

# CAP 5510: Introduction to Bioinformatics

**Giri Narasimhan**

ECS 254; Phone: x3748

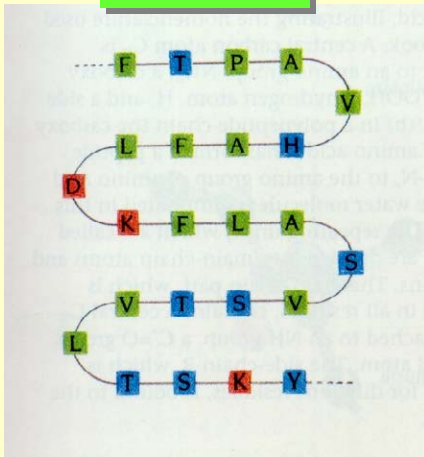
[giri@cis.fiu.edu](mailto:giri@cis.fiu.edu)

[www.cis.fiu.edu/~giri/teach/BioinfS07.html](http://www.cis.fiu.edu/~giri/teach/BioinfS07.html)

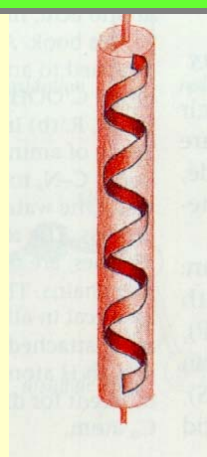
# Protein Structures

- Sequences of amino acid residues
- 20 different amino acids

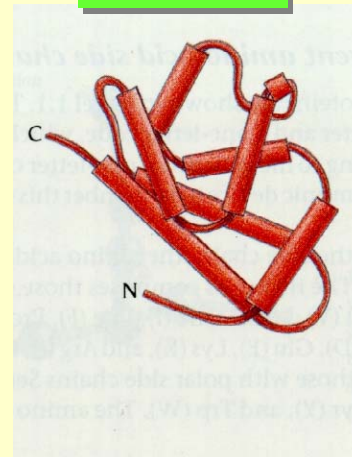
Primary



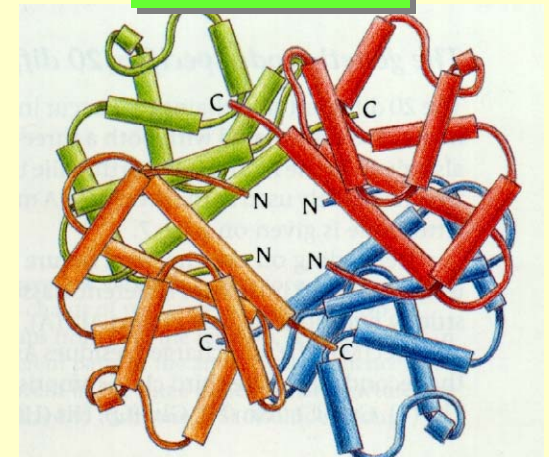
Secondary



Tertiary



Quaternary



# Proteins

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

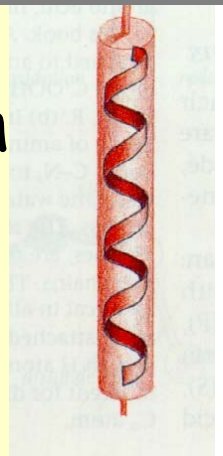
**Flavodoxin:**

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

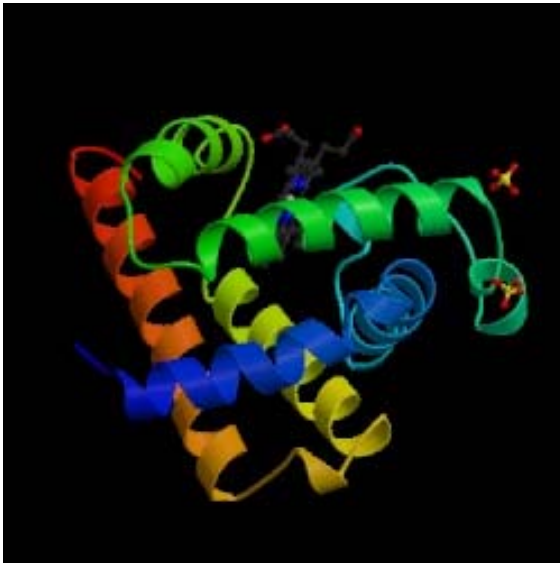
- Main chain with peptide bonds
- Side chains project outward from helix
- Stability provided by H-bonds between CO and NH groups of residues 4 locations away.

## □ $\beta$ -strand

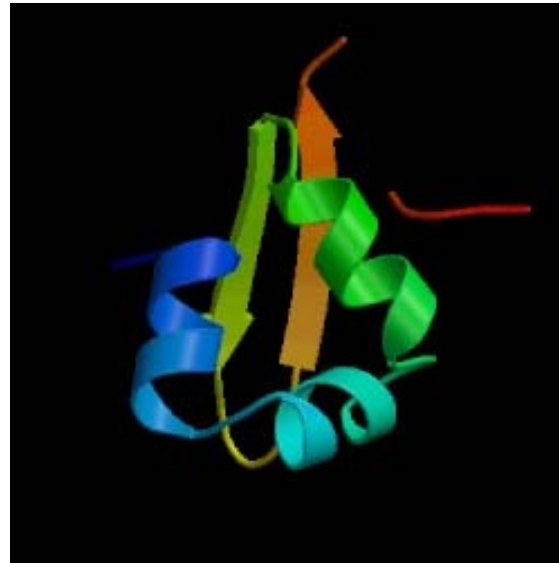
- Stability provided by H-bonds with one or more  $\beta$ -strands, forming  $\beta$ -sheets. Needs a  $\beta$ -turn.

# Proteins

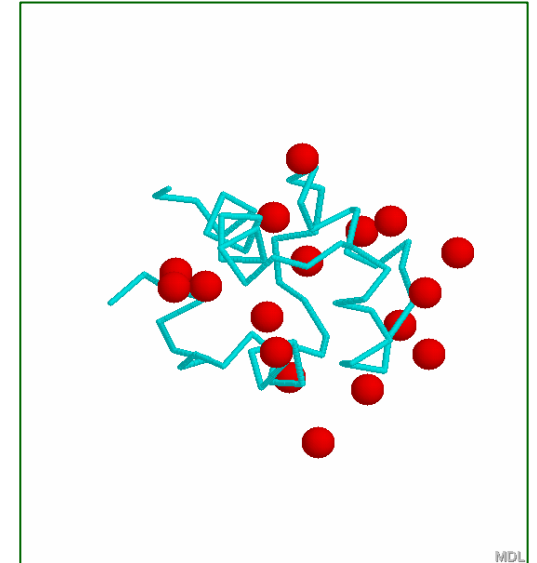
- **Tertiary structures** are formed by packing secondary structural elements into a globular structure.



Myoglobin



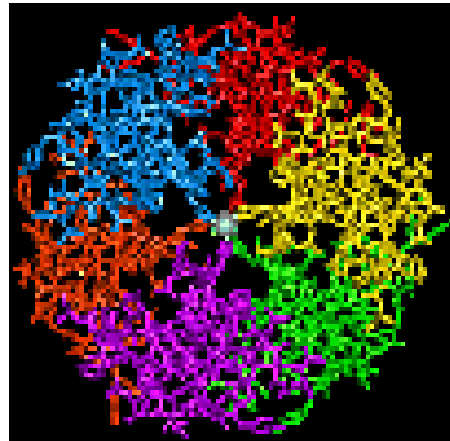
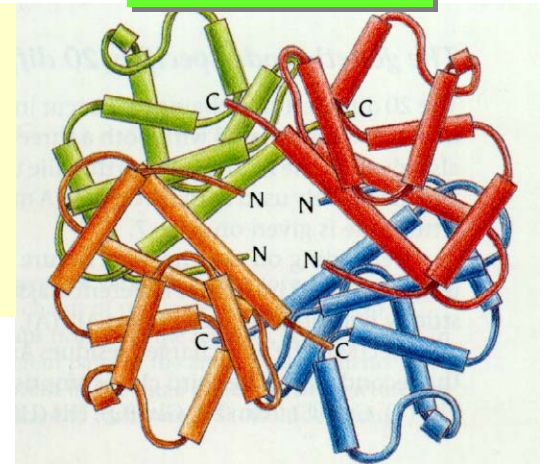
Lambda Cro



# Quaternary Structures in Proteins

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

Quaternary



Insulin Hexamer

# Amino Acid Types

<input type="checkbox"/> <b>Hydrophobic</b>	<b>I, L, M, V, A, F, P</b>
<input type="checkbox"/> <b>Charged</b>	
<input checked="" type="checkbox"/> <b>Basic</b>	<b>K, H, R</b>
<input checked="" type="checkbox"/> <b>Acidic</b>	<b>E, D</b>
<input type="checkbox"/> <b>Polar</b>	<b>S, T, Y, H, C, N, Q, W</b>
<input type="checkbox"/> <b>Small</b>	<b>A, S, T</b>
<input type="checkbox"/> <b>Very Small</b>	<b>A, G</b>
<input type="checkbox"/> <b>Aromatic</b>	<b>F, Y, W</b>

# Structure of a single amino acid

All 3 figures are cartoons of an amino acid residue.

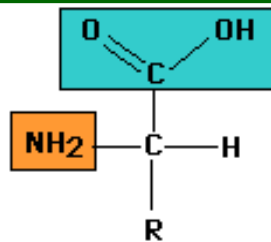
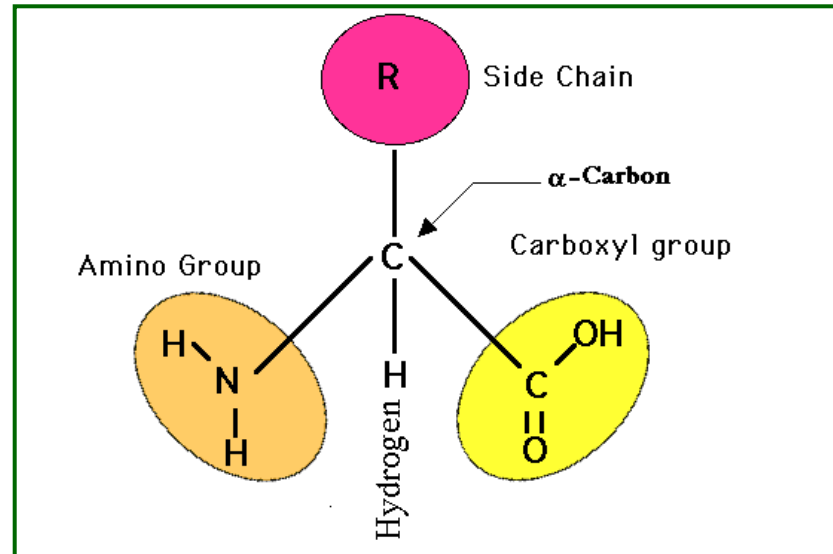
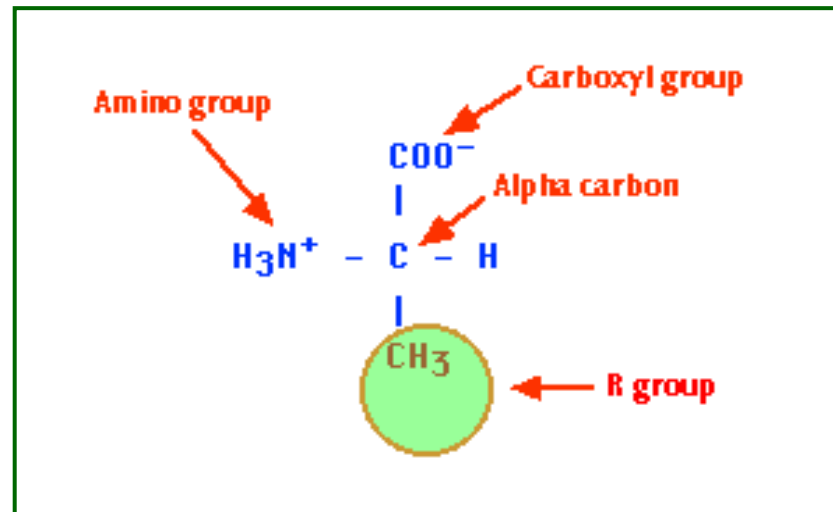
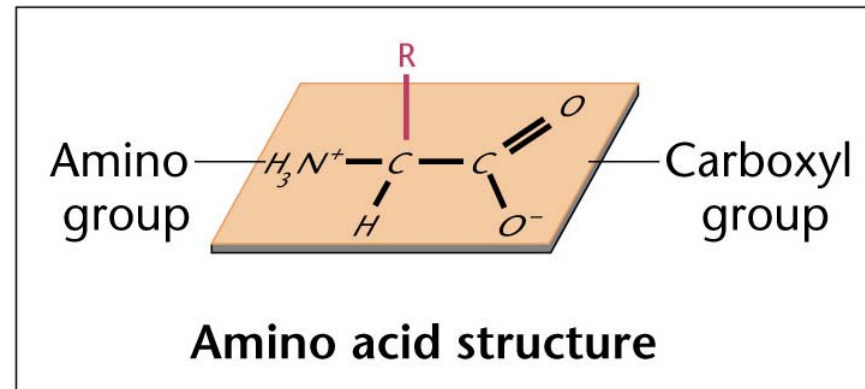


Fig. General formula for an amino acid molecule. "R" represents the variable groups that are attached to this basic molecule to make up the 20 common amino acids

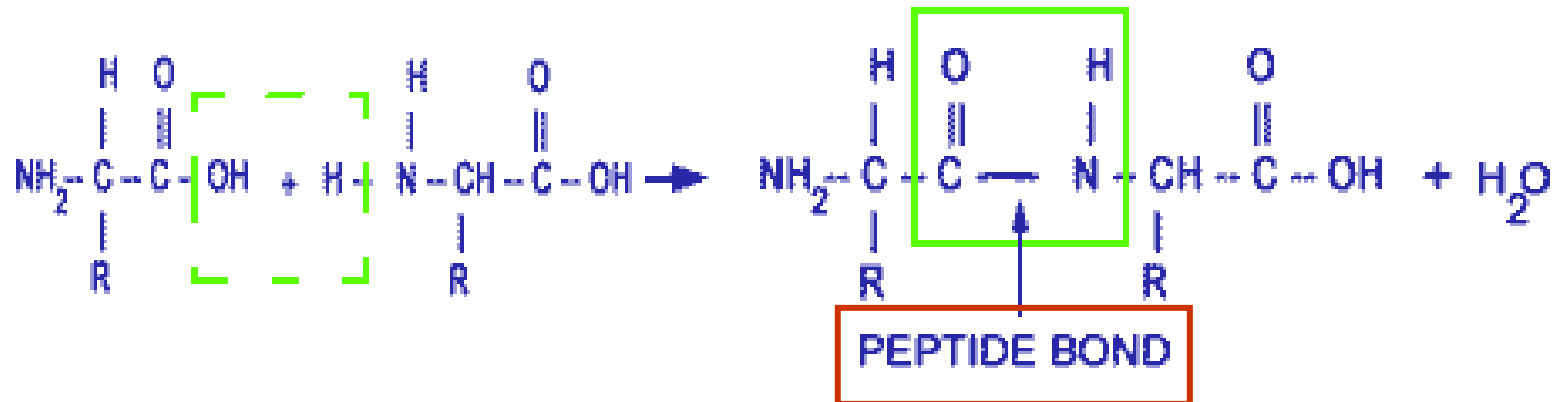




# Structure of a single amino acid

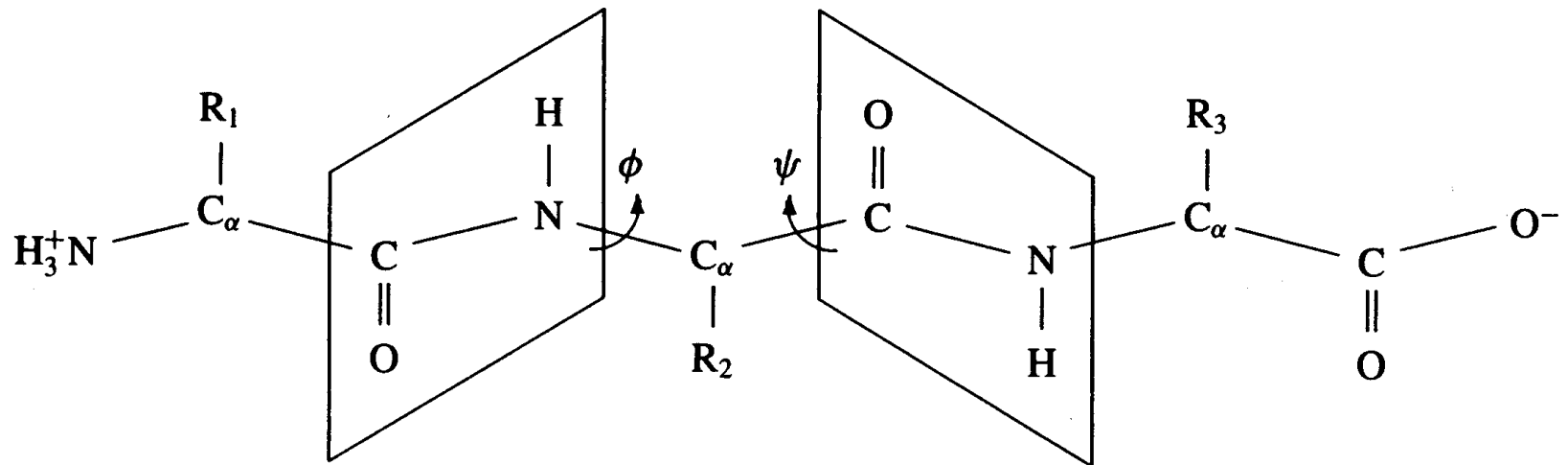


# Chains of amino acids



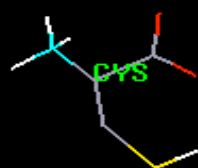
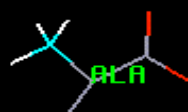
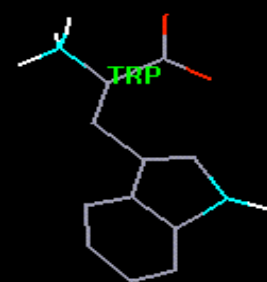
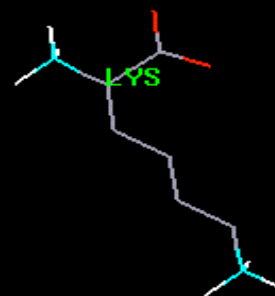
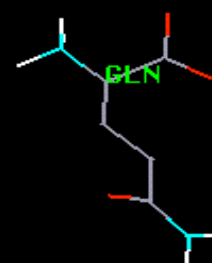
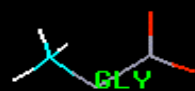
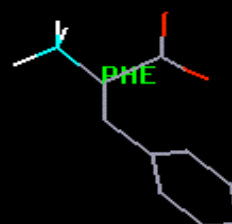
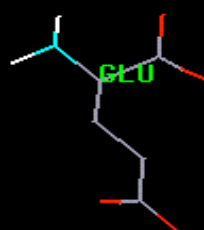
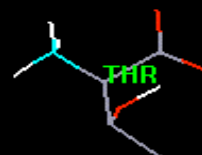
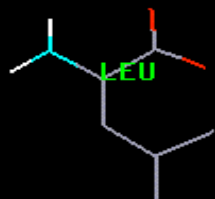
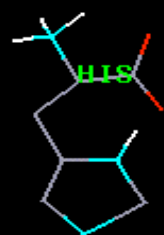
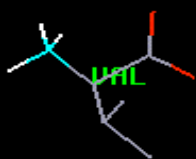
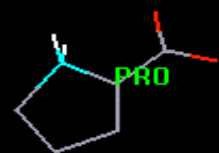
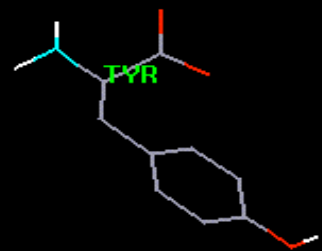
**Amino acids** vs **Amino acid residues**

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

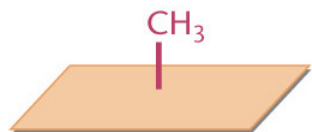


**FIGURE 1.2**

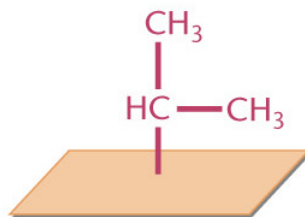
*A polypeptide chain. The  $\text{R}_i$  side chains identify the component amino acids. Atoms inside each quadrilateral are on the same plane, which can rotate according to angles  $\phi$  and  $\psi$ .*



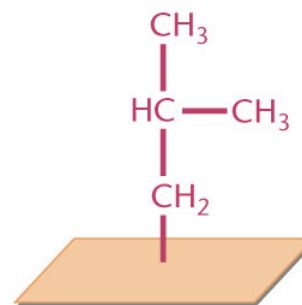
# 1. Nonpolar: Hydrophobic



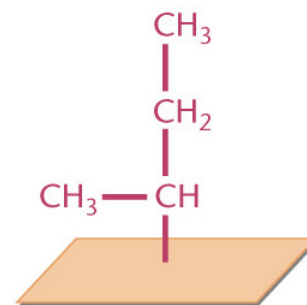
Alanine (ala-A)



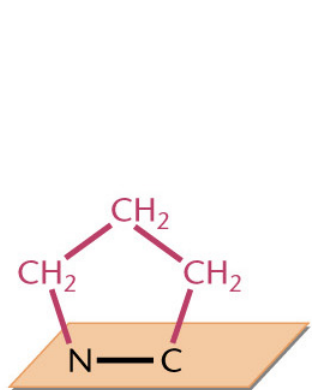
Valine (val-V)



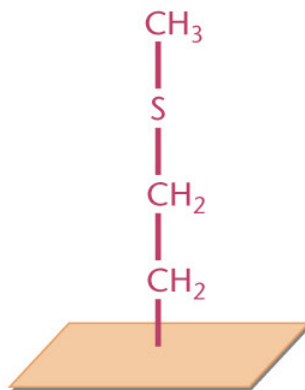
Leucine (leu-L)



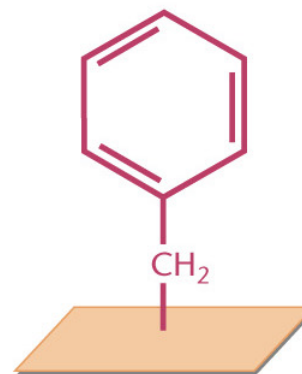
Isoleucine (ile-I)



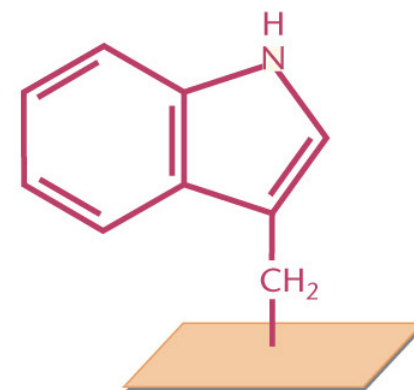
Proline (pro-P)



Methionine (met-M)



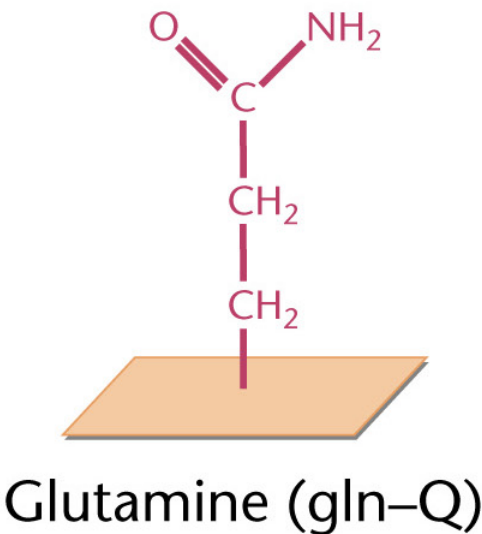
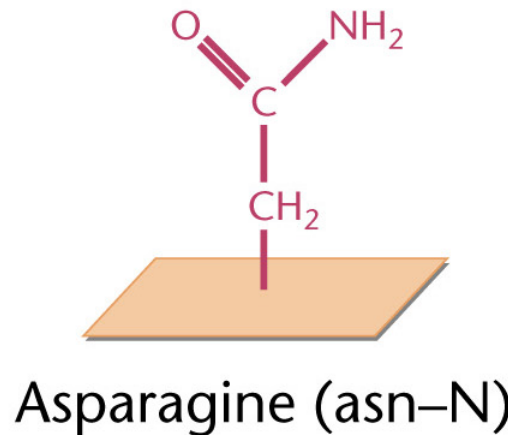
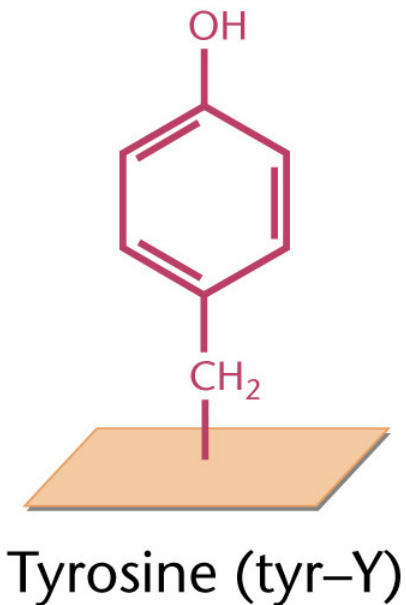
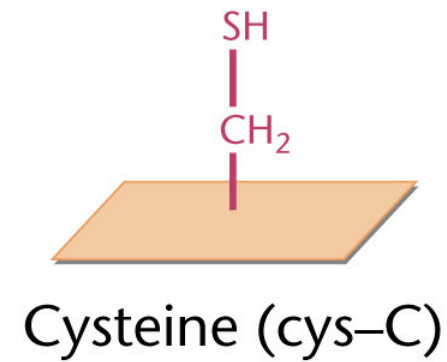
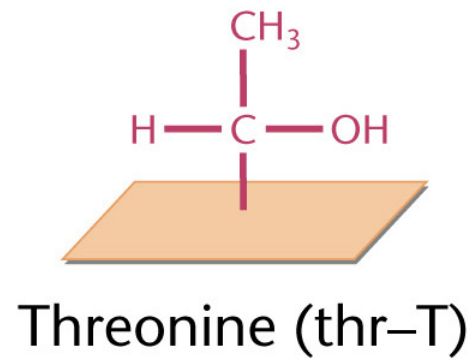
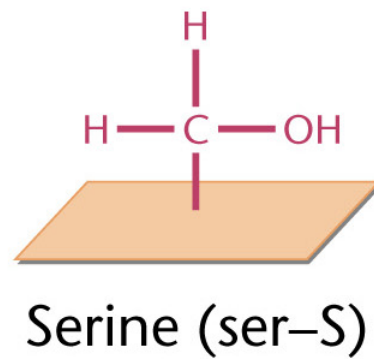
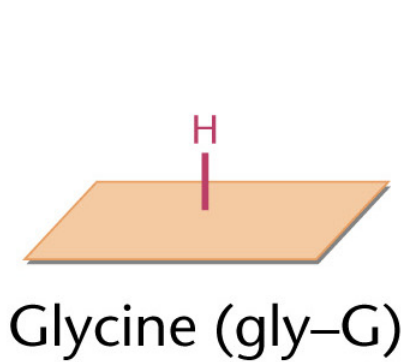
Phenylalanine (phe-F)



Tryptophan (trp-W)

Amino Acid Structures from Klug & Cummings

## 2. Polar: Hydrophilic



Amino Acid Structures from Klug & Cummings

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

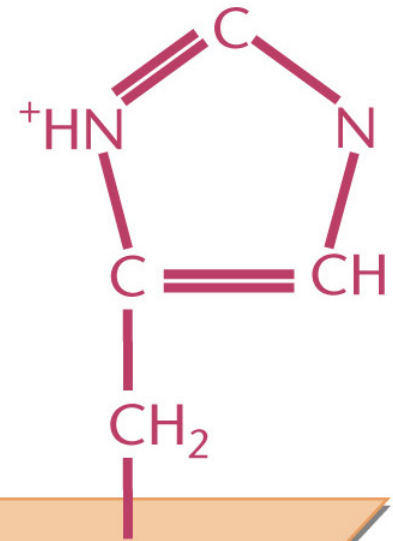
Amino Acid Structures  
from Klug & Cummings



Lysine (lys-K)

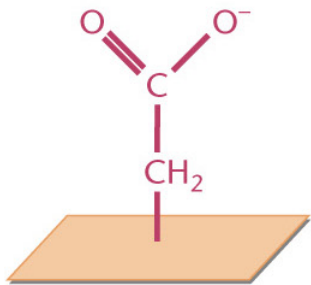


Arginine (arg-R)

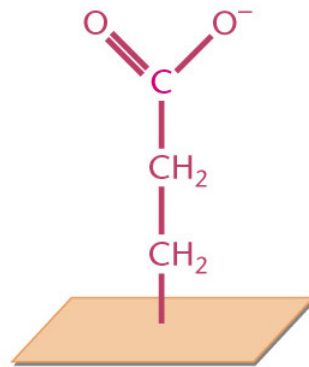


Histidine (his-H)

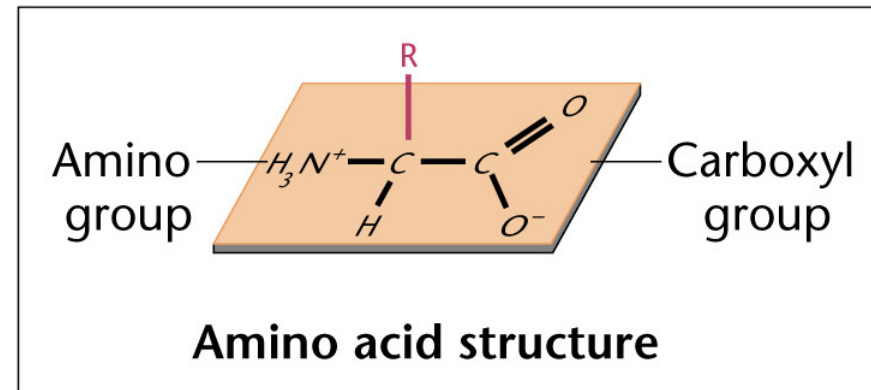
#### 4. Polar: negatively charged (acidic)



Aspartic acid (asp-D)



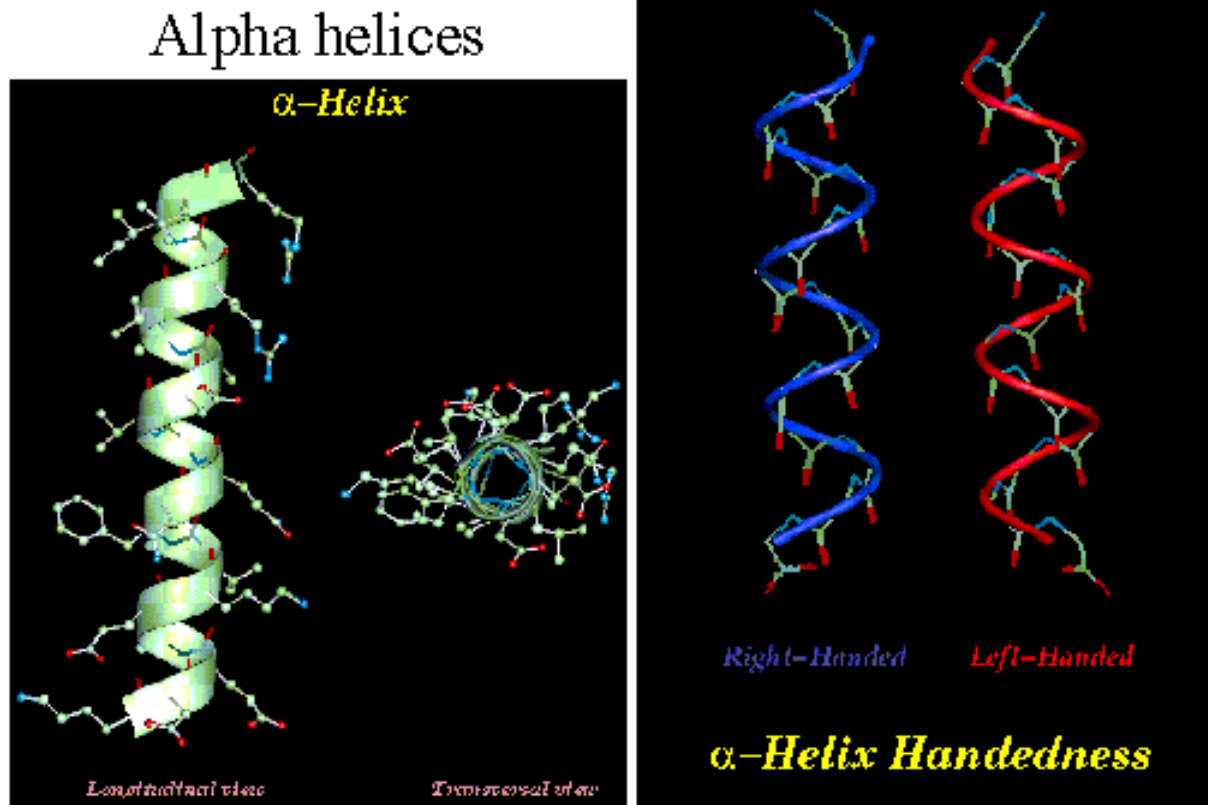
Glutamic acid (glu-E)



Amino Acid Structures from Klug & Cummings

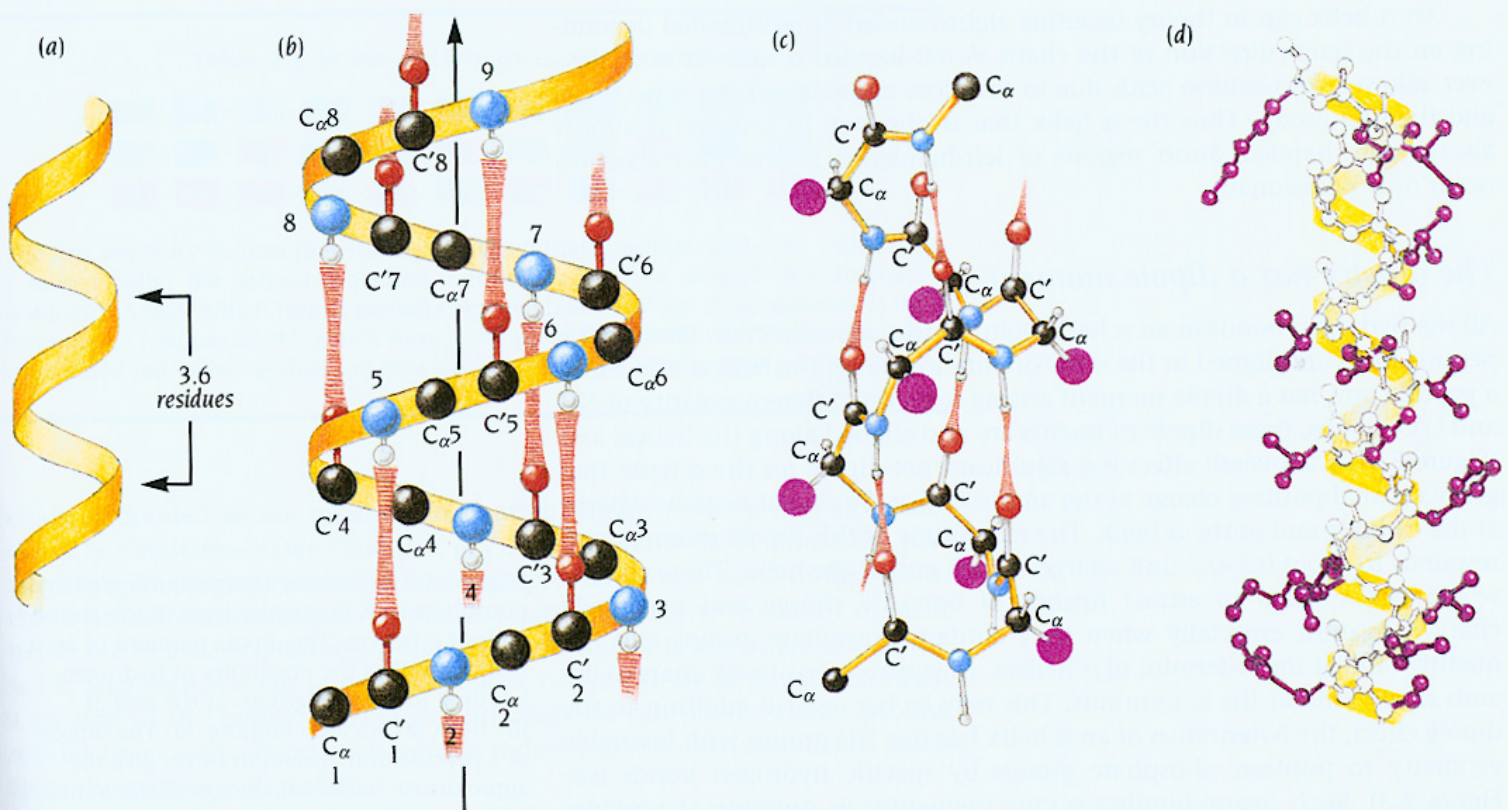


# Alpha helices



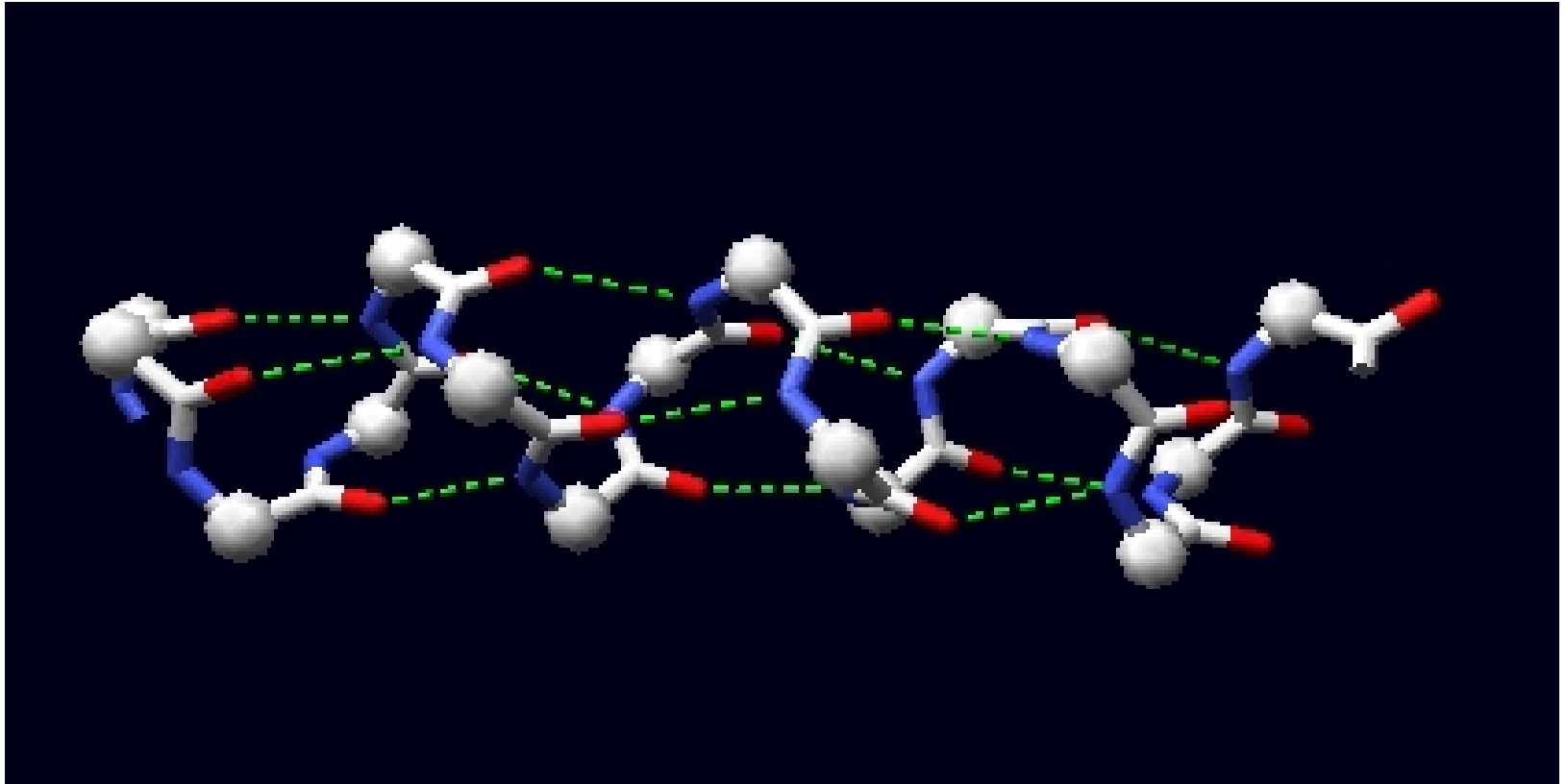
(c) David Gilbert, Aik Choon Tan, Gilleain Torrance and Mallika Veeramalai 2002

16



**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.

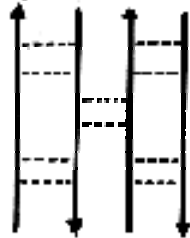
# Alpha Helix



# Beta Sheets

## Beta sheet

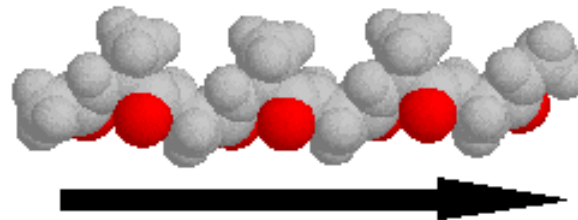
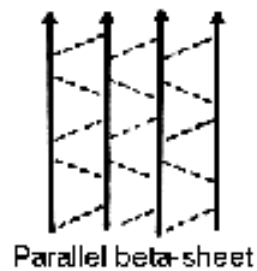
Antiparallel beta-sheet



The beta-hairpin turn.



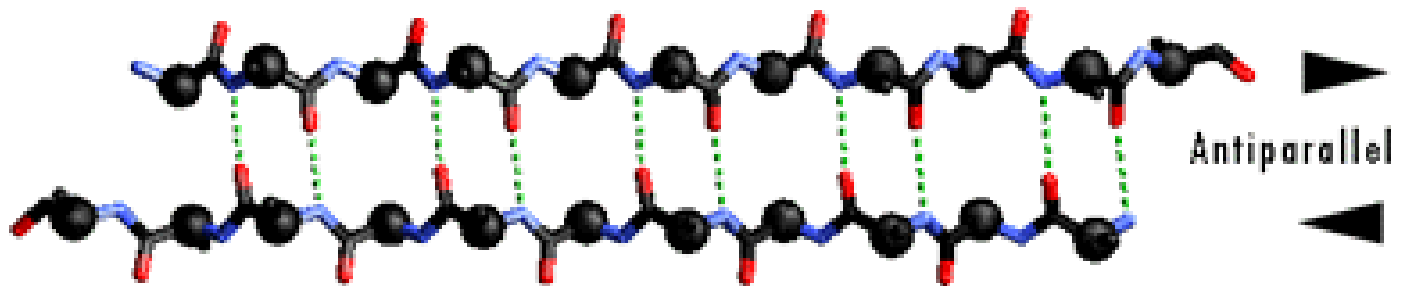
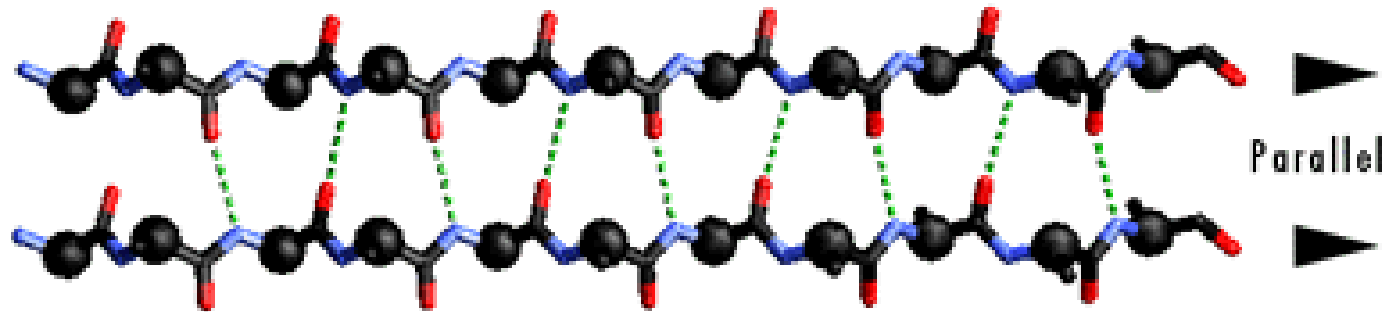
The dashed lines indicate main chain hydrogen bonds.



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

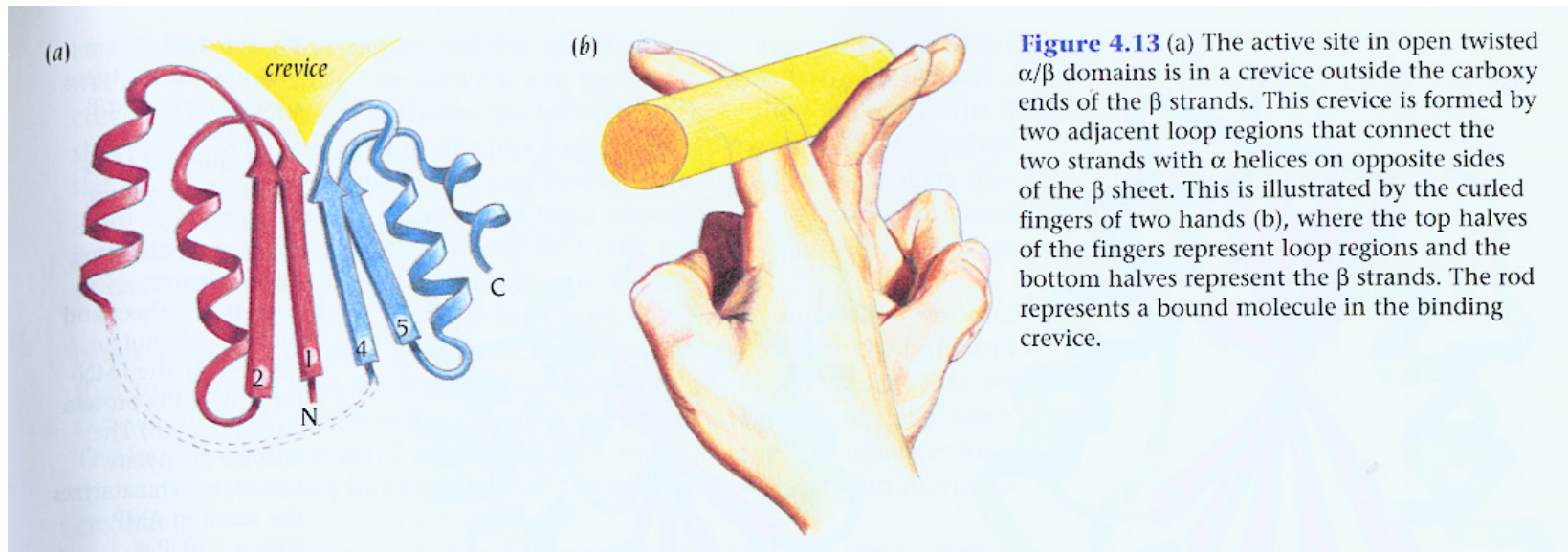
17

# Beta Sheets

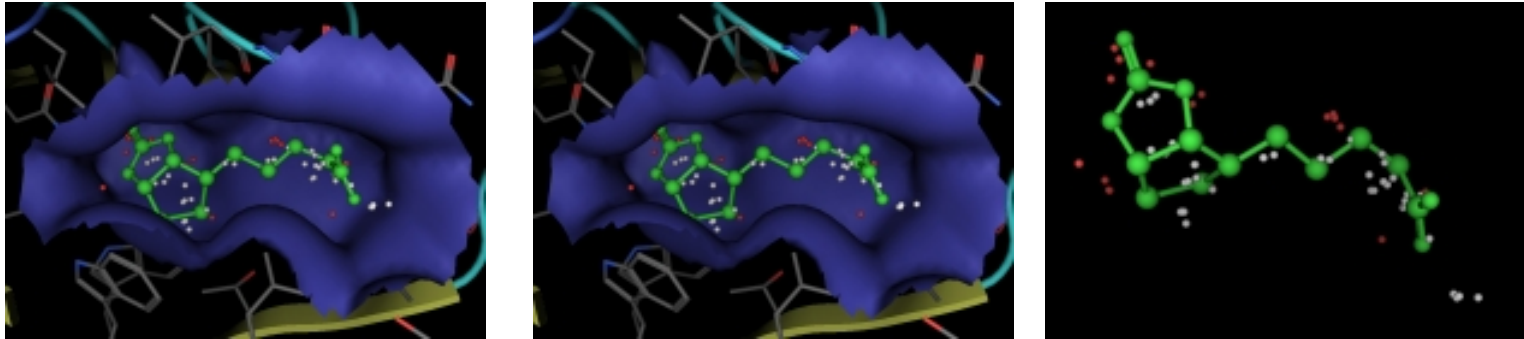


# Active Sites

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.



# Secondary Structure Prediction Software

254



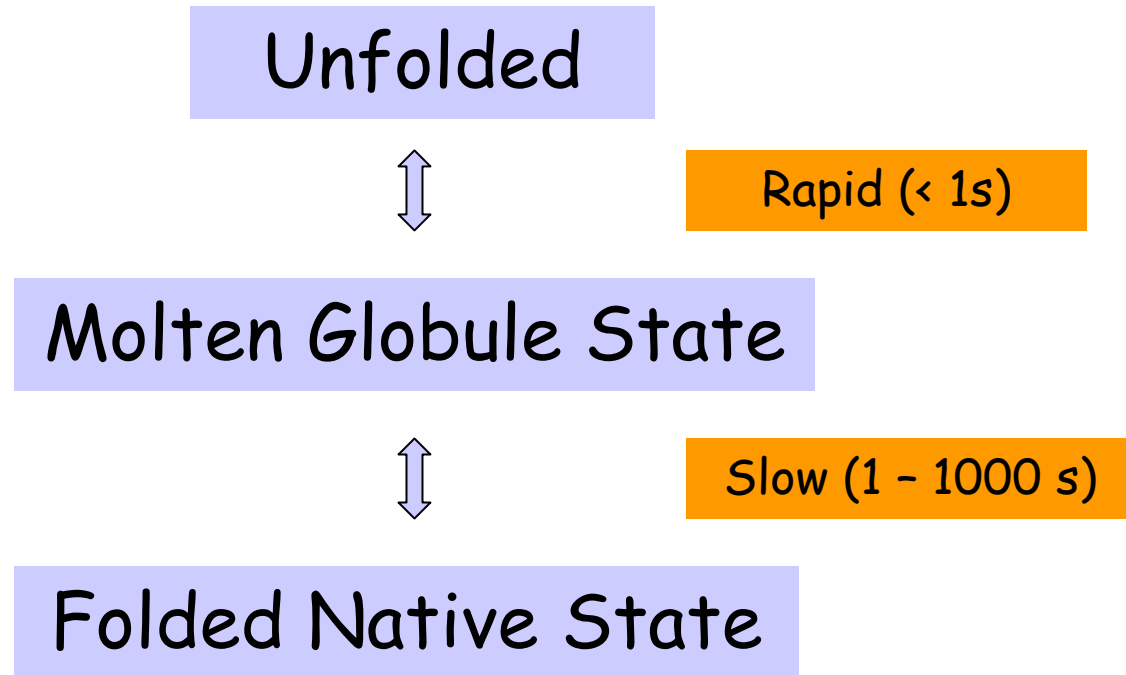
**Figure 11.3** Comparison of secondary structure predictions by various methods. The sequence of flavodoxin, an  $\alpha/\beta$  protein, was used as the query and is shown on the first line of the alignment. For each prediction, H denotes an  $\alpha$  helix, E a  $\beta$  strand, T a  $\beta$  turn; 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).



# 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 \[.\]](#)

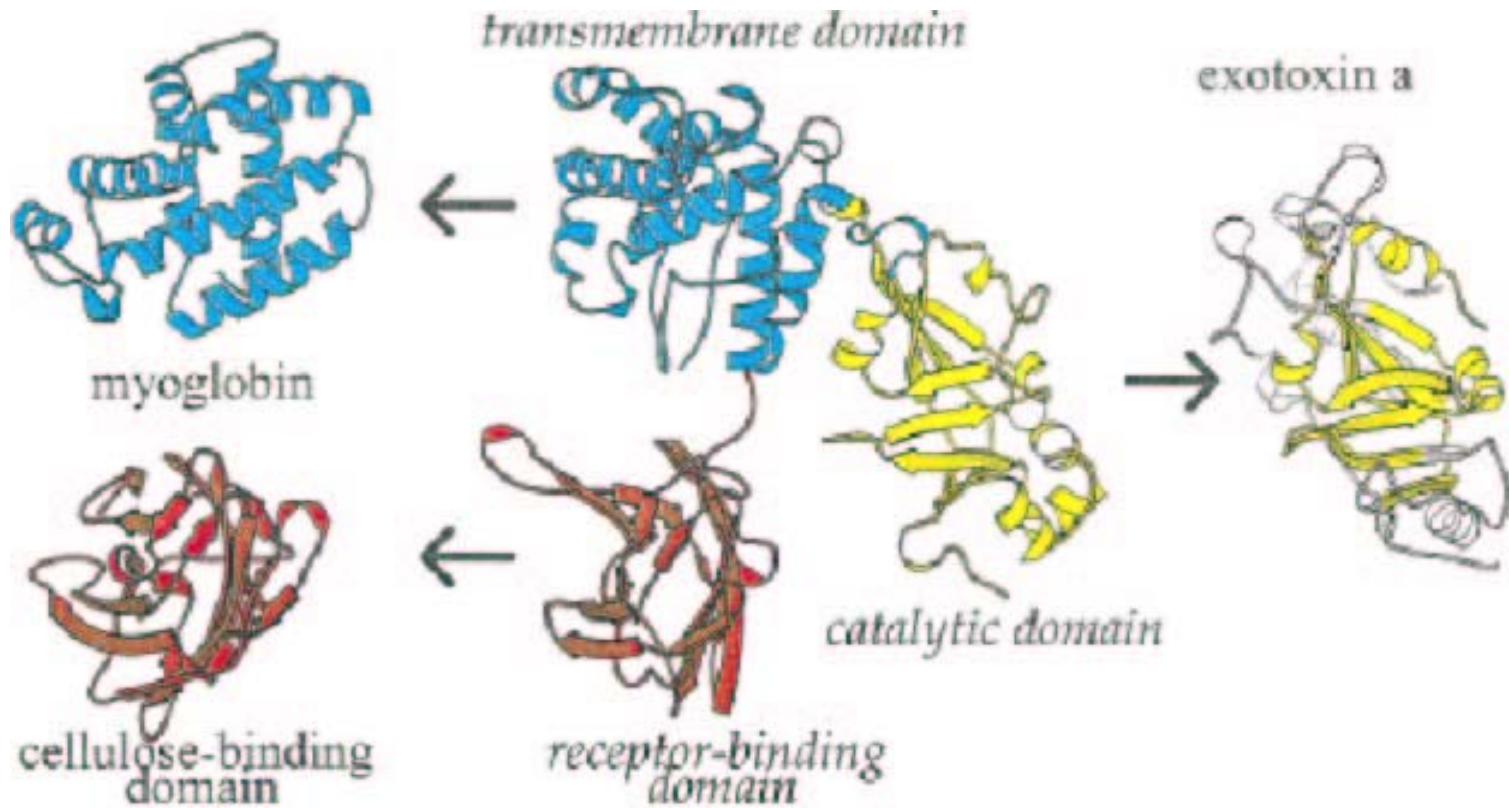
# Protein Folding



□ How to find minimum energy configuration?

# 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.

# 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. At the top, the document title is "SCOP: Family: Interleukin 8-like" and the URL is "http://scop.mrc-lmb.cam.ac.uk/scop/data/scop.0.004". The "Structural Classification of Proteins" section shows a lineage tree for the Interleukin 8-like family. The tree lists the following entries:

- 1. Root: [scop](#)
- 2. Class: [Alpha](#)
- 3. Fold: [Interleu](#)
- 4. Superfamily: [Interleukin 8-like](#)
- 5. Family: [Interleukin 8-like](#)

The "Proteins:" section lists the following entries:

- 1. Interleukin-8
  - 1. [human \(Homo sapiens\) \(8\)](#)
    - 1. [1JH8](#)
    - 2. [1JH8](#)
      - 1. chain a
      - 2. chain b
    - 3. [2JH8](#)
      - 1. chain a
      - 2. chain b
- 2. Platelet factor 4
  - 1. [bovine \(Bos taurus\) \(1\)](#)
    - 1. [1JH1](#)
      - 1. chain a
      - 2. chain b
      - 3. chain c
      - 4. chain d
- 3. Macrophage inflammatory protein 1beta has different generation mode
  - 1. [human \(Homo sapiens\) \(2\)](#)
    - 1. [1JH9](#)
      - 1. chain a
      - 2. chain b
    - 2. [1JH9](#)
      - 1. chain a
      - 2. chain b

**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

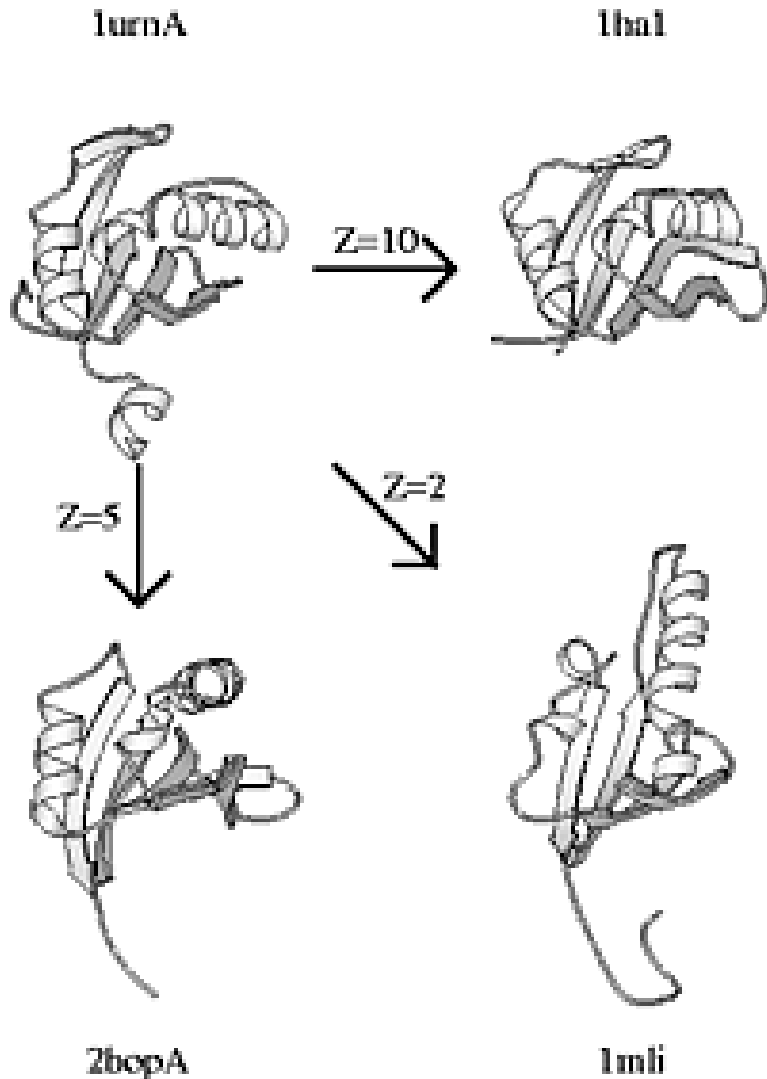
- ❑ **MMDB** contains groups of structurally related proteins
- ❑ **SARF** structurally similar proteins using secondary structure elements
- ❑ **VAST** Structure Neighbors
- ❑ **SSAP** uses double dynamic programming to structurally align proteins

# 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).

# 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----KKDELKKSLHAIFSRFG---QILDILV-SRS---LKM---
Z=10      *      *      *      *      *      *
1ha1    ahLTVKKIFVGGIKEDT-----EEHHLRDYFEOYG---KIEVIEI-MTDrgsGKK---
Z=5      *
2bopA   ----sCFALIS-GTANO-----vKCYRFRVKKNHRHR-----YENCTTtWFT---Vadnga
Z=2      *
1mli    ---mlFHVKMTVKLpvdmdpakatgkadeKELAQRLgregTWRHLWR-IAG-----

1urnA   ----RGQAFVIFKEV--SSATNALRSMQGFPFYDKPMRIQYAKTSDIIAKM-----
Z=10     **  ***  *      *      *
1ha1     ----RGFAFVTFDDH--DSVDKIVIO-kyHTVNGHNCEVRKAL-----
Z=5      *      *      *      *      *      *
2bopA   erggQAQILITFGSP--SORODFLKHVPLPP----GMNISGF-----tASLdf-----
Z=2      *      *      **      *      *
1mli     ----HYANYSVFDVpsvEALHDTLMQLpLFPY----MDIEVD-----gLCRHpssihsddr
    
```

# 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



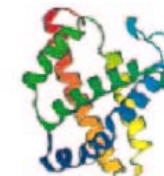
(18) 1oetC:3  
HTH-motif



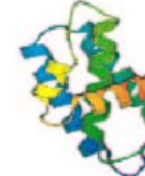
(18) 1ptf:1  
OB-fold



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



(14) 1mbd:1  
globin fold



(13) 1vin:3  
cyclin fold



(13) 1aozA:15  
blue copper protein



(13) 1lcf:17  
periplasmic binding protein



(12) 1eelA:3



(12) 1epaA:1  
lipocalin fold



(12) 2arcA:4  
beta-roll



(12) 2yhx: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)

