schematic diagram of an  $\alpha$  helix. Oxygen atoms are red, and N atoms are blue. Hydrogen bonds between O and N are red and striated. The side chains are represented as purple circles. (d) A ball-and-stick model of one  $\alpha$  helix in myoglobin. The path of the main chain is outlined in yellow; side chains are purple. Main-chain atoms are not colored. (e) One turn of an  $\alpha$  helix viewed down the helical axis. The purple side chains project out from the  $\alpha$  helix.

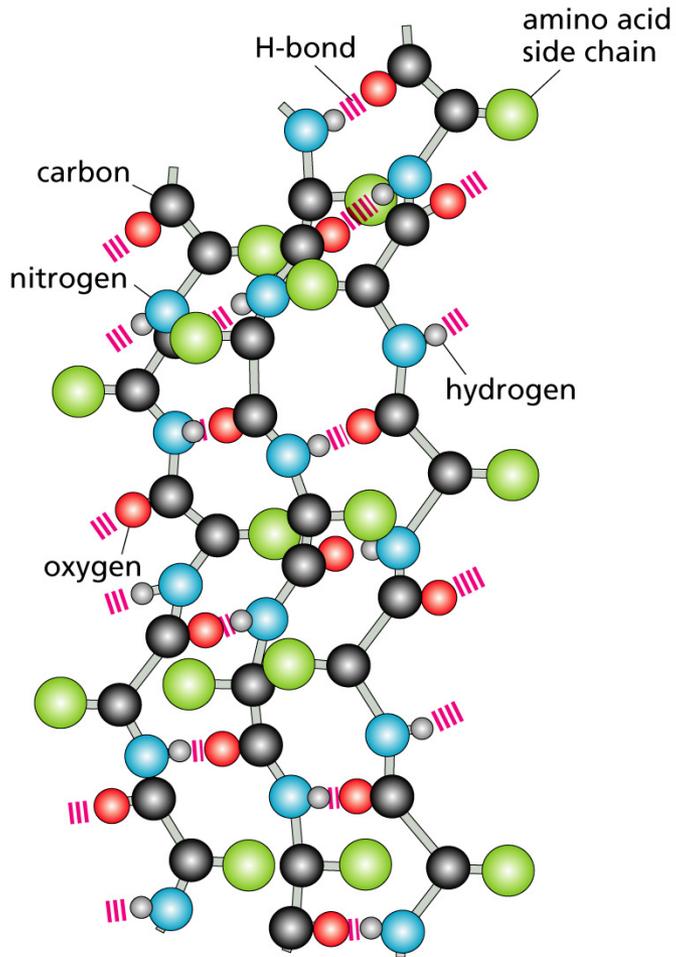


# Alpha Helix

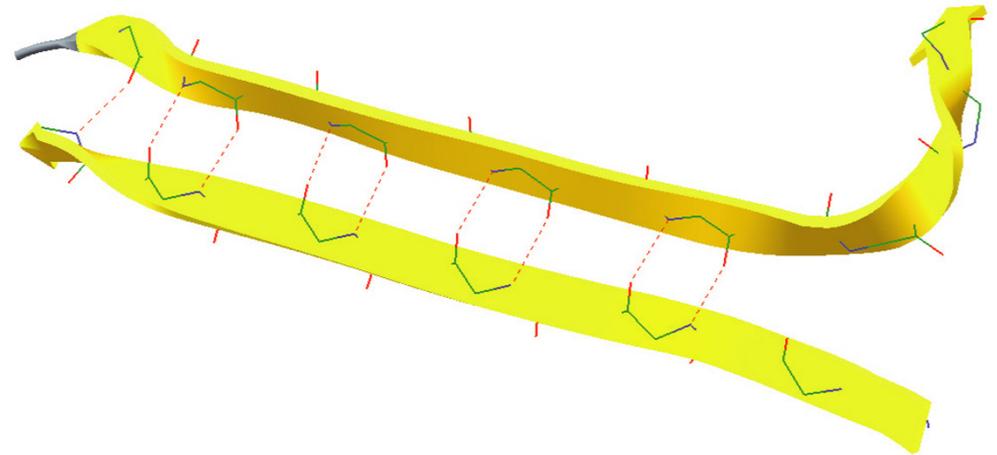


# Beta Strands and Sheets

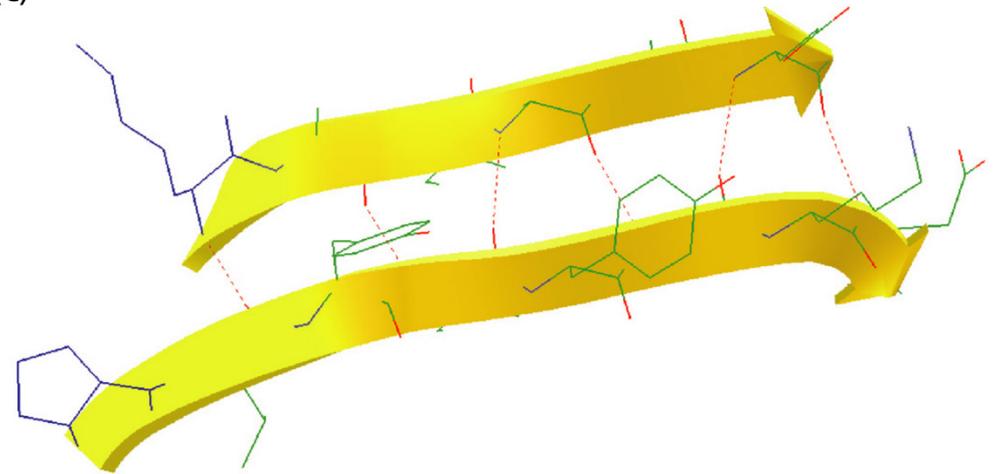
(A)



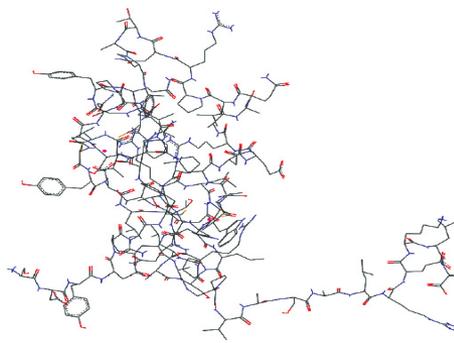
(B)



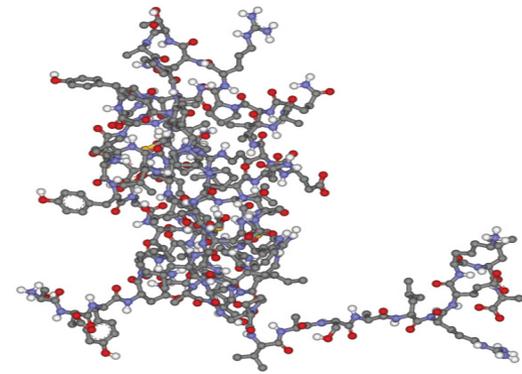
(C)



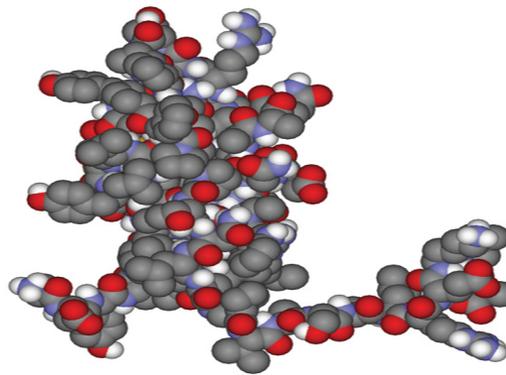
# Molecular Representations



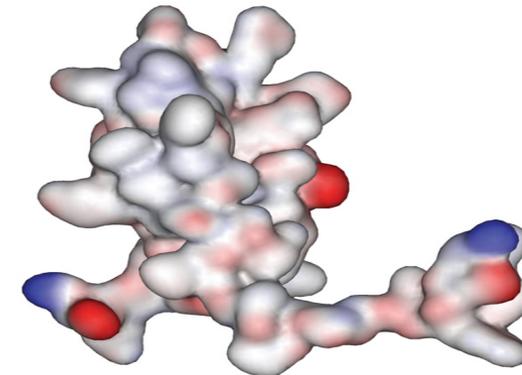
wire-frame



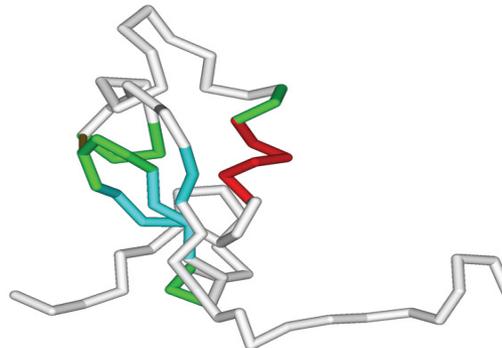
ball and stick



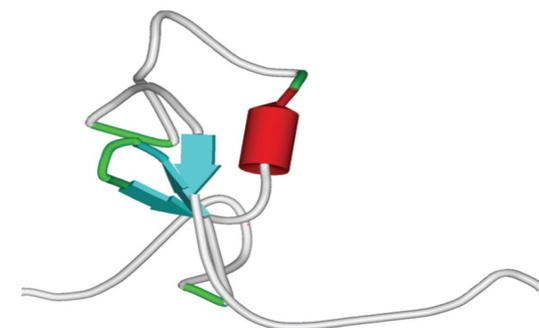
space-filling



surface



C $\alpha$  representation

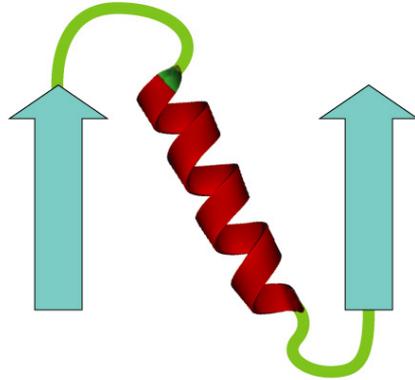


$\alpha/\beta$  schematic

CAP5510 / CGS5166

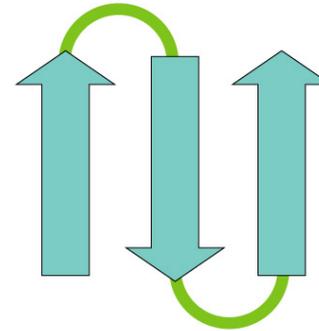
# Supersecondary structures

(A)



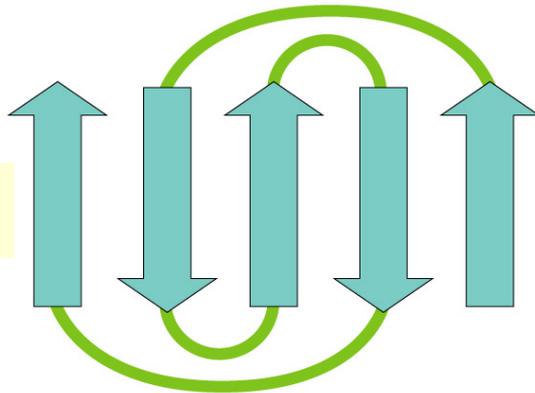
$\beta\alpha\beta$  repeat

(B)



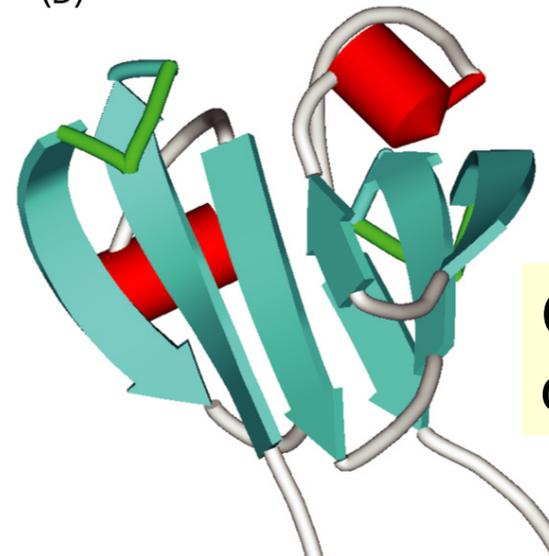
$\beta\alpha\beta$ -meander

(C)



Greek Key

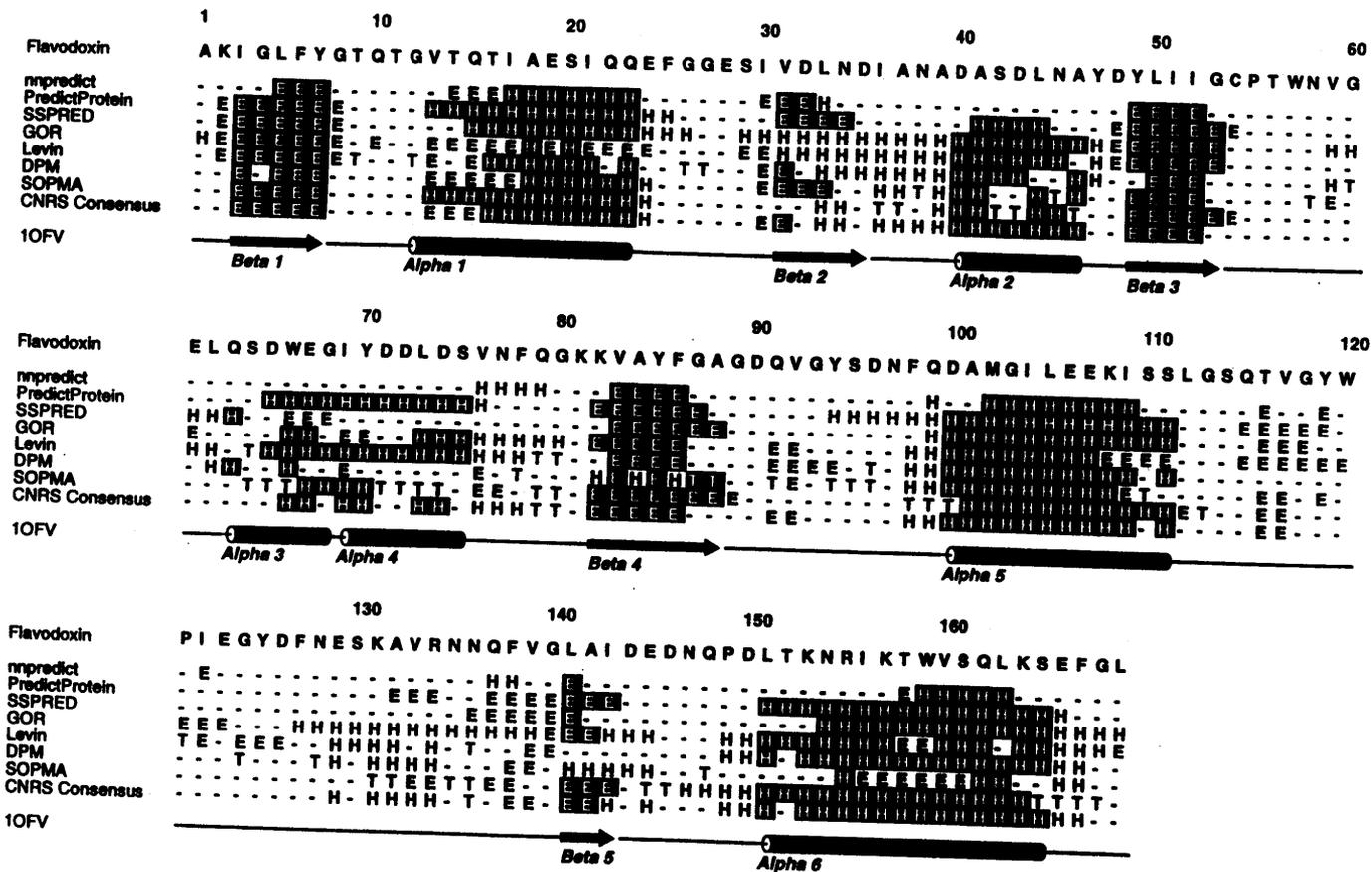
(D)



Gamma  $\beta$   
crystallin

# Secondary Structure Prediction Software

254

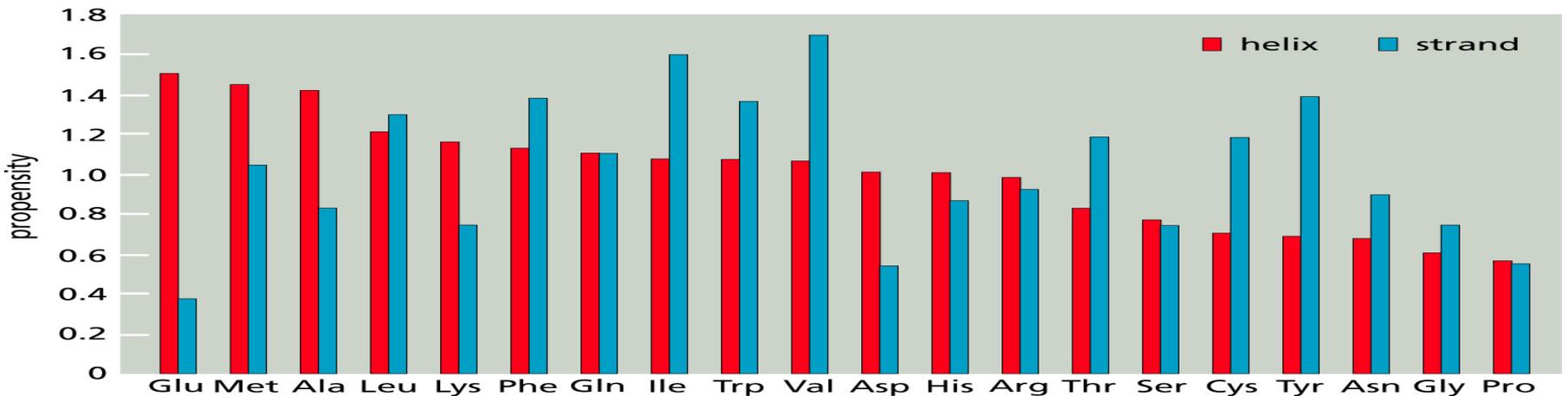


Recent Ones:  
 GOR V  
 PREDATOR  
 Zpred  
 PROF  
 NNSSP  
 PHD  
 PSIPRED  
 Jnet

**Figure 11.3** Comparison of secondary structure predictions by various methods. The sequence of flavodoxin, an  $\alpha/\beta$  protein, was used as the query and is shown on the first line of the alignment. For each prediction, H denotes an  $\alpha$  helix, E a  $\beta$  strand, T a  $\beta$  turn; all other positions are assumed to be random coil. Correctly assigned residues are shown in inverse type. The methods used are listed along the left side of the alignment and are described in the text. At the bottom of the figure is the secondary structure assignment given in the PDB file for flavodoxin (1OFV, Smith et al., 1983).

# Chou & Fasman Propensities

Amino Acid	helix		strand	
	Designation	<i>P</i>	Designation	<i>P</i>
Ala	F	1.42	b	0.83
Cys	l	0.70	f	1.19
Asp	l	1.01	B	0.54
Glu	F	1.51	B	0.37
Phe	f	1.13	f	1.38
Gly	B	0.61	b	0.75
His	f	1.00	f	0.87
Ile	f	1.08	F	1.60
Lys	f	1.16	b	0.74
Leu	F	1.21	f	1.30
Met	F	1.45	f	1.05
Asn	b	0.67	b	0.89
Pro	<b>B</b>	<b>0.57</b>	<b>B</b>	<b>0.55</b>
Gln	f	1.11	h	1.10
Arg	l	0.98	l	0.93
Ser	l	0.77	b	0.75
Thr	l	0.83	f	1.19
Val	f	1.06	F	1.70
Trp	f	1.08	f	1.37
Tyr	b	0.69	F	1.4



# GOR IV prediction for 1bbc

A FAGVLNDADIAAALEACKAADSFNHKAFFAKVGLTSKSADDVKKAF A I I  
C C C C C C H H H H H H H H H H H H C C C C C H H H H E E C C C C C H H H H H H H H H H  
A Q D K S G F I E E D E L K L F L Q N F K A D A R A L T D G E T K T F L K A G D S D G D G K I G V D  
H H C C C C H H H H H H H H H H H H H H H H C C C C C E E E E E E C C C C C C C E E E E C  
D V T A L V K A  
C E E E E E E C

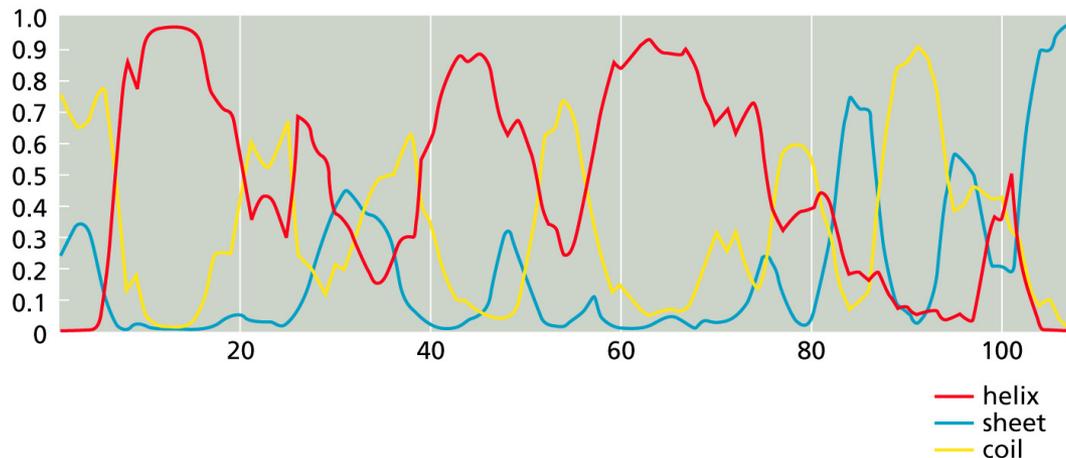
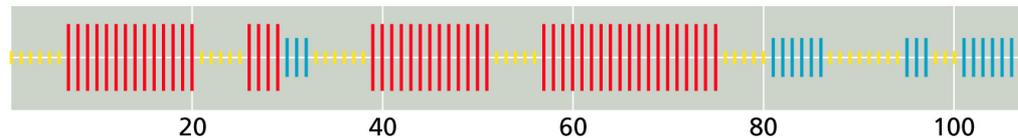
sequence length: 108

GOR IV:

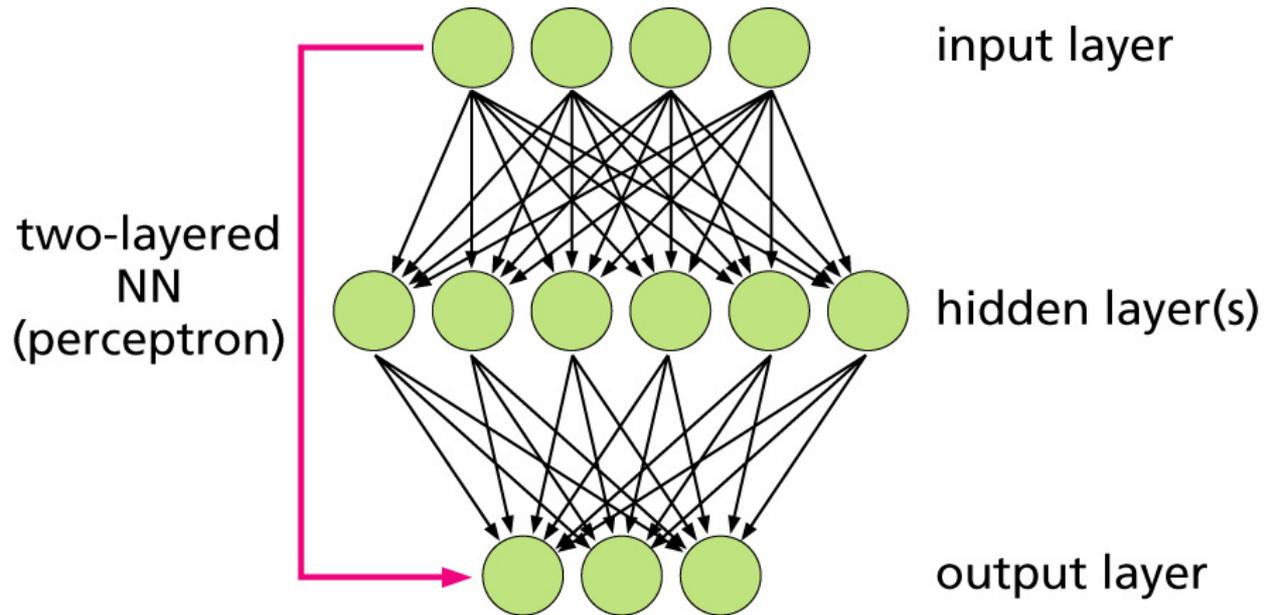
alpha helix (Hh) : 50 is 46.30%

beta sheet (Ee) : 18 is 16.67%

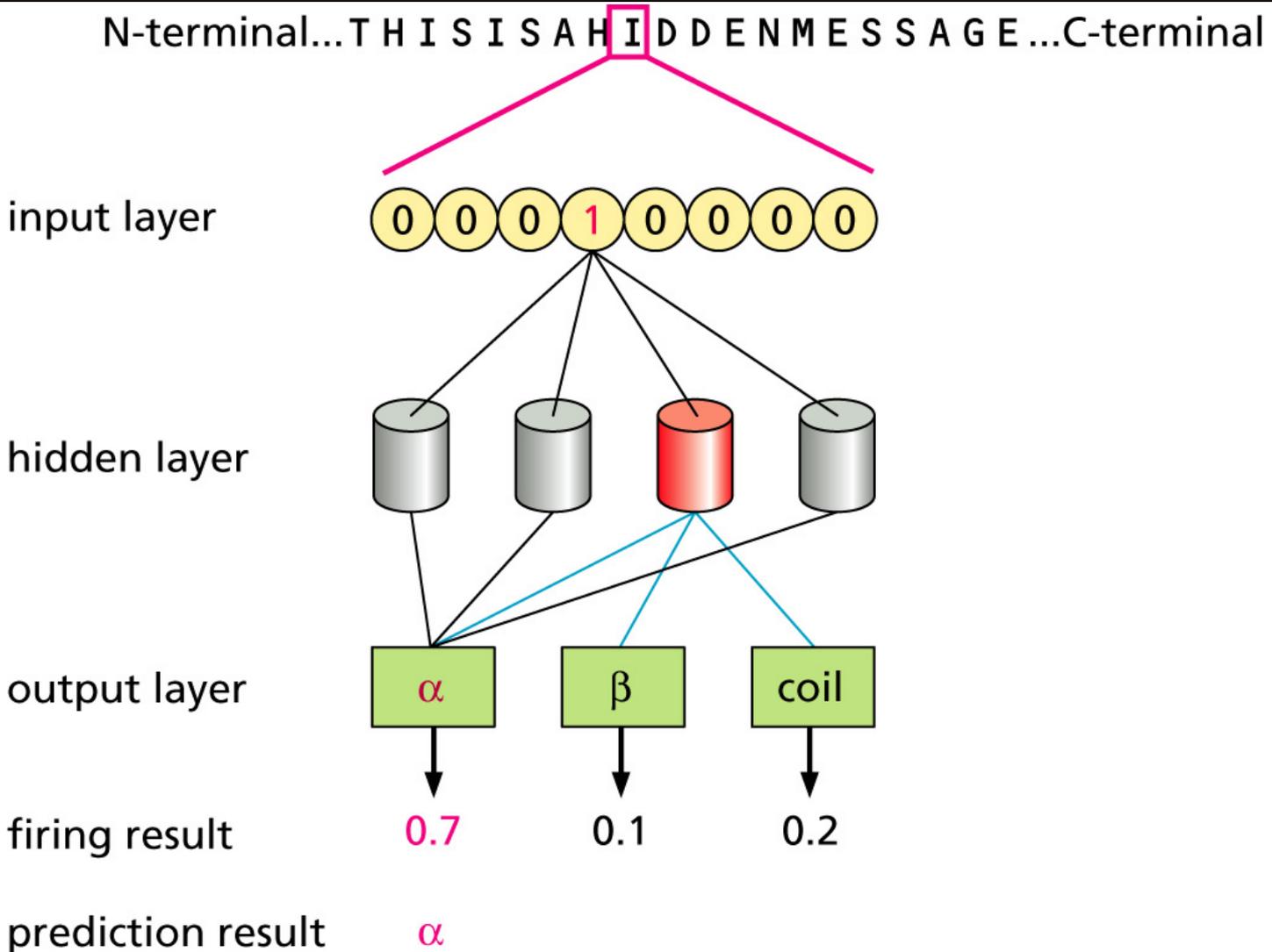
random coil (Cc) : 40 is 37.04%



# Neural Networks



# Neural Network Prediction of SS



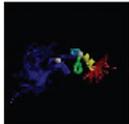
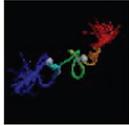
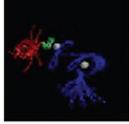
# PDB: Protein Data Bank

- ❑ Database of protein tertiary and quaternary structures and protein complexes. <http://www.rcsb.org/pdb/>
- ❑ Over 29,000 structures as of Feb 1, 2005.
- ❑ Structures determined by
  - NMR Spectroscopy
  - X-ray crystallography
  - Computational prediction methods
- ❑ Sample PDB file: [Click here](#) [ [\\_](#) ]

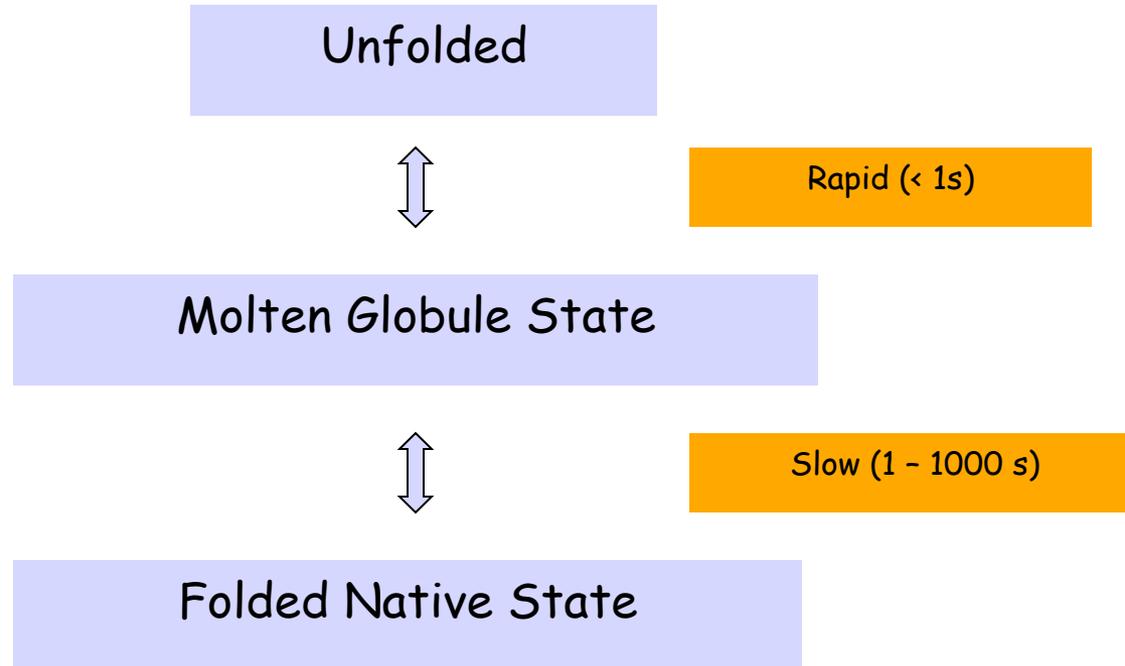
# PDB Search Results

- Results (1-10 of 91)
- Results ID List
- Refine this Search
- 1 Structures Awaiting Release
- Select All
- Deselect All
- Download Selected
- ▶ Tabulate
- ▶ Narrow Query
- ▶ Sort Results
- ▶ Results per Page
- Show Query Details
- Results Help

1 2 3 4 5 .. 10 ↩

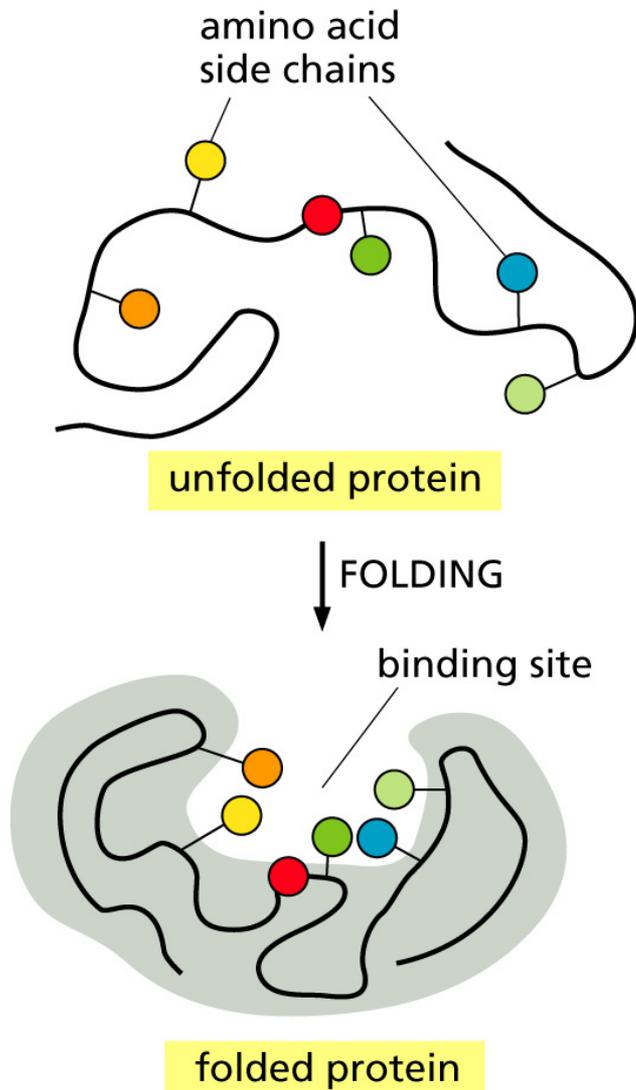
- |                                     |                        |   |   |  |
|-------------------------------------|------------------------|---|---|--|
| <input checked="" type="checkbox"/> | <b>1X6Z</b>            |    |          | <b>Solution structure of the LIM domain of carboxyl terminal LIM domain protein 1</b>          |
|                                     | <i>Characteristics</i> |   |   | <b>Release Date:</b> 17-Nov-2005 <b>Exp. Method:</b> NMR 20 Structures                         |
|                                     | <i>Classification</i>  |   |   | <b>Structural Protein</b>  |
|                                     | <i>Compound</i>        |   |   | <b>Mol. Id:</b> 1 <b>Molecule:</b> C Terminal Lim Domain Protein 1 <b>Fragment:</b> Lim Domain |
|                                     | <i>Authors</i>         |   |   | <b>Qin, X.R., Nagashima, T., Hayashi, F., Yokoyama, S.</b>                                     |
| <hr/>                               |                        |   |   |  |
| <input checked="" type="checkbox"/> | <b>1X4K</b>            |   |          | <b>Solution structure of LIM domain in LIM-protein 3</b>                                       |
|                                     | <i>Characteristics</i> |   |   | <b>Release Date:</b> 14-Nov-2005 <b>Exp. Method:</b> NMR 20 Structures                         |
|                                     | <i>Classification</i>  |   |   | <b>Metal Binding Protein</b>   |
|                                     | <i>Compound</i>        |   |   | <b>Mol. Id:</b> 1 <b>Molecule:</b> Skeletal Muscle Lim Protein 3 <b>Fragment:</b> Lim Domain   |
|                                     | <i>Authors</i>         |   |   | <b>He, F., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama,</b>            |
| <hr/>                               |                        |   |   |  |
| <input checked="" type="checkbox"/> | <b>1X4L</b>            |  |    | <b>Solution structure of LIM domain in Four and a half LIM domains protein 2</b>               |
|                                     | <i>Characteristics</i> |   |   | <b>Release Date:</b> 14-Nov-2005 <b>Exp. Method:</b> NMR 20 Structures                         |
|                                     | <i>Classification</i>  |   |   | <b>Metal Binding Protein</b>   |
|                                     | <i>Compound</i>        |   |   | <b>Mol. Id:</b> 1 <b>Molecule:</b> Skeletal Muscle Lim Protein 3 <b>Fragment:</b> Lim Domain   |
|                                     | <i>Authors</i>         |   |   | <b>He, F., Muto, Y., Inoue, M., Kigawa, T., Shirouzu, M., Terada, T., Yokoyama,</b>            |

# Protein Folding

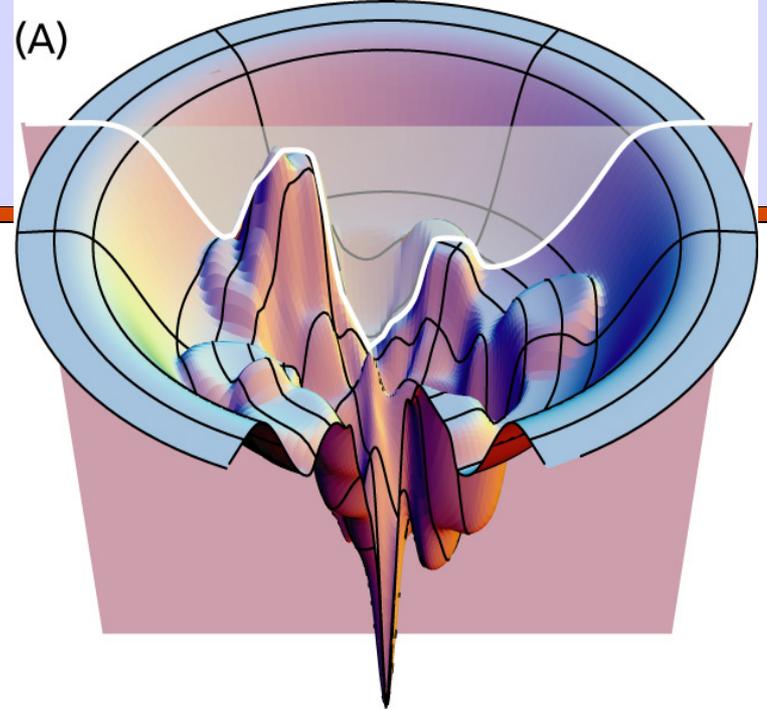


□ How to find minimum energy configuration?

# Protein Folding



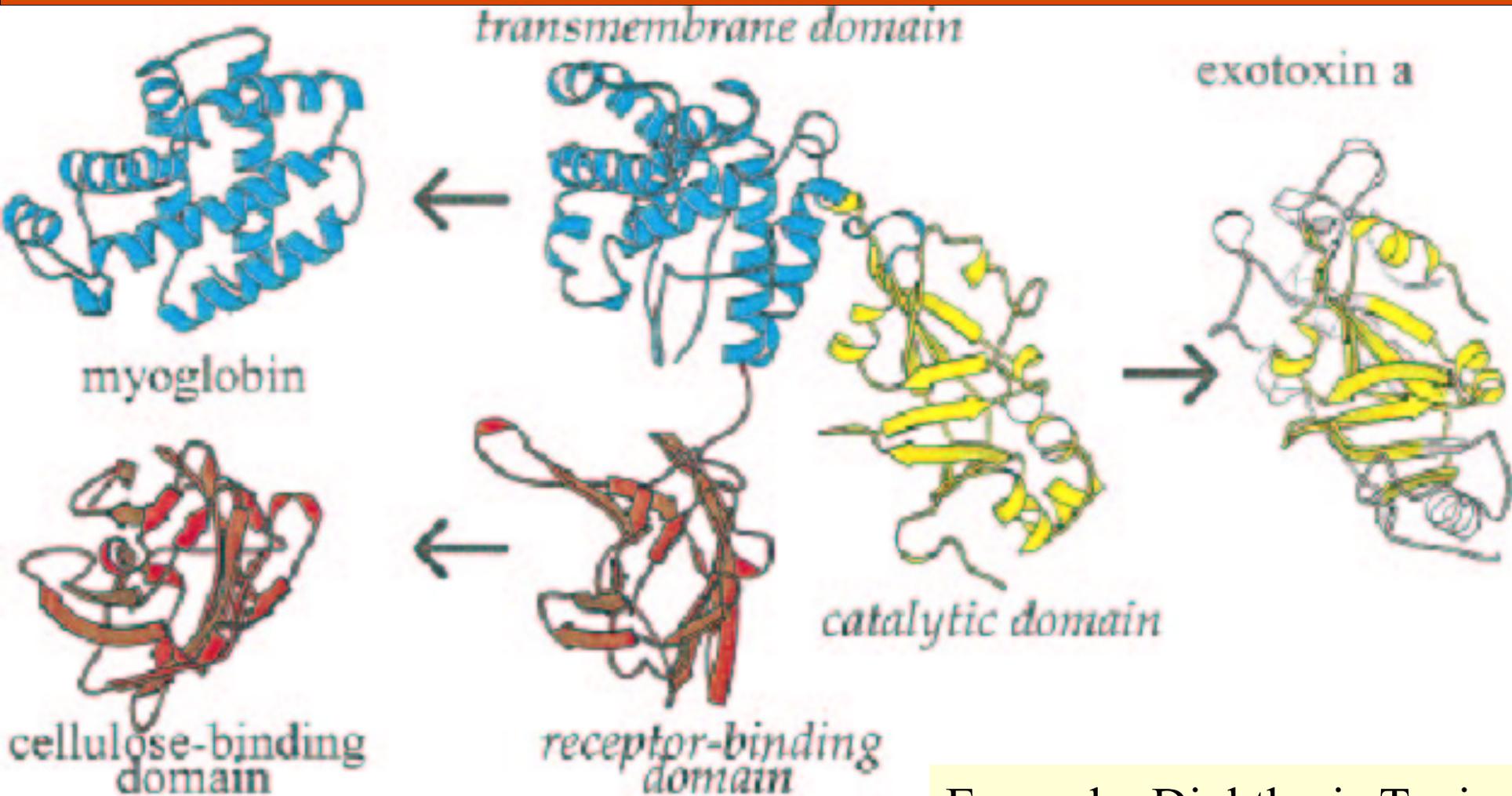
# Energy Landscape



(B)



# Modular Nature of Protein Structures



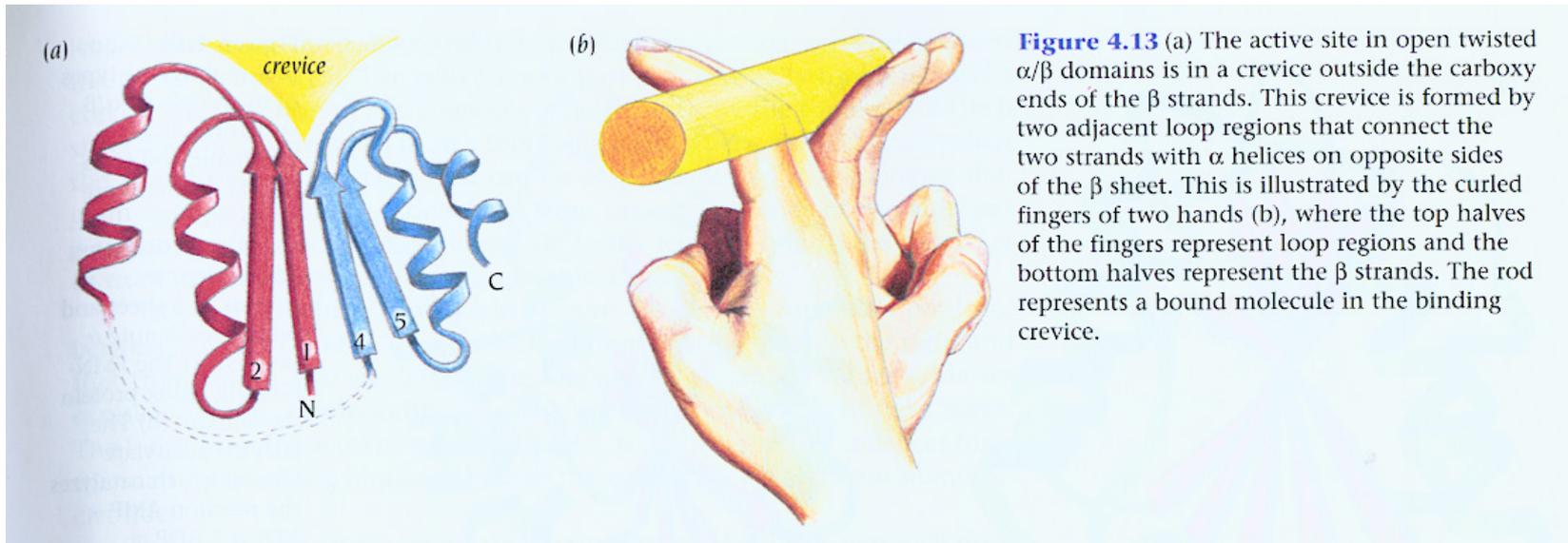
Example: Diphtheria Toxin

# Protein Structures

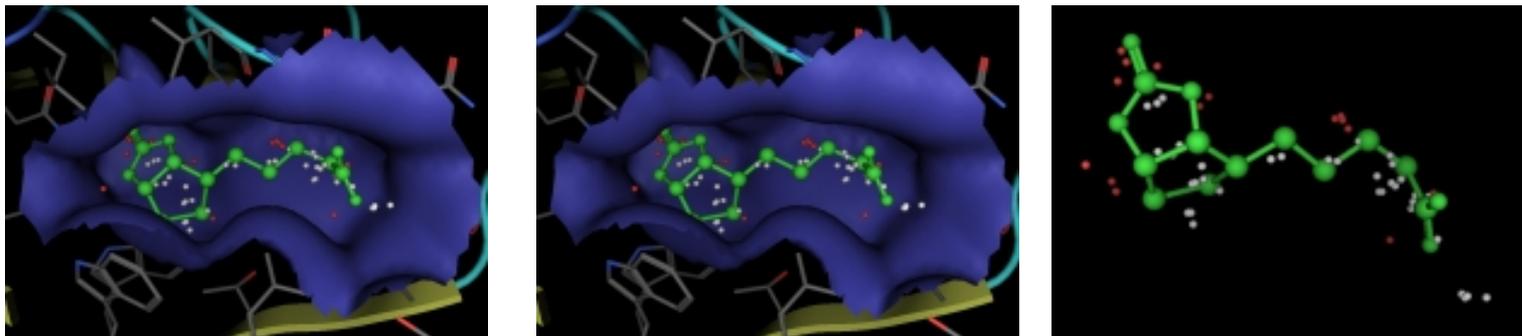
- ❑ Most proteins have a **hydrophobic core**.
- ❑ Within the core, specific **interactions** take place between amino acid side chains.
- ❑ Can an amino acid be replaced by some other amino acid?
  - Limited by space and available contacts with nearby amino acids
- ❑ Outside the core, proteins are composed of loops and structural elements in contact with water, solvent, other proteins and other structures.

# Active Sites

Active sites in proteins are usually hydrophobic pockets/crevices/troughs that involve sidechain atoms.



# Active Sites



**Left** PDB 3RTD (streptavidin) and the first site located by the MOE Site Finder. **Middle** 3RTD with complexed ligand (biotin). **Right** Biotin ligand overlaid with calculated alpha spheres of the first site.

# Viewing Protein Structures

- SPDBV
- RASMOL
- CHIME

# Structural Classification of Proteins

- Over 1000 protein families known
  - Sequence alignment, motif finding, block finding, similarity search
- **SCOP** (Structural Classification of Proteins)
  - Based on structural & evolutionary relationships.
  - Contains ~ 40,000 domains
  - Classes (groups of folds), Folds (proteins sharing folds), Families (proteins related by function/evolution), Superfamilies (distantly related proteins)

## SCOP Family View

The screenshot displays the NCSA Mosaic WWW browser interface. The main window shows the SCOP Family View for the Interleukin 8-like family. The interface includes a menu bar (File, Options, Navigate, Annotate, Help), document title (SCOP: Family: Interleukin 8-like), and document URL (http://scop.mrc-lmb.cam.ac.uk/scop/data/scop.0.004). The structural classification of proteins is shown, including the lineage (Root: scop, Class: Alpha, Fold: Interleu, Superfamily: Interleukin 8-like, Family: Interleukin 8-like) and a list of proteins (Interleukin-8, Platelet factor 4, Macrophage inflammatory protein 1beta). Annotations with arrows point to specific elements: 'click here to display protein in 3D-viewer' points to the 'Interleu' link; 'click here for sequence and references (NCBI)' points to the 'Interleukin 8-like' link; 'PDB entry names' points to the 'human (Homo sapiens) (3)' entry; 'click here to fetch image' points to the '1.hum' entry; and 'keyword search facility' points to the search box at the bottom. On the right, the RasMol 2.4 window shows a 3D wireframe model of a protein structure, and the xv 3.00 window shows a static image of Human MIP-1β and Interleukin 8 Dimers.

**Figure 2.** A typical scop session is shown on a unix workstation. A scop page, of the Interleukin 8-like family, is displayed by the WWW browser program (NCSA Mosaic) (Schatz & Hardin, 1994). Navigating through the tree structure is accomplished by selecting any underlined entry, by clicking on buttons (at the top of each page) and by keyword searching (at the bottom of each page). The static image comparing two proteins in this family was downloaded by clicking on the icon indicated and is displayed by image-viewer program xv. By clicking on one of the green icons, commands were sent to a molecular viewer program (RasMol) written by Roger Sayle (Sayle, 1994), instructing it to automatically display the relevant PDB file and colour the domain in question by secondary structure. Since sending large PDB files over the network can be slow, this feature of scop can be configured to use local copies of PDB files if they are available. Equivalent WWW browsers, image-display programs and molecular viewers are also available free for Windows-PC and Macintosh platforms.

# CATH: Protein Structure Classification

- Semi-automatic classification; ~36K domains
- 4 levels of classification:
  - Class (C), depends on sec. Str. Content
    - $\alpha$  class,  $\beta$  class,  $\alpha/\beta$  class,  $\alpha+\beta$  class
  - Architecture (A), orientation of sec. Str.
  - Topology (T), topological connections &
  - Homologous Superfamily (H), similar str and functions.

# DALI/FSSP Database

- ❑ Completely automated; 3724 domains
- ❑ Criteria of compactness & recurrence
- ❑ Each domain is assigned a Domain Classification number DC\_l\_m\_n\_p representing fold space attractor region (l), globular folding topology (m), functional family (n) and sequence family (p).

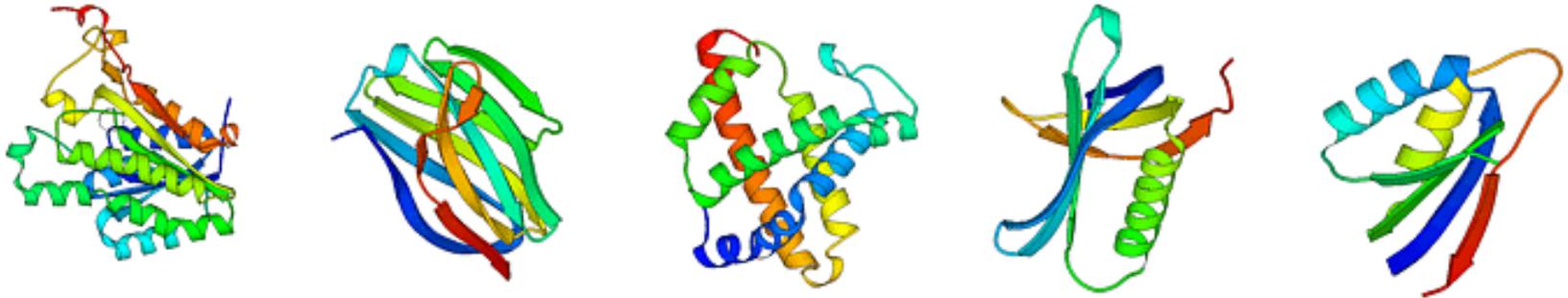
# Structural Alignment

- What is structural alignment of proteins?
  - 3-d superimposition of the atoms as "best as possible", i.e., to minimize RMSD (root mean square deviation).
  - Can be done using **VAST** and **SARF**
- Structural similarity is common, even among proteins that do not share sequence similarity or evolutionary relationship.

# Other databases & tools

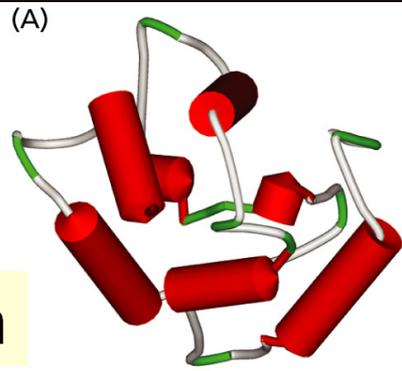
- ❑ **MMDB** contains groups of structurally related proteins
- ❑ **SARF** structurally similar proteins using secondary structure elements
- ❑ **VAST** Structure Neighbors
- ❑ **SSAP** uses double dynamic programming to structurally align proteins

# 5 Fold Space classes

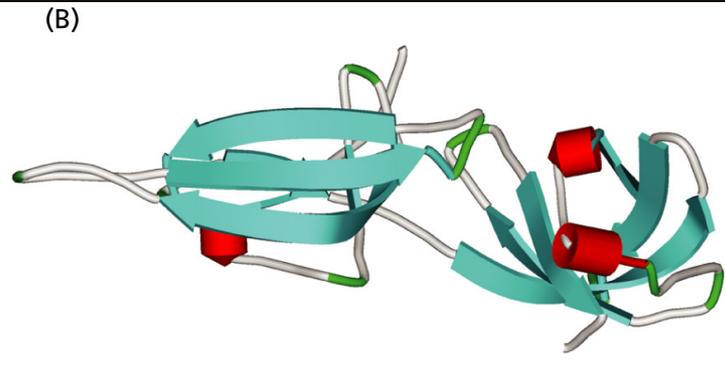


Attractor 1 can be characterized as alpha/beta, attractor 2 as all-beta, attractor 3 as all-alpha, attractor 5 as alpha-beta meander (1mli), and attractor 4 contains antiparallel beta-barrels e.g. OB-fold (1prtF).

# Examples of protein classes

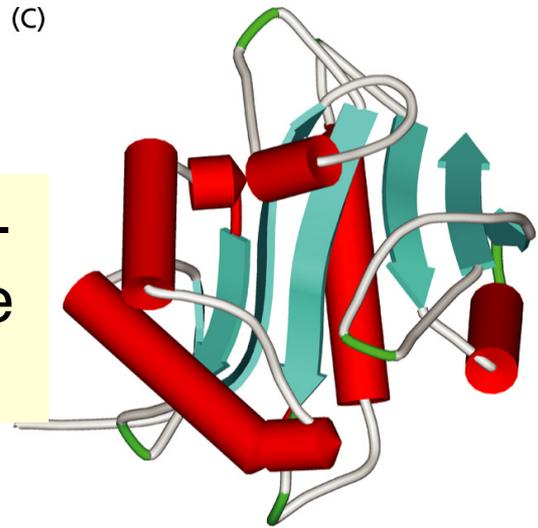


1B8C

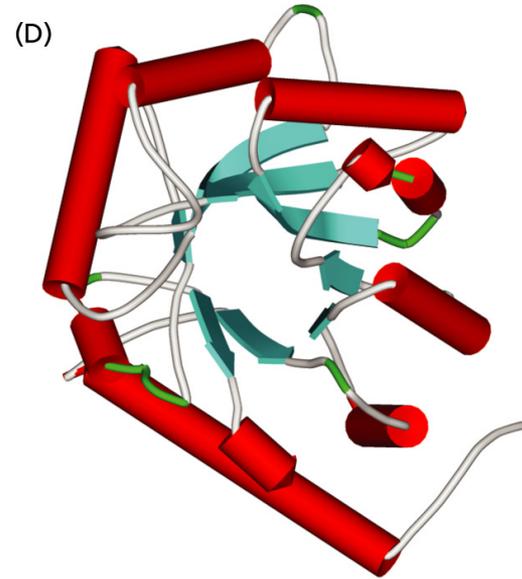


1BKB

Translation  
Initiation  
Factor 5A



1CJW



1CT5

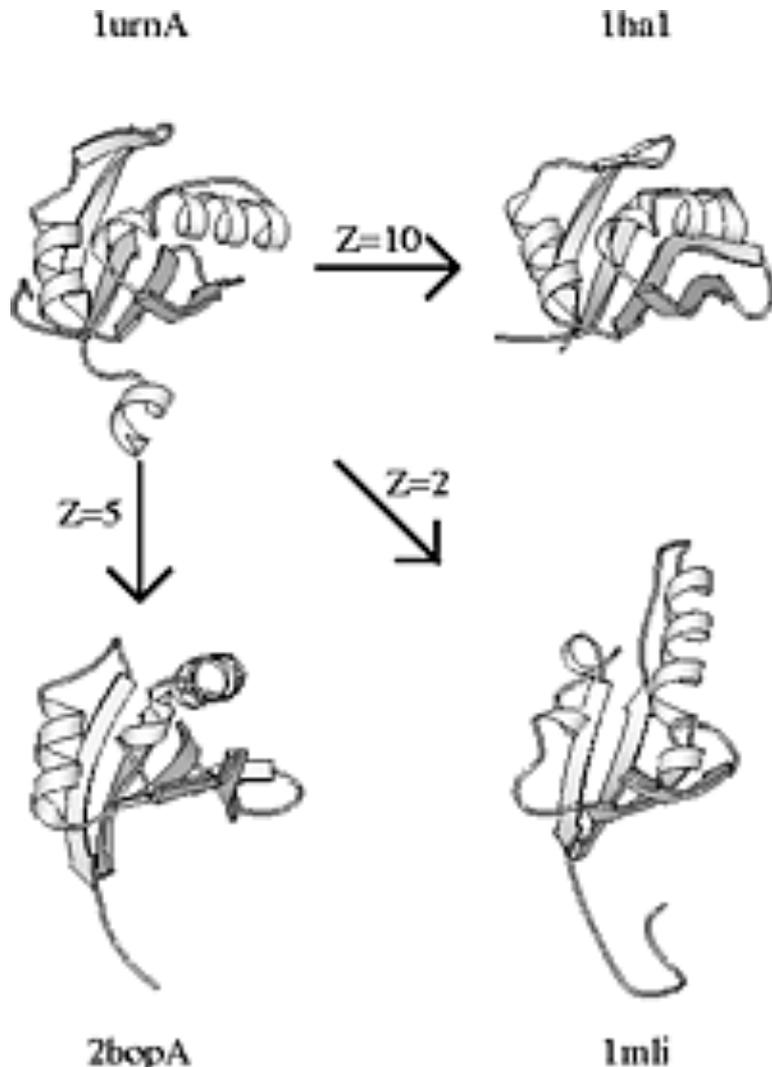
Hypothetical  
protein from  
yeast  
( $\alpha/\beta$  alternating fold)

Parvalbumin

Serotonin N-  
acetyltransfe  
rase

CAP5510 / CGS5166

# Fold Types & Neighbors



Structural neighbours of 1urnA (top left). 1mli (bottom right) has the same topology even though there are shifts in the relative orientation of secondary structure elements.

# Sequence Alignment of Fold Neighbors

**B**

```

1urnA  --RPNHTIYINNLNEKI-----KKDELKKSLSLHAIFSRFG---QILDILV-SRS---LKM---
Z=10      *          *              *  *          *  *          *
1ha1    ahLTVKKIFVGGIKEDT-----EEHHLRDYFEOYG---KIEVIEI-MTDrgsGKK---
Z=5      *
2bopA   ----sCFALIS-GTANQ-----vKCYRFRVKKNHRHR-----YENCTTtWFT---Vadnga
Z=2      *
1mli    ---mLFHVKMTVKLpvdmdpakatgkadeKELAQRlgregTWRHLWR-IAG-----

1urnA   ----RGQAFVIFKEV--SSATNALRSMQGFPPFYDKPMRIQYAKTDSDIIAKM-----
Z=10     **  ***  *          *              *
1ha1    ----RGFAFVTFDDH--DSVDKIVIO-kyHTVNGHNCSEVRKAL-----
Z=5      *  *          *  *          *  *  *
2bopA   erggQAQILITFGSP--SORODFLKHVPLPP---GMNISGF-----tASLdf-----
Z=2      *          *  **          *  *
1mli    ----HYANYSVFDVpsvEALHDTLMQLpLFPY---MDIEVD-----gLCRHpsasihsddr
    
```

# Frequent Fold Types



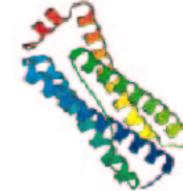
(141) 1hdcA:1  
alpha/beta domain



(85) 1mfA:3  
immunoglobulin fold



(63) 1ceo:2  
TIM barrel



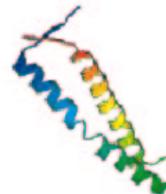
(43) 1bcfA:1  
helical bundle



(36) 2pii:2  
alpha/beta-meander



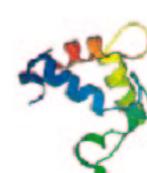
(33) 1vdfA:1  
single helix



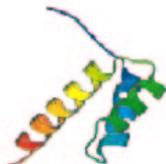
(27) 1grj:2  
coiled coil



(25) 1bbt2:1  
beta-meander



(19) 1rro:2  
EF-hand



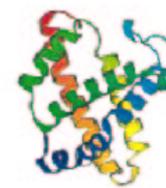
(18) 1octC:3  
HTH-motif



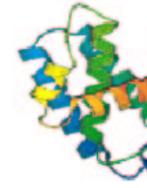
(18) 1prtF:1  
OB-fold



(17) 3grs:2  
FAD/NAD binding domain



(14) 1mbd:1  
globin fold



(13) 1vin:3  
cyclin fold



(13) 1aozA:15  
blue copper protein



(13) 1lef:17  
periplasmic binding protein



(12) 1eelA:3  
lectin fold



(12) 1epaA:1  
lipocalin fold



(12) 2arcA:4  
beta-roll



(12) 2yhxc:3  
actin fold

# Protein Structure Prediction

- **Holy Grail** of bioinformatics
- **Protein Structure Initiative** to determine a set of protein structures that span protein structure space sufficiently well. **WHY?**
  - Number of folds in natural proteins is limited. Thus a newly discovered proteins should be within modeling distance of some protein in set.
- **CASP**: Critical Assessment of techniques for structure prediction
  - To stimulate work in this difficult field

# PSP Methods

- *homology*-based modeling
- methods based on *fold recognition*
  - *Threading* methods
- *ab initio* methods
  - From first principles
  - With the help of databases

# ROSETTA

- ❑ Best method for PSP
- ❑ As proteins fold, a large number of partially folded, low-energy conformations are formed, and that local structures combine to form more global structures with minimum energy.
- ❑ Build a database of known structures (I-sites) of short sequences (3-15 residues).
- ❑ Monte Carlo simulation assembling possible substructures and computing energy

# Threading Methods

□ See p471, Mount

● [http://www.bioinformaticsonline.org/links/ch\\_10\\_t\\_7.html](http://www.bioinformaticsonline.org/links/ch_10_t_7.html)

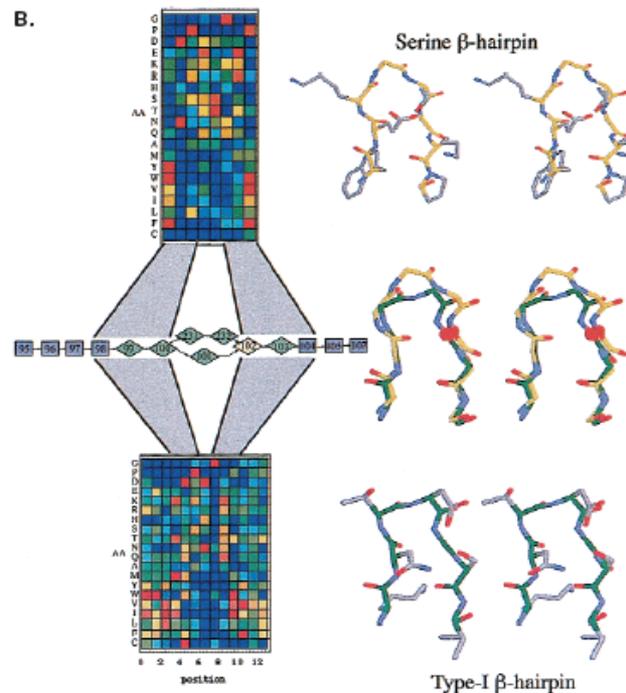


FIGURE 10.30. A hidden Markov model (discrete state-space model) of protein three-dimensional structure. (B) HMM called HMMSTR based on I-sites, 3- to 15-amino-acid patterns that are associated with three-dimensional structural features. The two matrices with colored squares represent alignment of sets of patterns that are found to be associated with a structure, in this case the hairpin turns shown on the right. Each column in the table corresponds to the amino acid variation found for one structural position in one of the turns. (*Blue* side chains) Conserved nonpolar residues; (*green*) conserved polar residues; (*red*) conserved proline; and (*orange*) conserved glycine. The two hairpins are aligned structurally in the middle structure on the right and the observed variation in the corresponding amino acid positions is represented by the HMM between the matrices on the left. The HMM represents an alignment of the two hairpin structural motifs in three-dimensional space and an alignment of the sequences. A short mismatch in the turn is represented by splitting the model into two branches. The shaped icons represent states, each of which represents a structure and a sequence position. Each state contains probability distributions about the sequence and structural attributes of a single position in the motif, including the probability of observing a particular amino acid, secondary structure,  $\Phi$ - $\Psi$  backbone angles, and structural context, e.g., location of  $\beta$  strand in a  $\beta$  sheet. Rectangles are predominantly  $\beta$ -strand states, and diamonds are predominantly turns. The color of the icon indicates a sequence preference as follows: (*blue*) hydrophobic; (*green*) polar; and (*yellow*) glycine. Numbers in icons are arbitrary identification numbers for the HMM states. There is a transition probability of moving from each state in the model to the next, as in HMMs that represent *msa*'s. This model is a small component of the main HMMSTR model that represents a merging of the entire I-sites library. Three different models, designated  $\lambda^P$ ,  $\lambda^C$ , and  $\lambda^R$ , are included in HMMSTR, which differ in details as to how the alignment of the I-sites was obtained to design the branching patterns (topology) of the model and which structural data were used to train the model. HMMSTR may be used for a variety of different predictions, including secondary structure prediction, structural context prediction, and  $\Phi$ - $\Psi$  dihedral angle prediction. Predictions are made by aligning the model with a sequence, finding if there is a high-scoring alignment, and deciphering the highest-scoring path through the model. The HMMSTR program may be downloaded or used on a server that can be readily located by a Web search. (B, reprinted, with permission, from Bystroff et al. 2000 [©2000 Elsevier].)

# Modeling Servers

- SwissMODEL
- 3DJigsaw
- CPHModel
- ESyPred3D
- Geno3D
- SDSC1
- Rosetta
- MolIDE
- SCWRL
- PSIPred
- MODELLER
- LOOPY

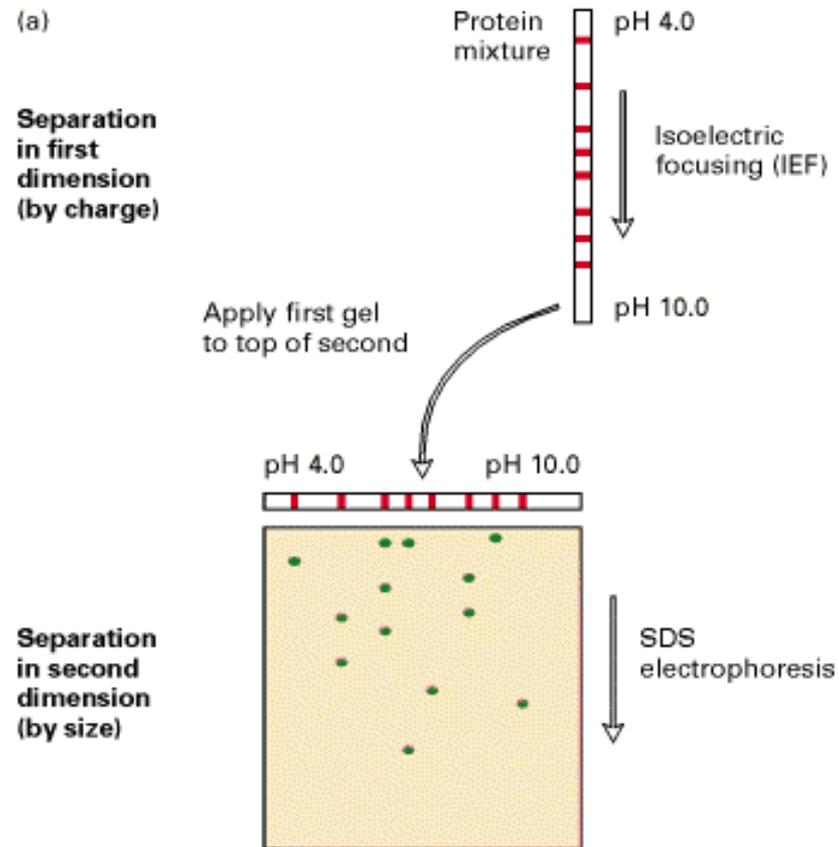
# Genomics

- Study of all genes in a genome
  - All aspects of total gene content
  - Gene Expression
    - Microarray experiments & analysis
    - RNA-Seq

# Proteomics

- Study of all **proteins** in a genome, or comparison of whole genomes.
  - Whole genome annotation & Functional proteomics
  - Whole genome comparison
  - Protein Expression: **2D Gel Electrophoresis**

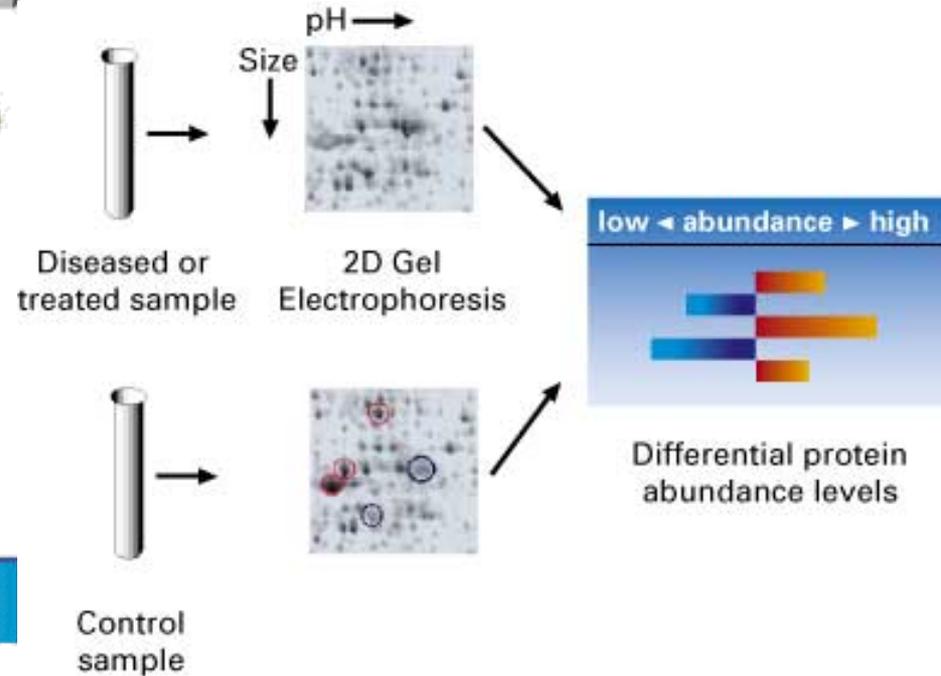
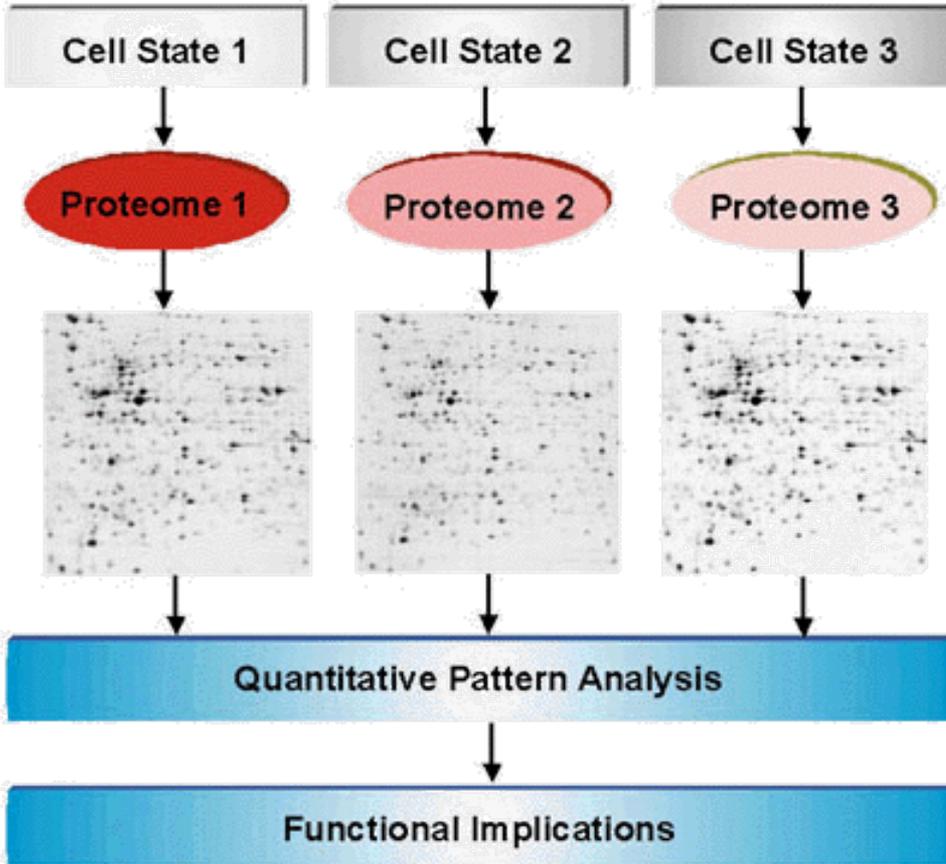
# 2D-Gels



# 2D Gel Electrophoresis



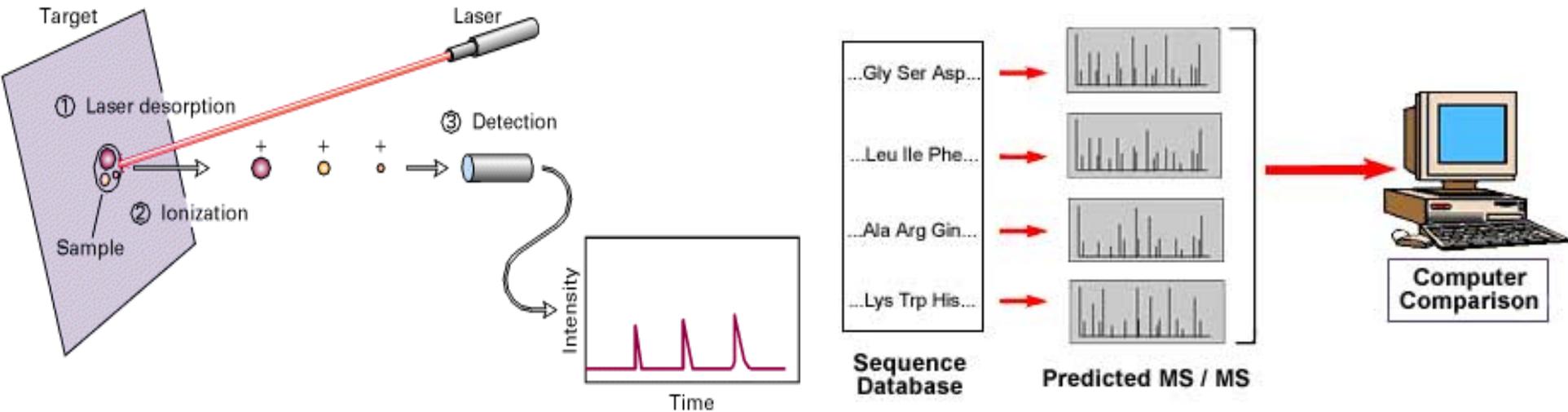
# 2D-gels



Comparing Proteomes For Differences in Protein Expression

Comparing Different Sample Types For Changes in Protein Levels

# Mass Spectrometry



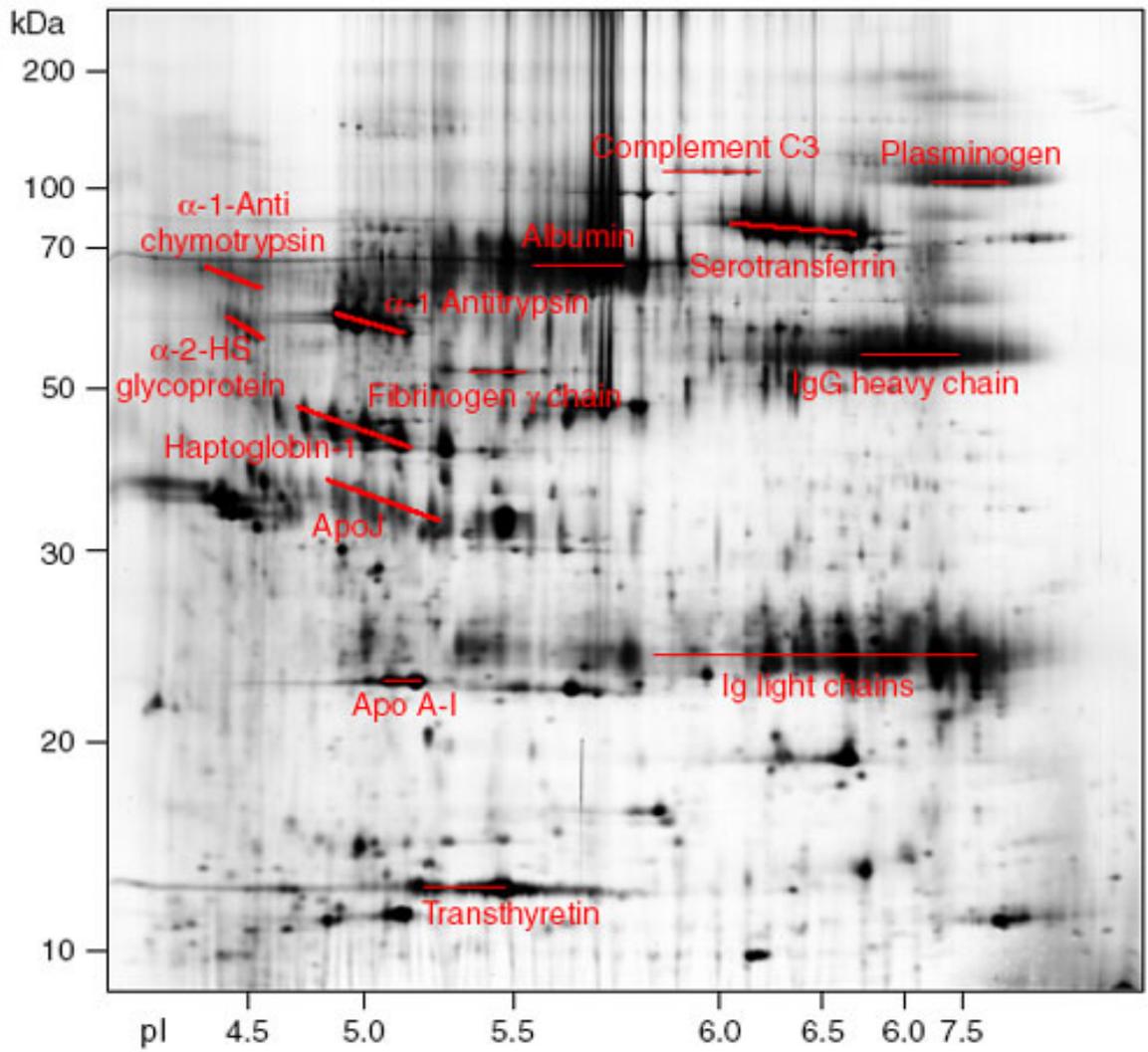
## ❑ Mass measurements By Time-of-Flight

- ❑ Laser ionizes protein
- ❑ Electric field accelerates molecules in sample toward detector
- ❑ Time to detector is inversely proportional to mass of molecule
- ❑ Infer molecular weights of proteins and peptides

# Mass Spectrometry (MS)

## □ Using Peptide Masses to Identify Proteins

- Peptide mass fingerprint is a compilation of molecular weights of peptides
- Use molecular weight of native protein and MS signature to search database for similarly-sized proteins with similar MS maps
- Fairly easy to sequence proteins using MS



TRENDS in Biotechnology

# Other Proteomics Tools

## From ExPASy/SWISS-PROT:

- ❑ **AACompIdent** identify proteins from aa composition  
[Input: aa composition, isoelectric point, mol wt., etc. Output: proteins from DB]
- ❑ **AACompSim** compares proteins aa composition with other proteins
- ❑ **MultIdent** uses mol wt., mass fingerprints, etc. to identify proteins
- ❑ **PeptIdent** compares experimentally determined mass fingerprints with theoretically determined ones for all proteins
- ❑ **FindMod** predicts post-translational modifications based on mass difference between experimental and theoretical mass fingerprints.
- ❑ **PeptideMass** theoretical mass fingerprint for a given protein.
- ❑ **GlycoMod** predicts oligosaccharide modifications from mass difference
- ❑ **TGREASE** calculates hydrophobicity of protein along its length