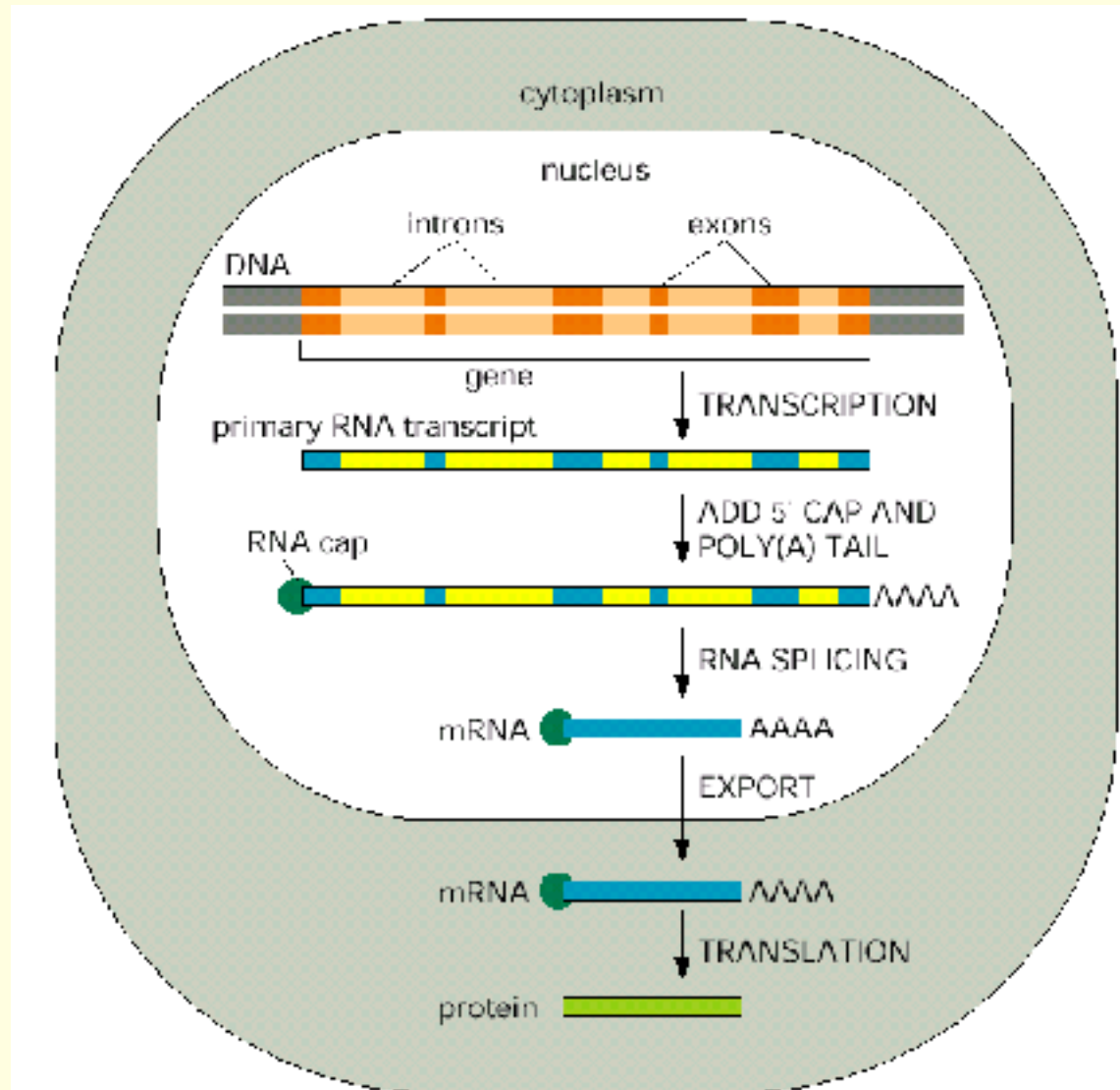
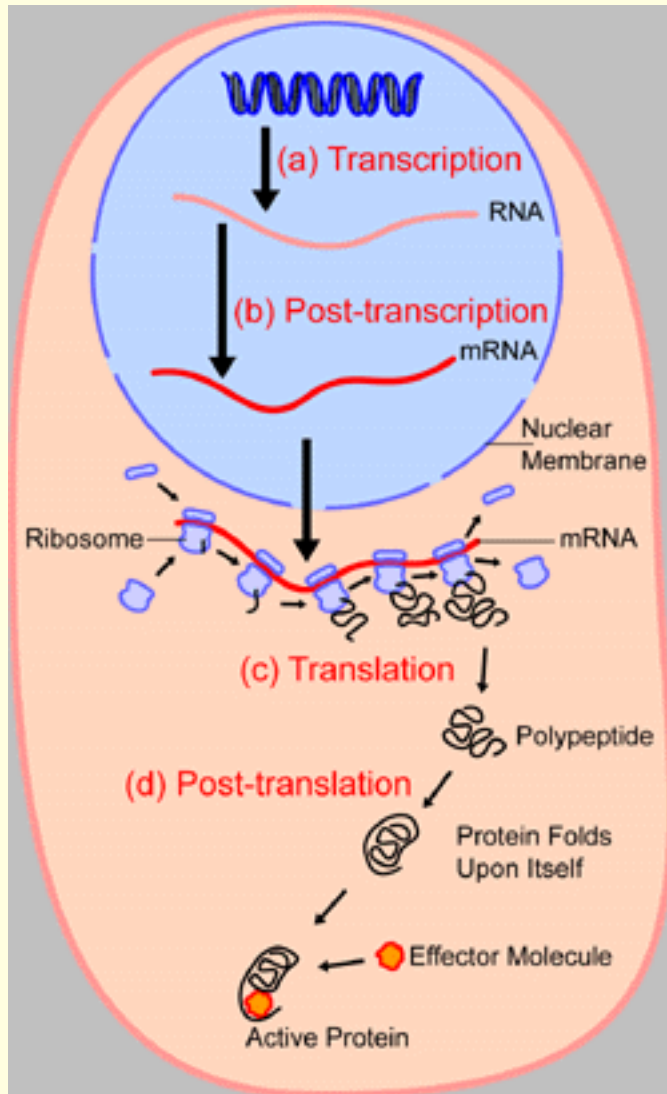
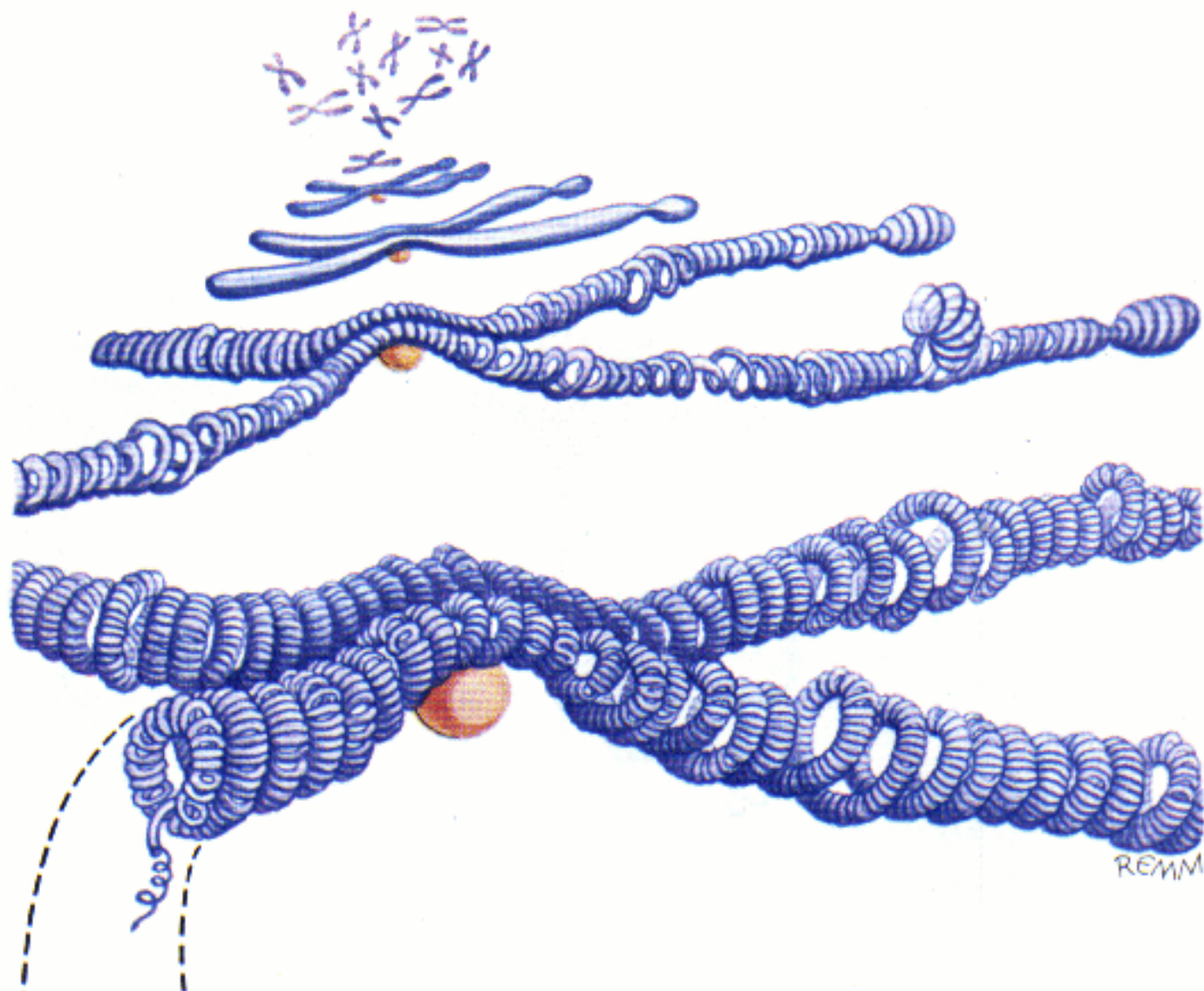


# CMPS 6630: Introduction to Computational Biology and Bioinformatics

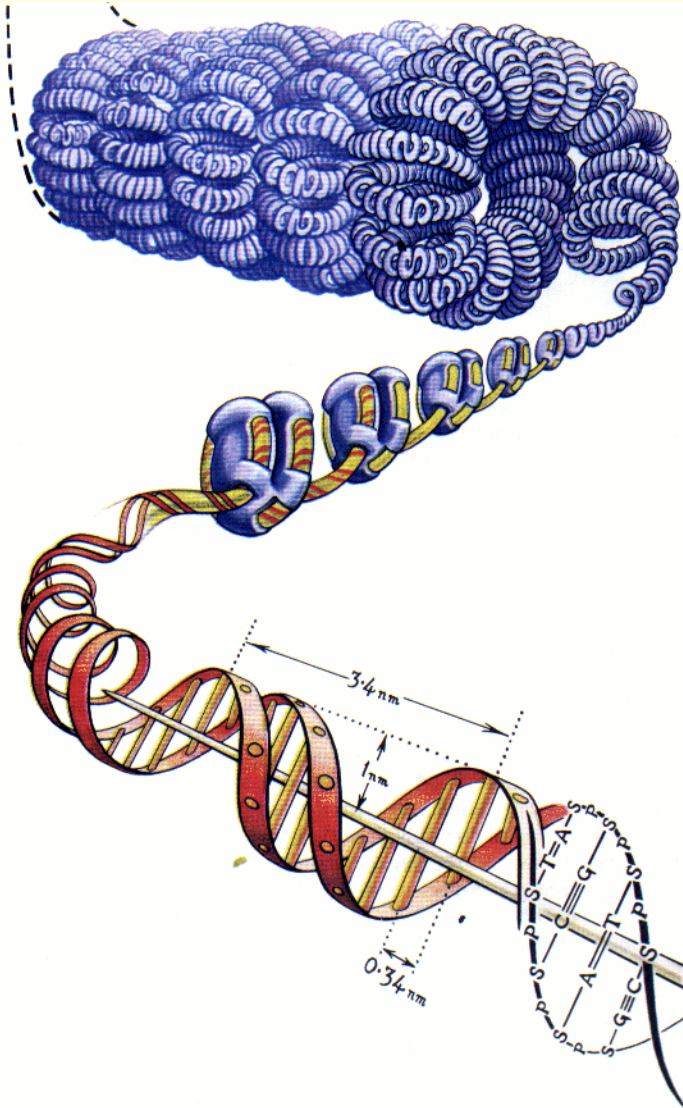
Protein Structure

# Recap: Central “Dogma”





# DNA Structure



Franklin/Watson/Crick showed that DNA has a structure that is stable, and facilitates replication.

Enzymes that bind DNA and RNA must have a “compatible” structure (e.g. ribosomes).

What is the (molecular) nature of this regulation?

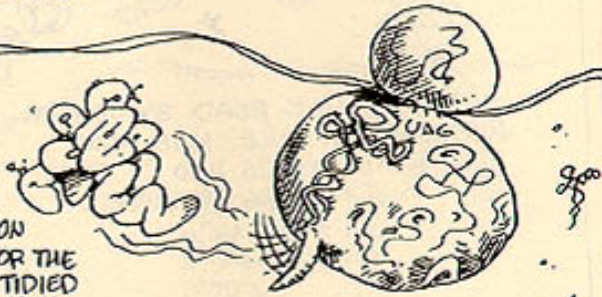
THIS PROCESS CONTINUES UNTIL THE RIBOSOME REACHES ONE OF THE 'STOP' SIGNALS.



IT STOPS BECAUSE THERE IS NO tRNA WITH AN ANTICODON TO MATCH.

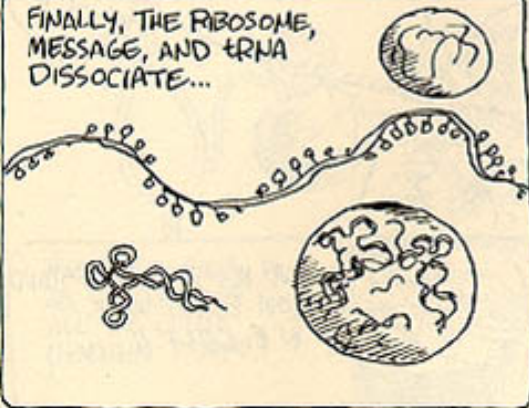


THE COMPLETED PROTEIN IS CLIPPED OFF BY ANOTHER RIBOSOMAL ENZYME.



IT IS ALSO COMMON AT THIS POINT FOR THE PROTEIN TO BE TIDIED UP IN VARIOUS WAYS.

FINALLY, THE RIBOSOME, MESSAGE, AND tRNA DISSOCIATE...



...AND THE NEW MACROMOLECULE GOES OFF TO DO ITS JOB: STRUCTURE, ENZYME, OR WHATEVER...

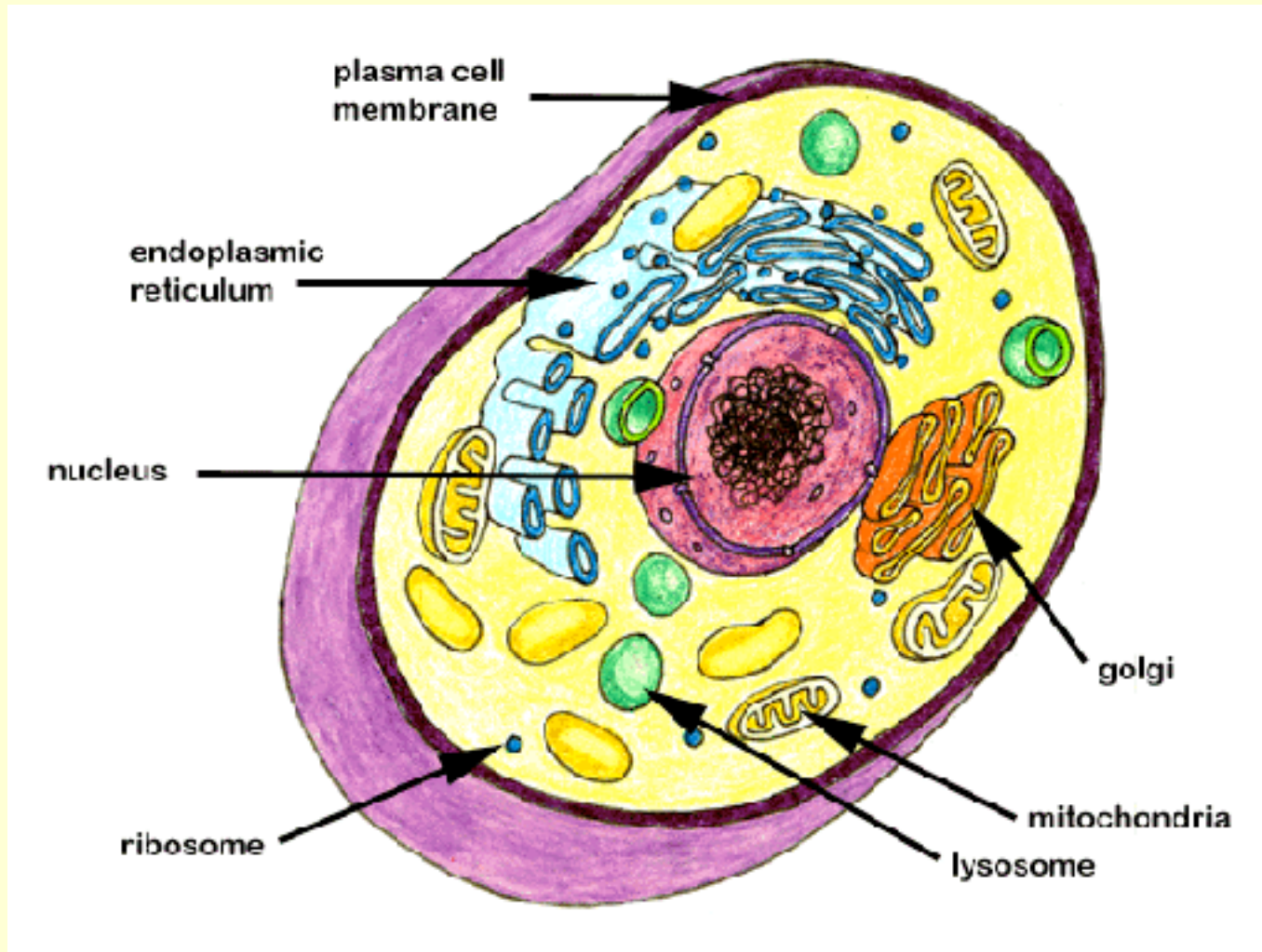


Proteins have 3D shape that is determined by a sequence of amino acids.

Structure is function (e.g. ribosome, hemoglobin, transcription factors).

# Cell Composition

---



# Cell Composition

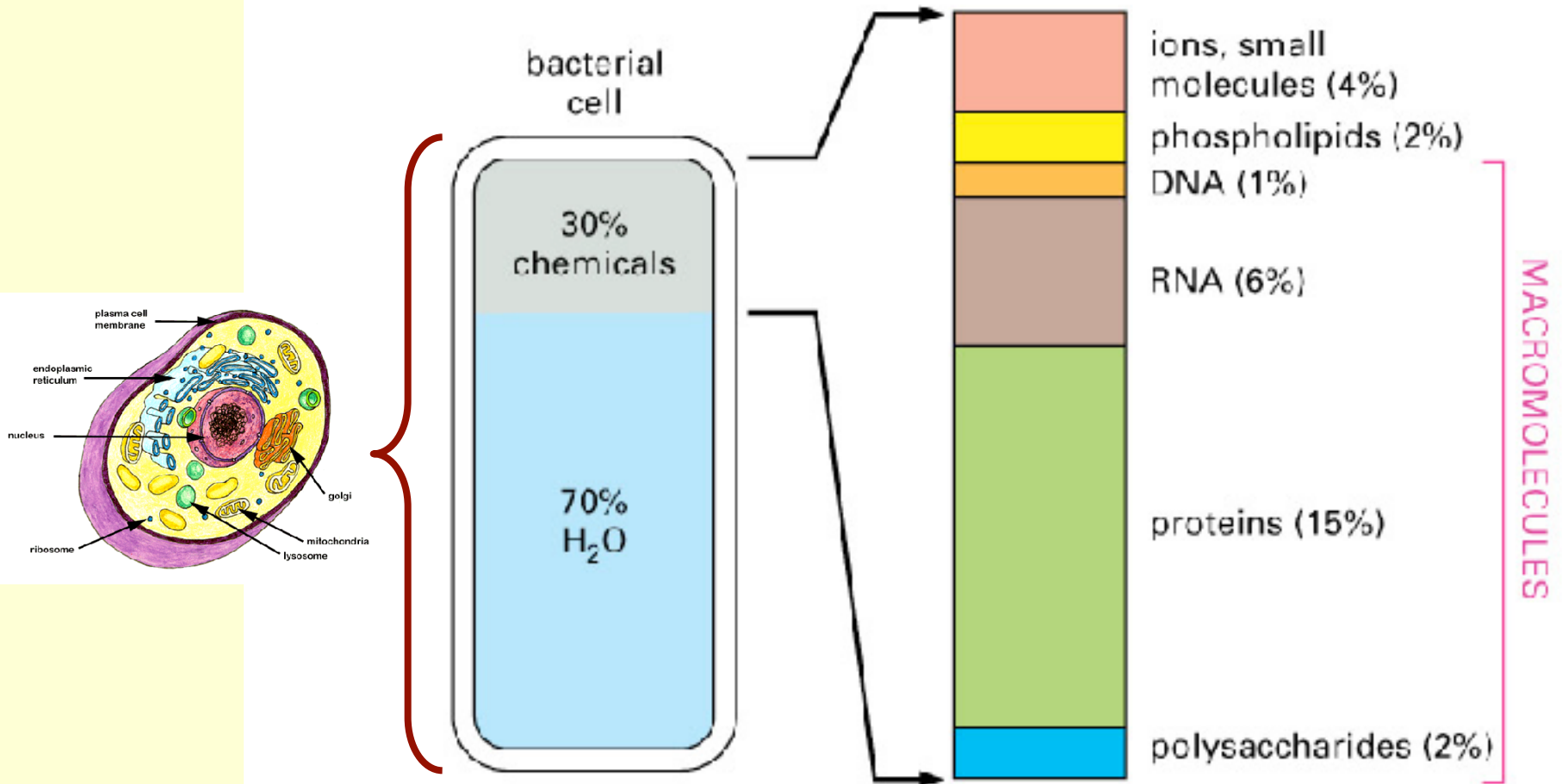
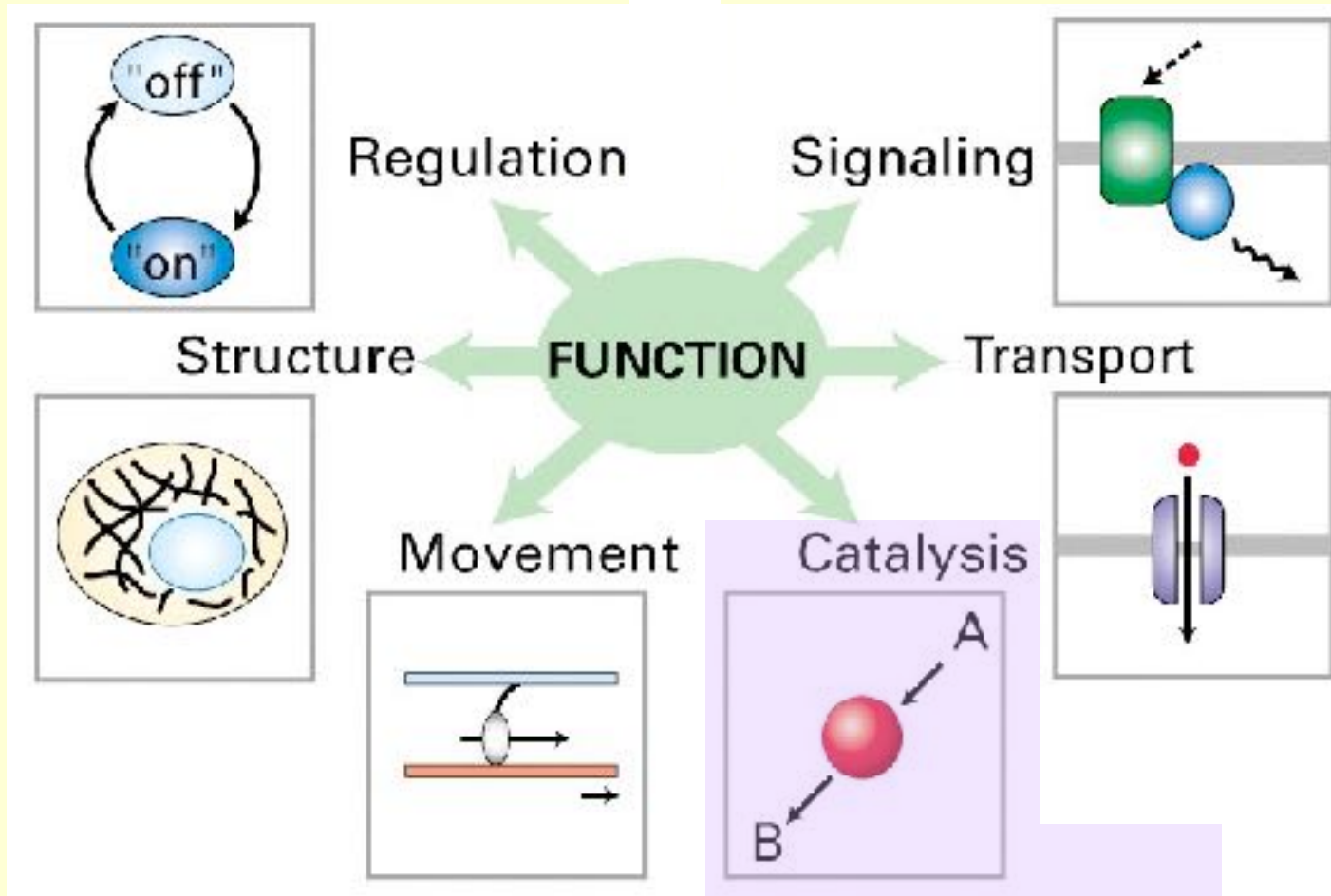


Figure 2-26 Essential Cell Biology, 2/e. (© 2004 Garland Science)

# Protein Functions

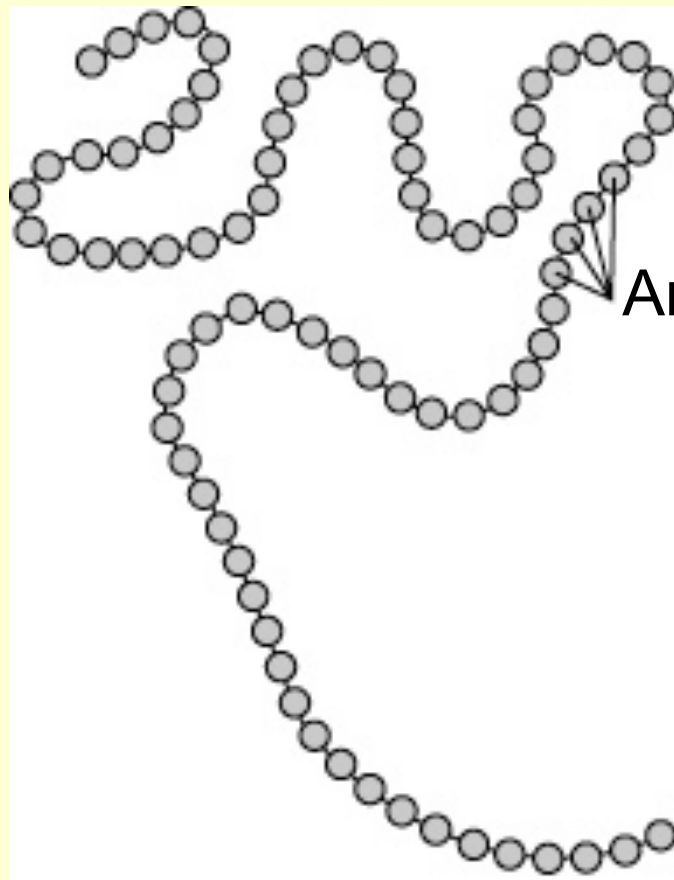
Activity  
Specificity  
Regulation



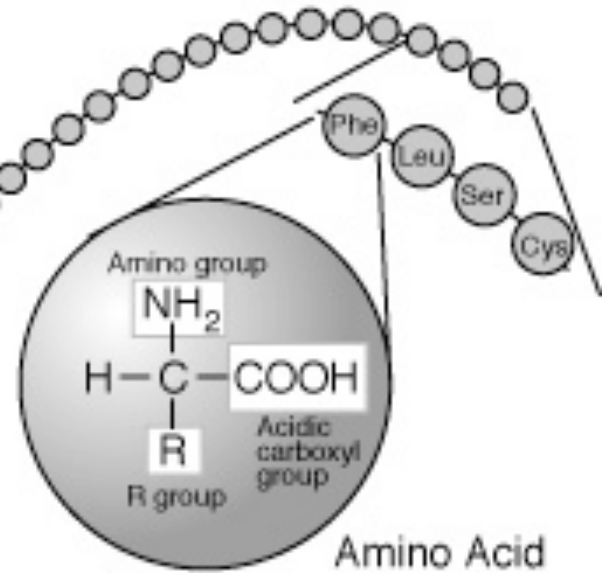
**Enzymes**



Primary protein structure  
is sequence of a chain of amino acids



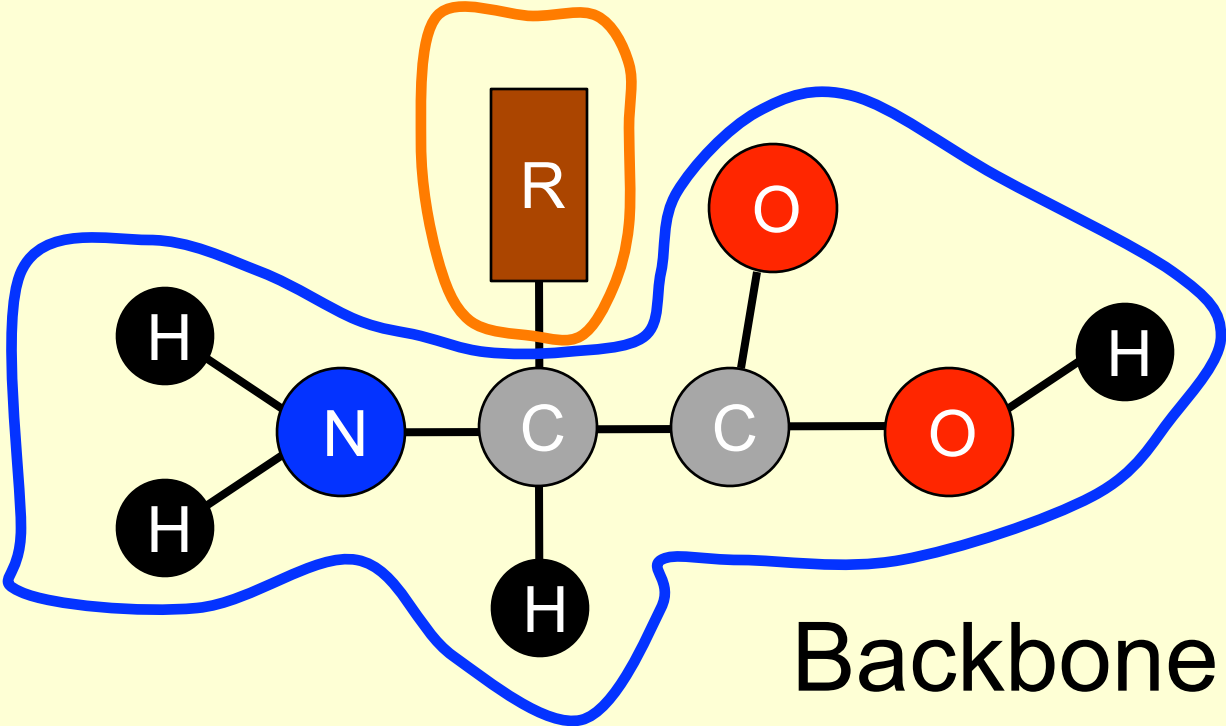
Amino Acids (or Residues)



# Amino Acid

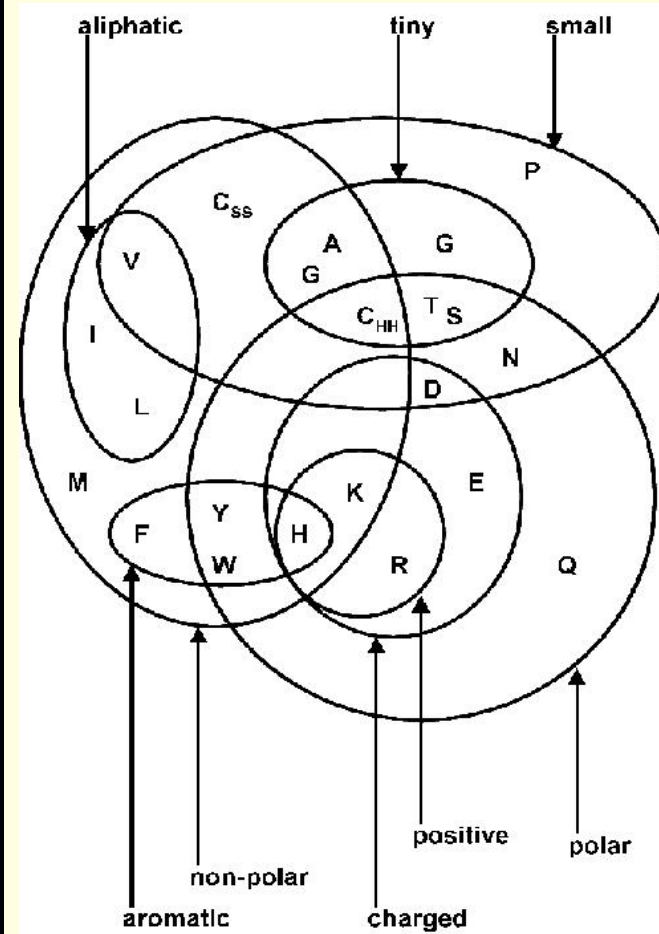
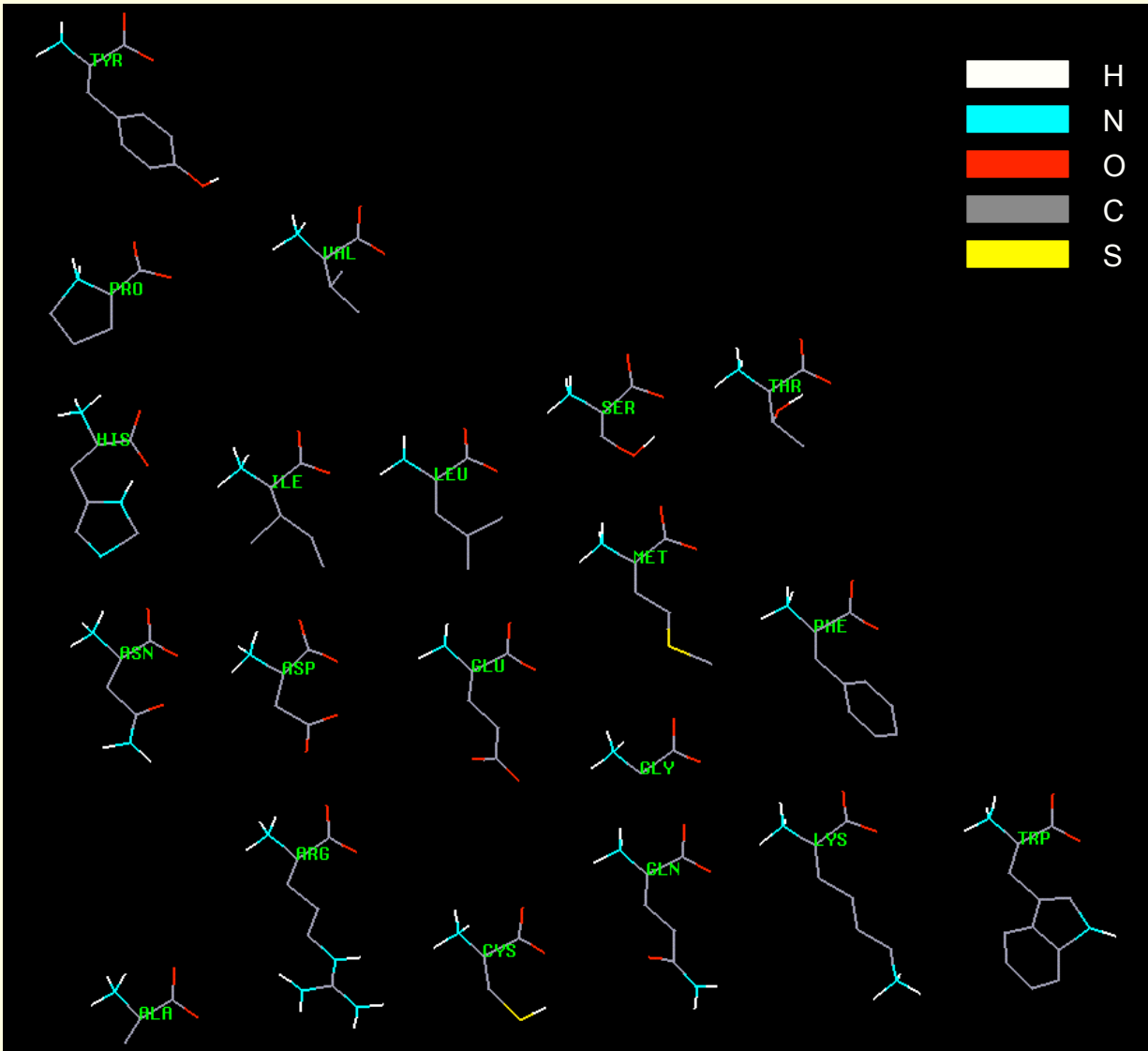
---

Sidechain



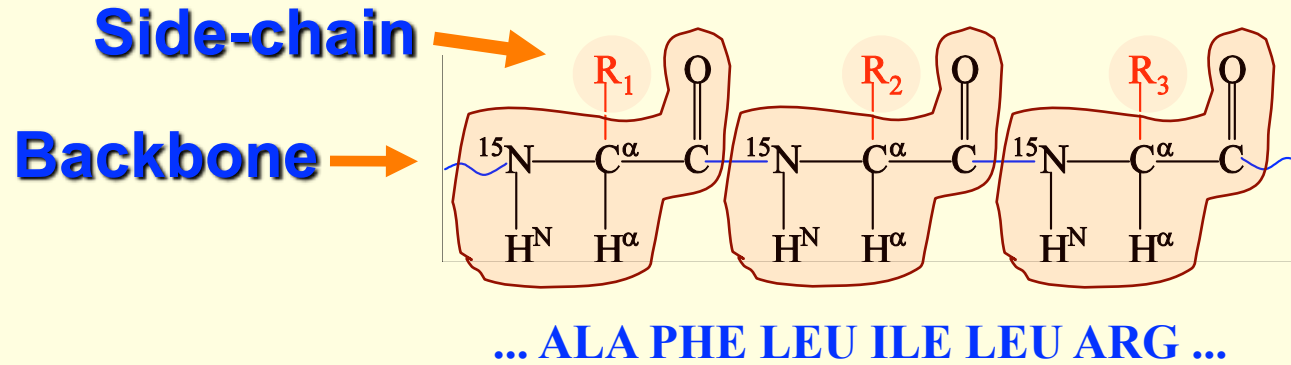
Backbone

# Amino Acids

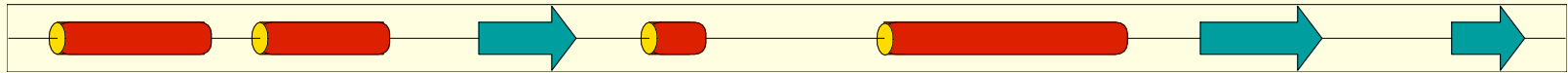


# Protein Structure

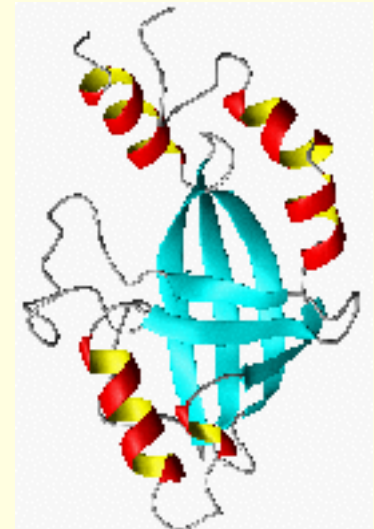
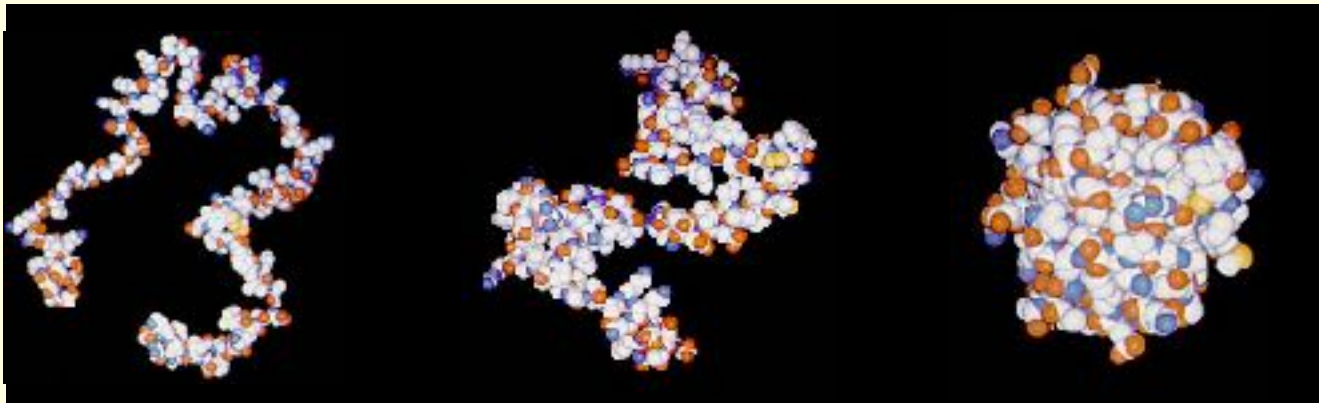
Primary Sequence: Linear String of Amino Acids



Secondary structure: regular  $\alpha$ -helices and  $\beta$ -strands



Global Fold

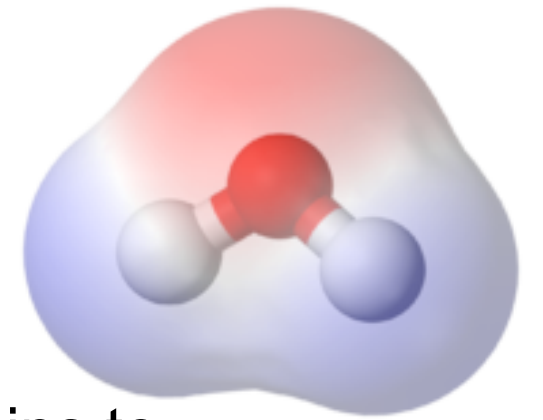


# Protein Scale and Size

- Size is measured in Daltons (Da). An average residue is ~135 Da.
- Interatomic distances are measured in Angstroms (Å),  $1 \times 10^{-10}$  m.
- Concentration is measured in mol/L (1 mol =  $6.022 \times 10^{23}$ ).
- Proteins fold to “native state” in microseconds to seconds.

# Hydrophobic Residues

- **Nonpolar** and **uncharged**
- Tend to avoid water
- Tend to interact with other nonpolar sidechains to minimize contact with water
- Tend to be buried in the protein core

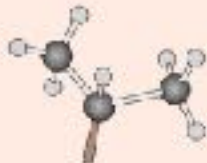


Water

## Hydrophobic



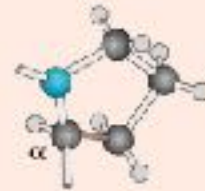
Alanine  
Ala  
A



Valine  
Val  
V



Phenylalanine  
Phe  
F



Proline  
Pro  
P



Leucine  
Leu  
L

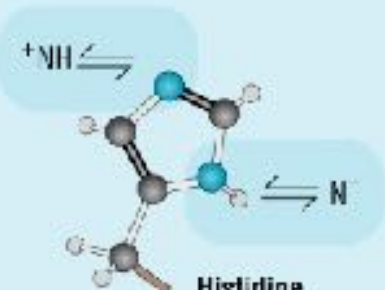
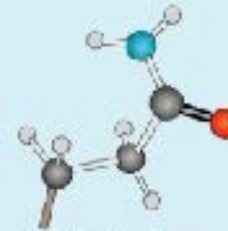
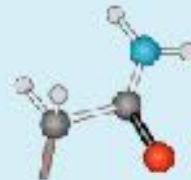
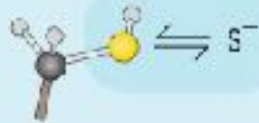
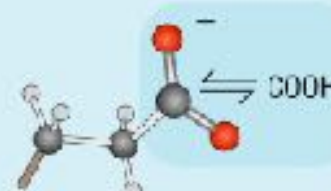
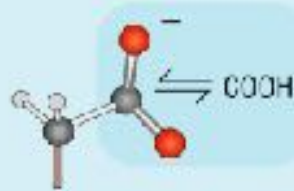
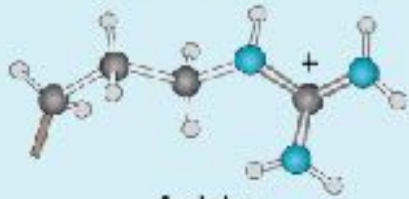


Isoleucine  
Ile  
I

# Hydrophilic Residues

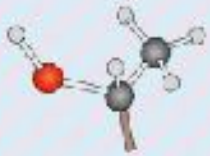
- **Polar** or **charged**
- Tend to interact with water or other hydrophilic sidechains

Hydrophilic



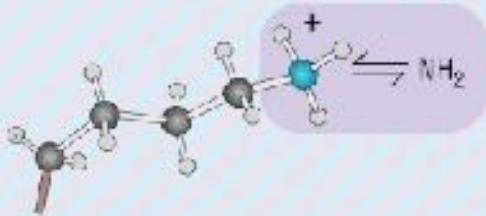
# Amphipathic Residues

- Have *both polar and nonpolar* characteristics
- Tend to form interfaces between hydrophobic and hydrophilic residues

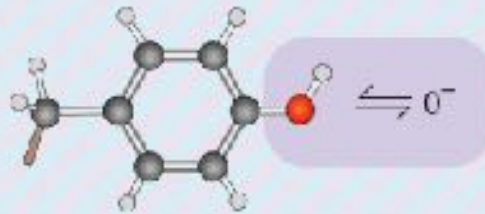


Threonine  
Thr  
T

Amphipathic



Lysine  
Lys  
K



Tyrosine  
Tyr  
Y



Methionine  
Met  
M

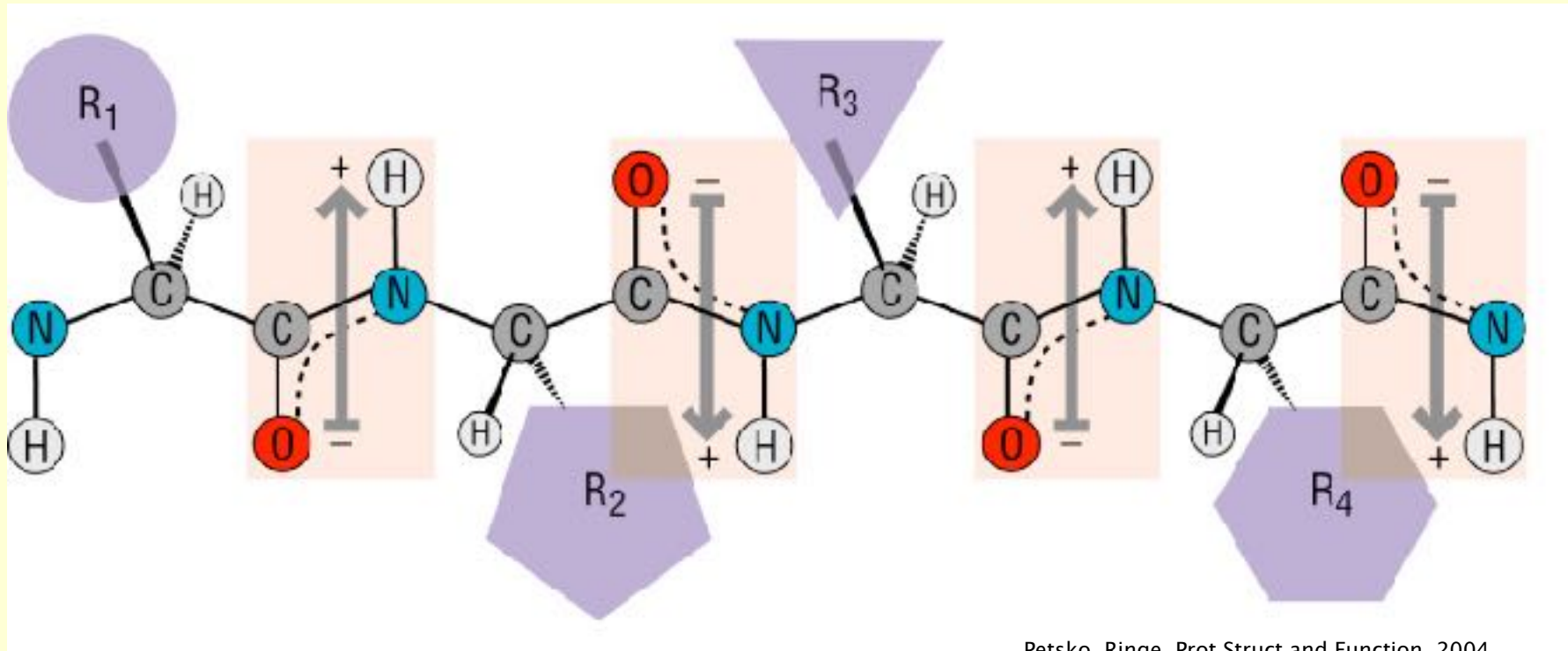
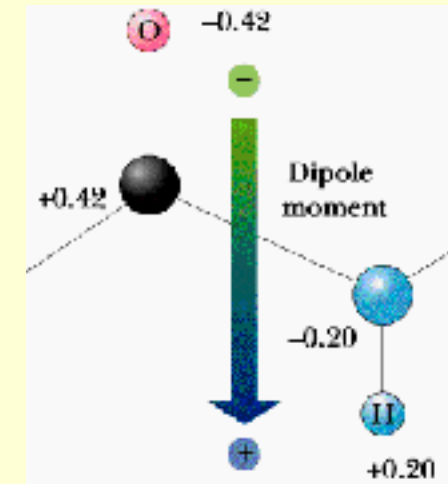


Tryptophan  
Trp  
W



# Peptide Backbone

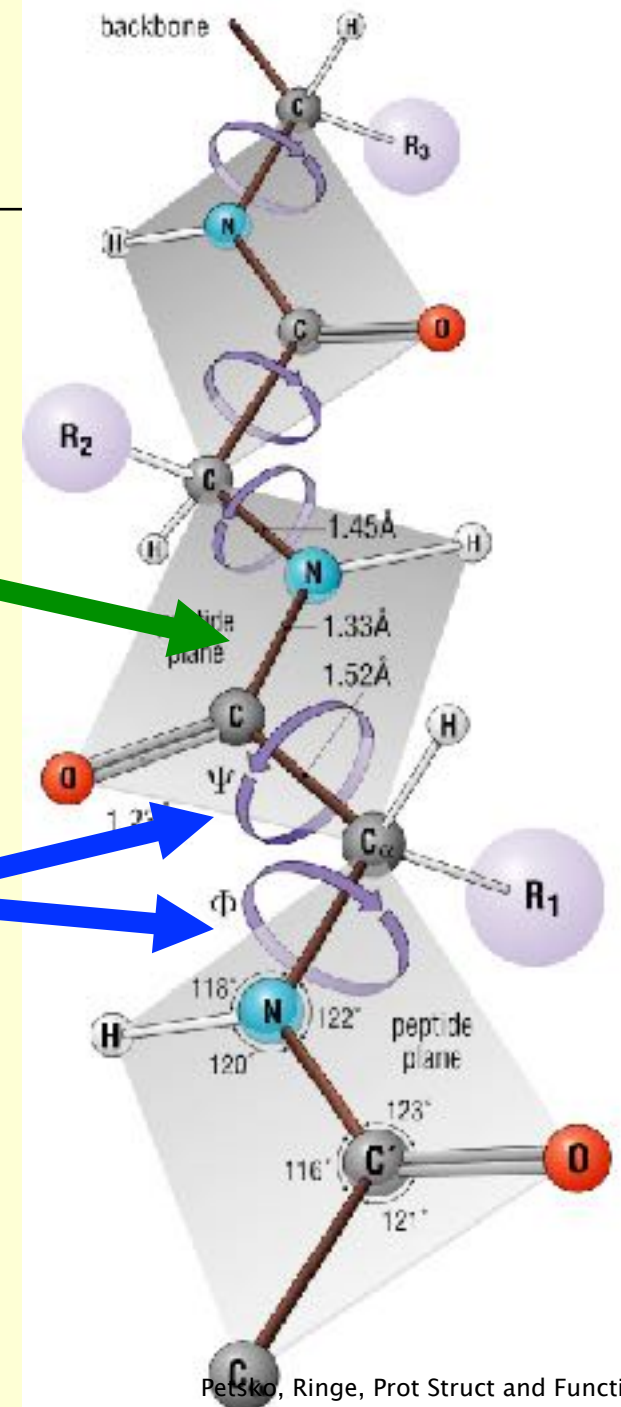
- Peptide Backbone is polar (dipolar)
- Double bonds are not very flexible, but single bonds....

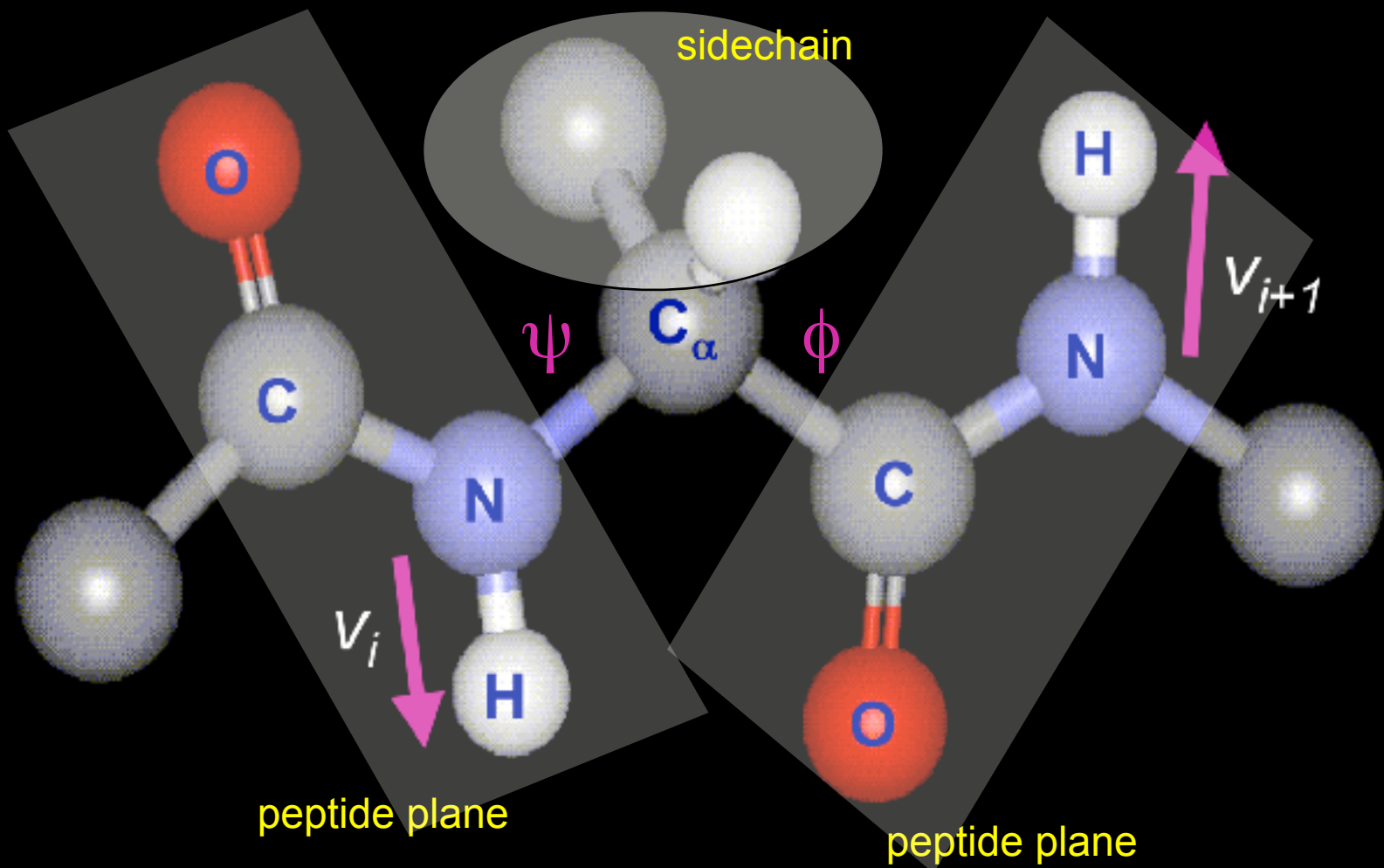


# Peptide Backbone

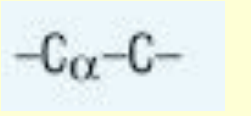
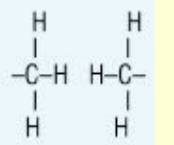
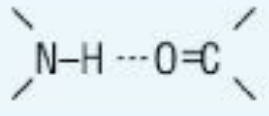
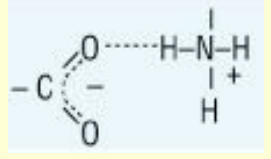
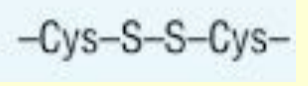
Peptide bond has partial double bond character and is rigid  
(ie. it doesn't rotate)

Other backbone bonds are flexible  
psi torsion angle  
phi torsion angle

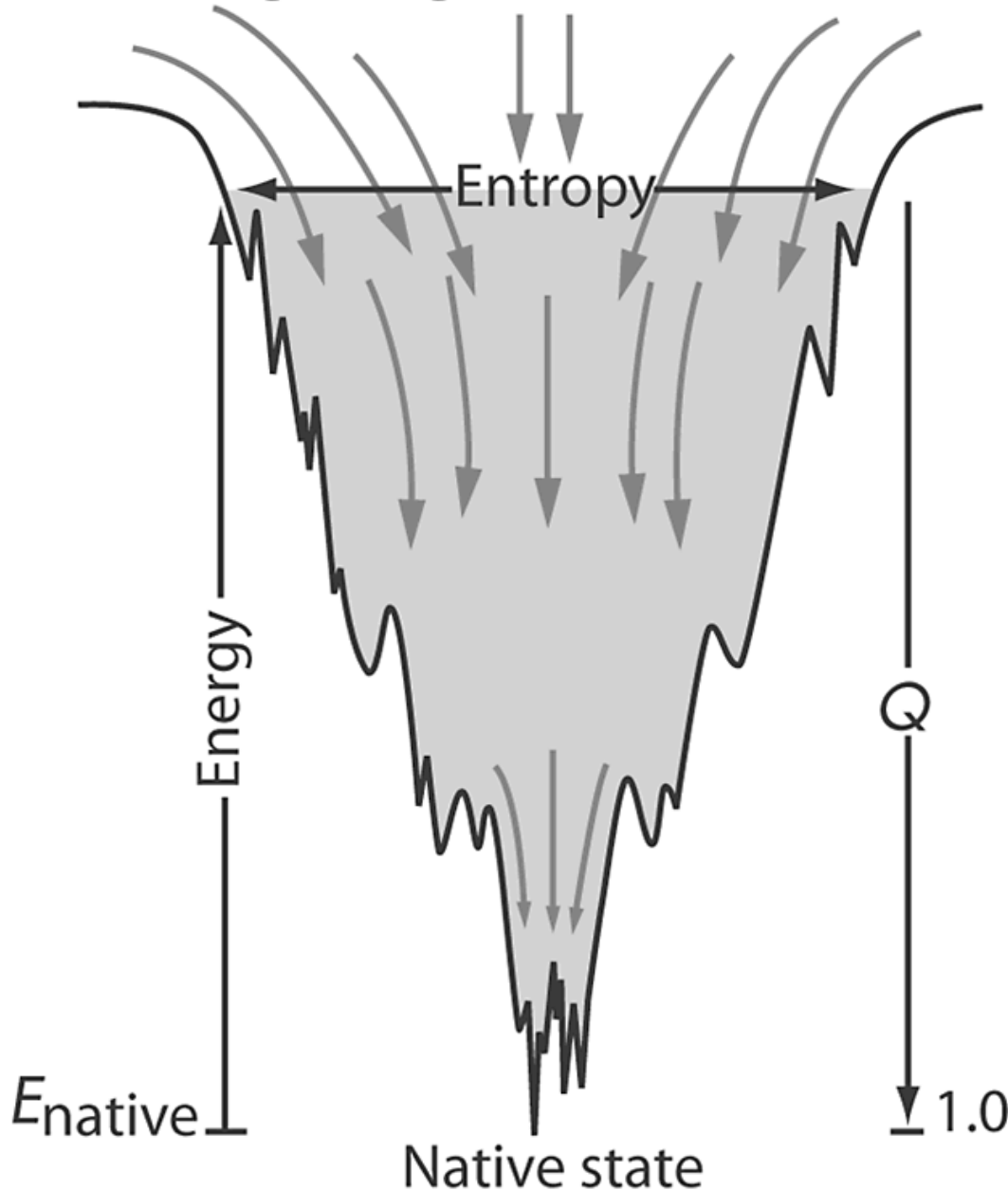




# Stabilizing Forces

| Interaction   | Typical Distance | Free Energy  |
|---|------------------|--|
| * Covalent Bond   | 1.5 Å            | 356 kJ/mol   |
|    |                  |  |
| * Van der Waals Interaction   | 3.5 Å            | 4-15 kJ/mol  |
|    |                  |  |
| * Hydrogen Bond   | 3.0 Å            | 2-6 kJ/mol<br>Up to 20 kJ/mol if one atom is charged |
|    |                  |  |
| Salt Bridge   | 2.8 Å            | 13-17 kJ/mol   |
|   |                  |  |
| Disulfide Bond  | 2.2 Å            | 167 kJ/mol   |
|  |                  |  |
| * Long Range Electrostatic  | variable         | variable   |

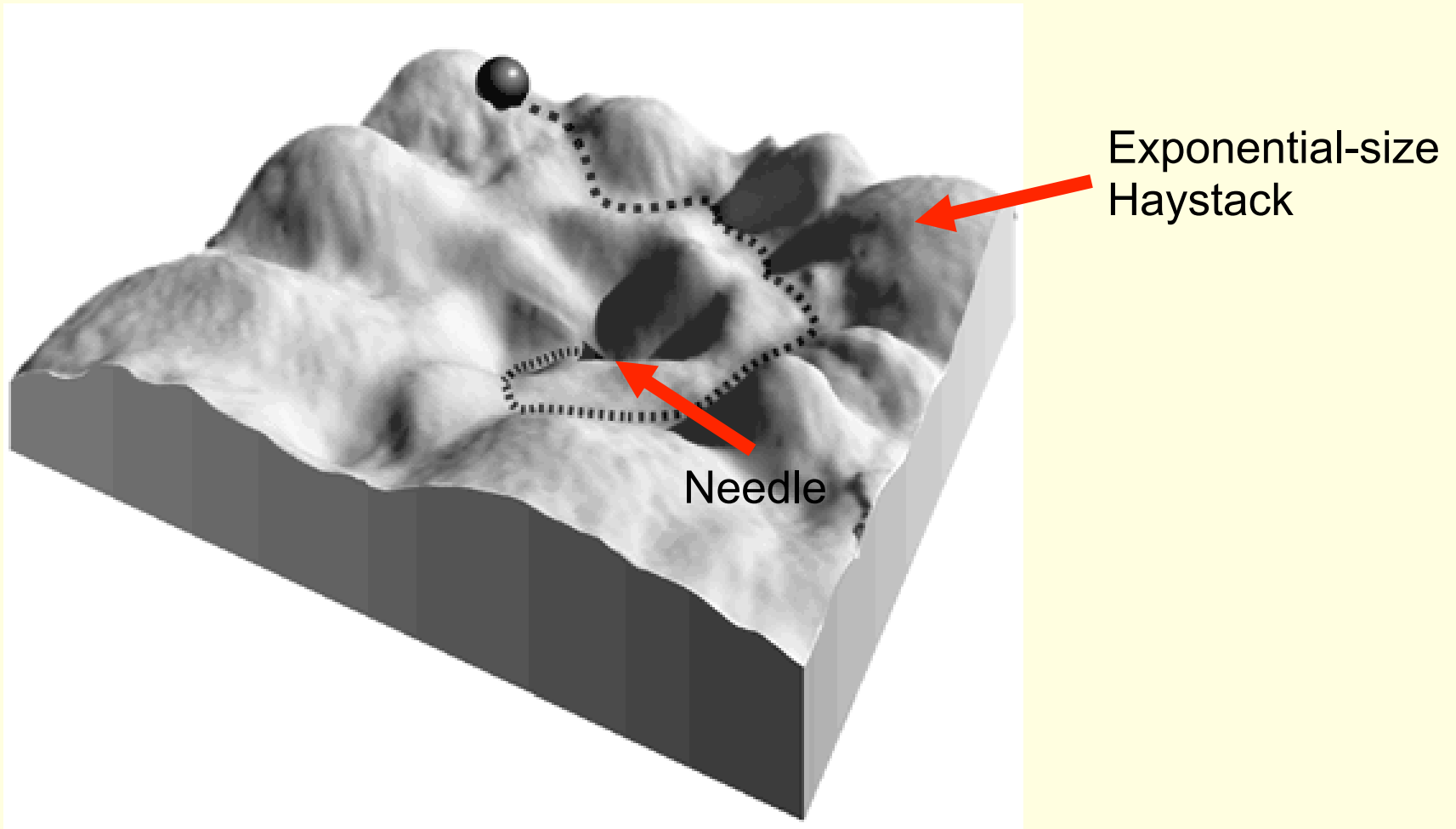
Beginning of helix formation



Protein folding is a physical/chemical process.

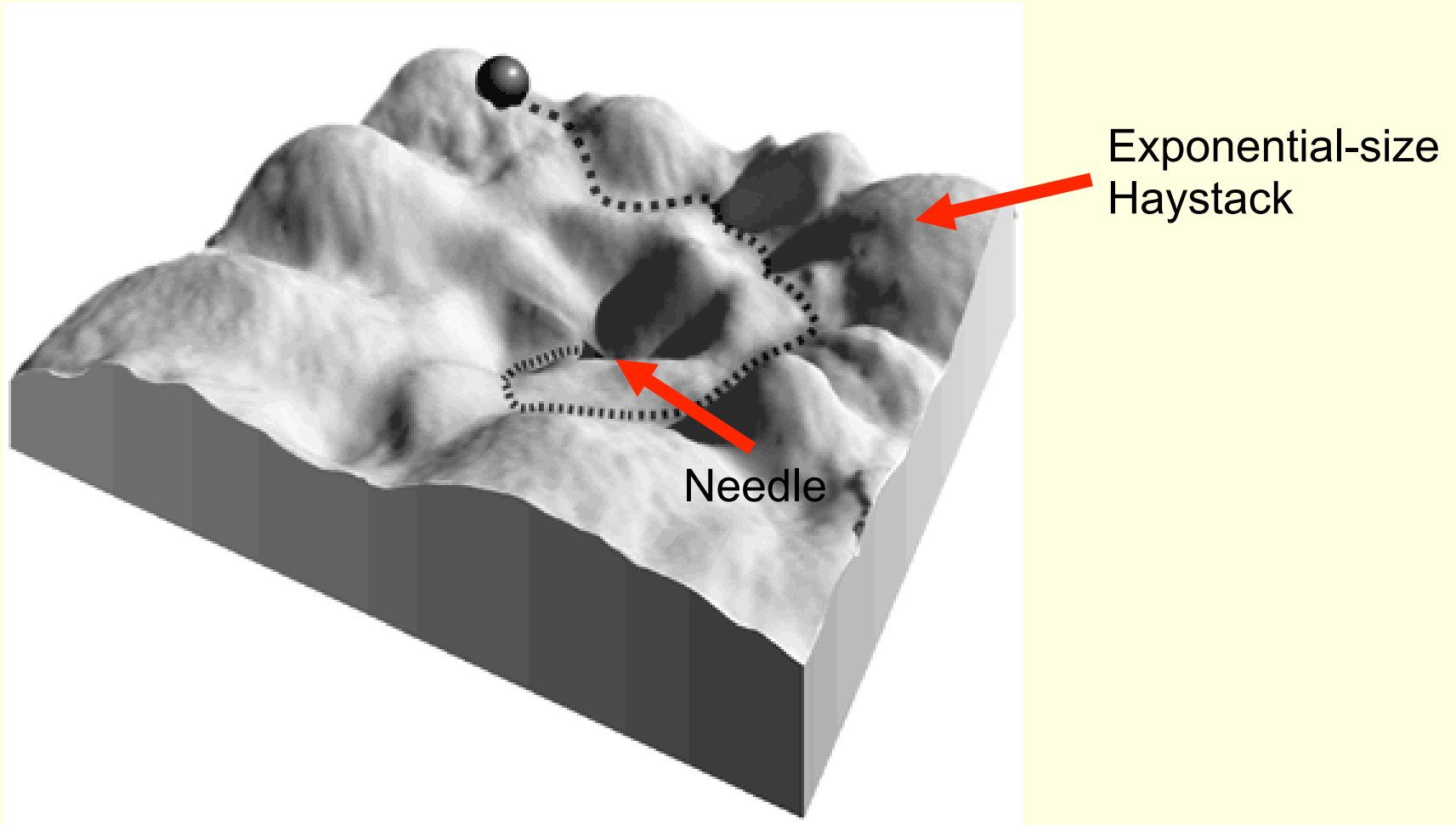
**Thesis:** Proteins adopt low energy conformations in their native state.

Conformational energy is determined by solvent (water), van der Waals forces, charge, electrostatic forces, etc.



Levinthal's Paradox (1969): The conformation space is exponential in the number of amino acids, but the folding pathway must be relatively short.

# *De novo* Structure Prediction



Does this landscape have a more compact representation?  
Can we quickly find low-energy conformations?

# Backbone Flexibility

---

- Protein with  $n$  residues has  $2n$  rotatable backbone dihedrals
- $2n$  degrees-of-freedom (DOF)
- Number of protein conformations is exponential in its length (consider kinematic chain)

Large Search Spaces!





# How Many Proteins Are There?

^  
Possible

## Sequences

Avg Protein 300 Amino Acids (AAs) long

$$20^{300} \approx 10^{390}$$

There are only  $10^{80}$  particles estimated in the universe

## Human Proteins

Current Estimate ~20,000 Genes \*conservative

Each gene has alternate splicing (say 2 / protein)

Each protein can be post-synthetically modified (say 2 / protein)  
cleavage (insulin), phosphorylated (kinase), glycosylated (sugars), ...

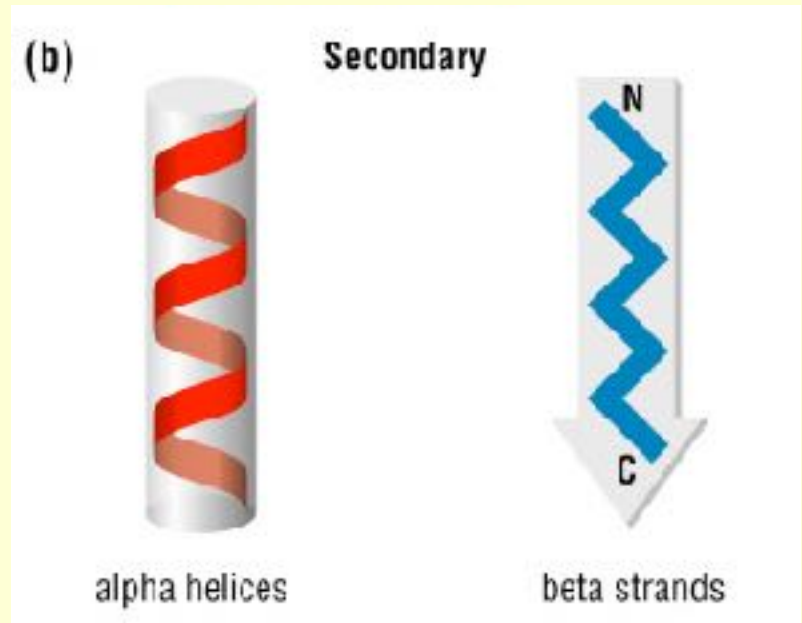
SO

$$20,000 \times 2 \times 2 = 80,000$$

# Levels of Protein Structure

---

```
GLY SER MET SER
GLY ILE ALA LEU
SER ARG LEU ALA
GLU GLY GLY LEU
PHE LYS LEU ARG
MET LEU LEU ASN
VAL GLU TYR GLU
LYS ARG VAL ARG
ALA GLN ALA LYS
```



Petsko, Ringe, Prot Struct and Function, 2004

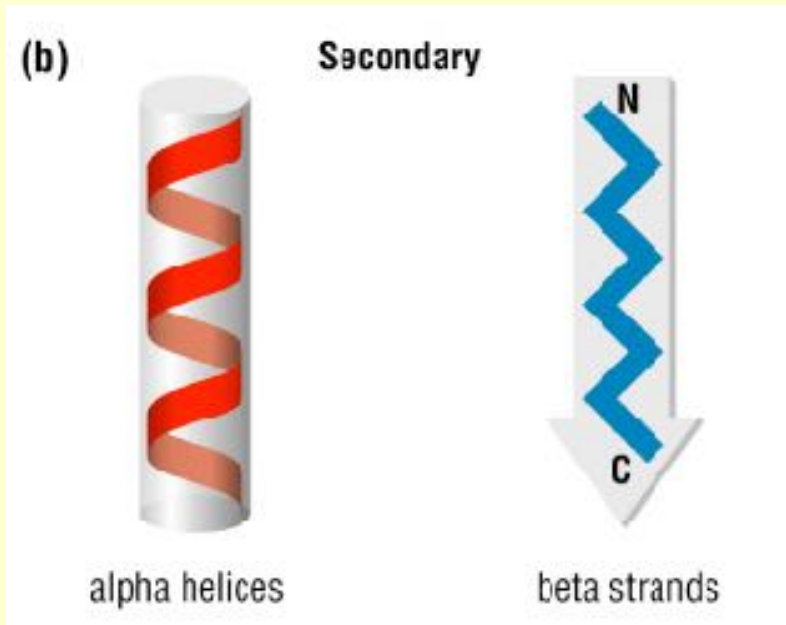
**Primary Structure**

Aka Primary Sequence

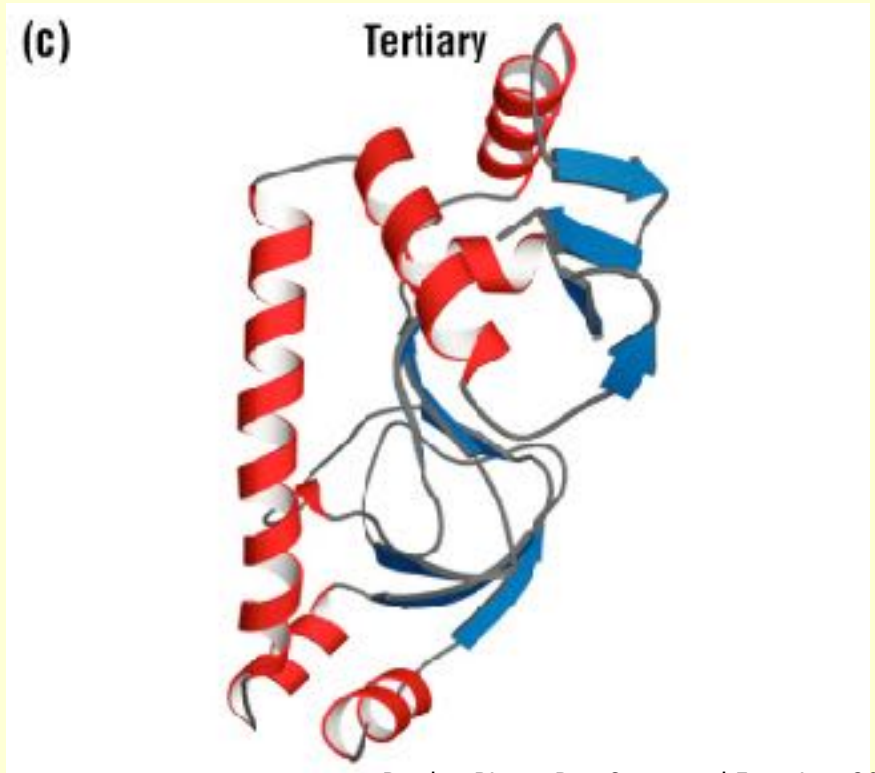
**Secondary Structure**

# Levels of Protein Structure

---



Secondary Structure



Petsko, Ringe, Prot Struct and Function, 2004

Tertiary Structure

# Levels of Protein Structure

---

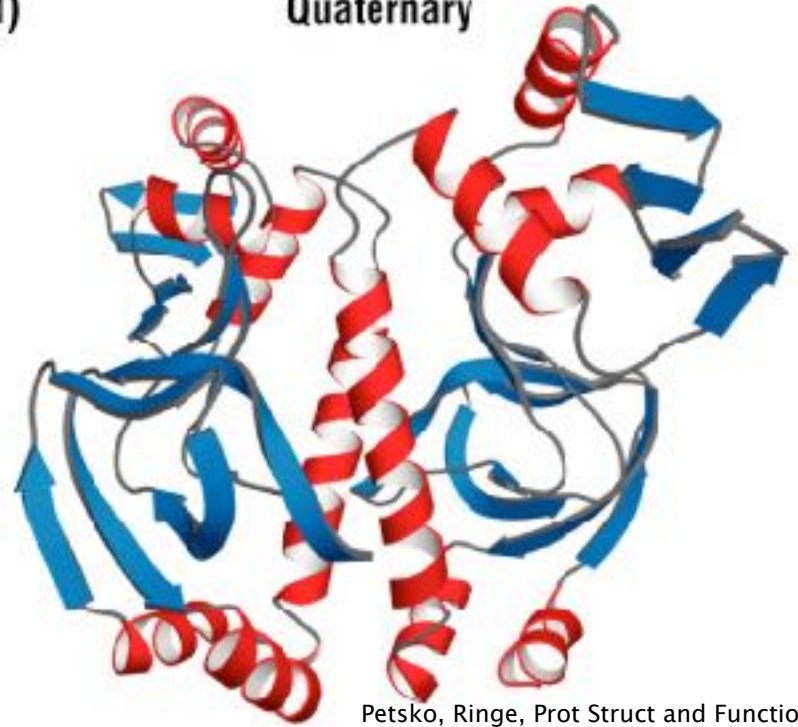
(c)

Tertiary



(d)

Quaternary

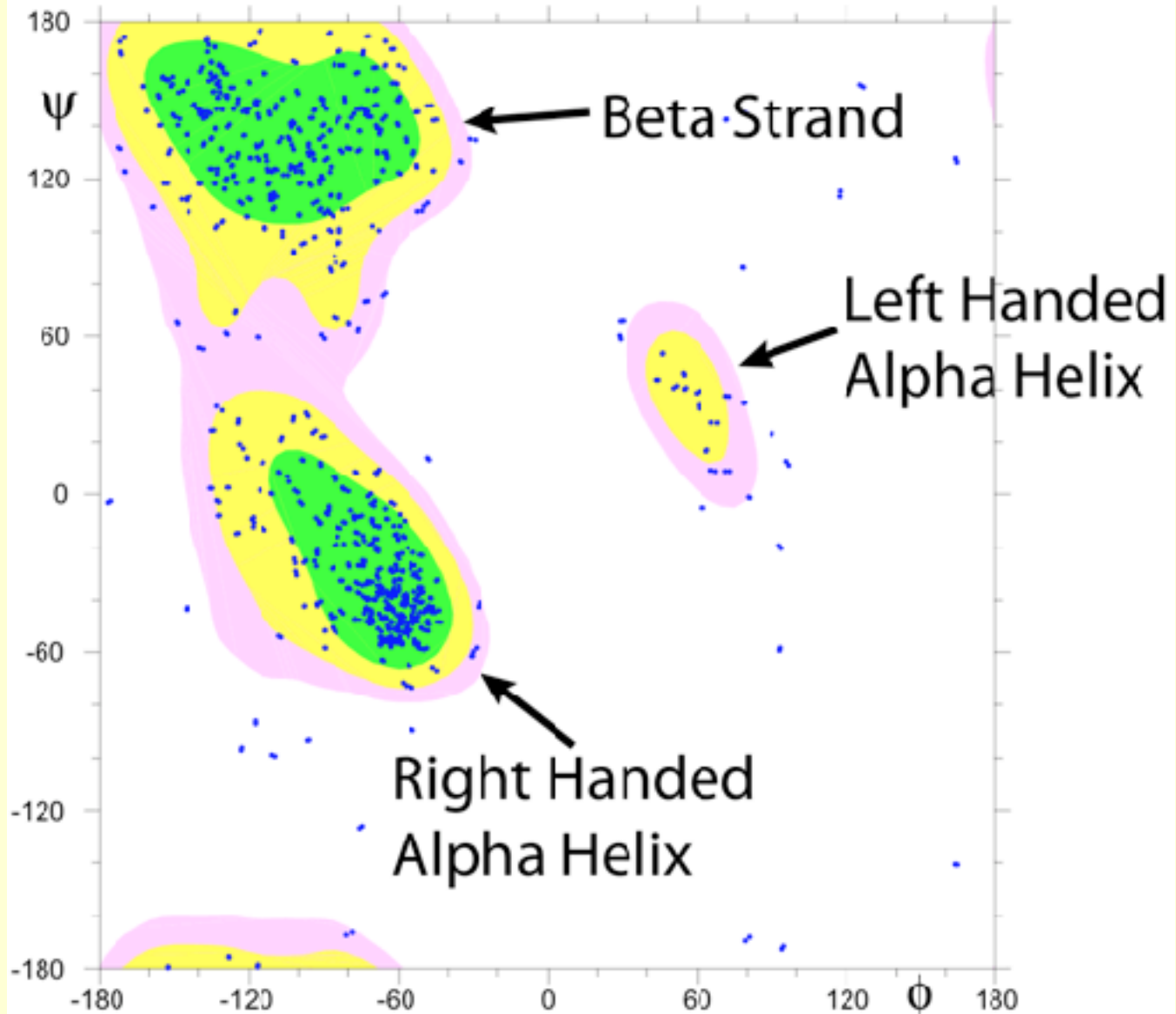


Petsko, Ringe, Prot Struct and Function, 2004

Tertiary Structure

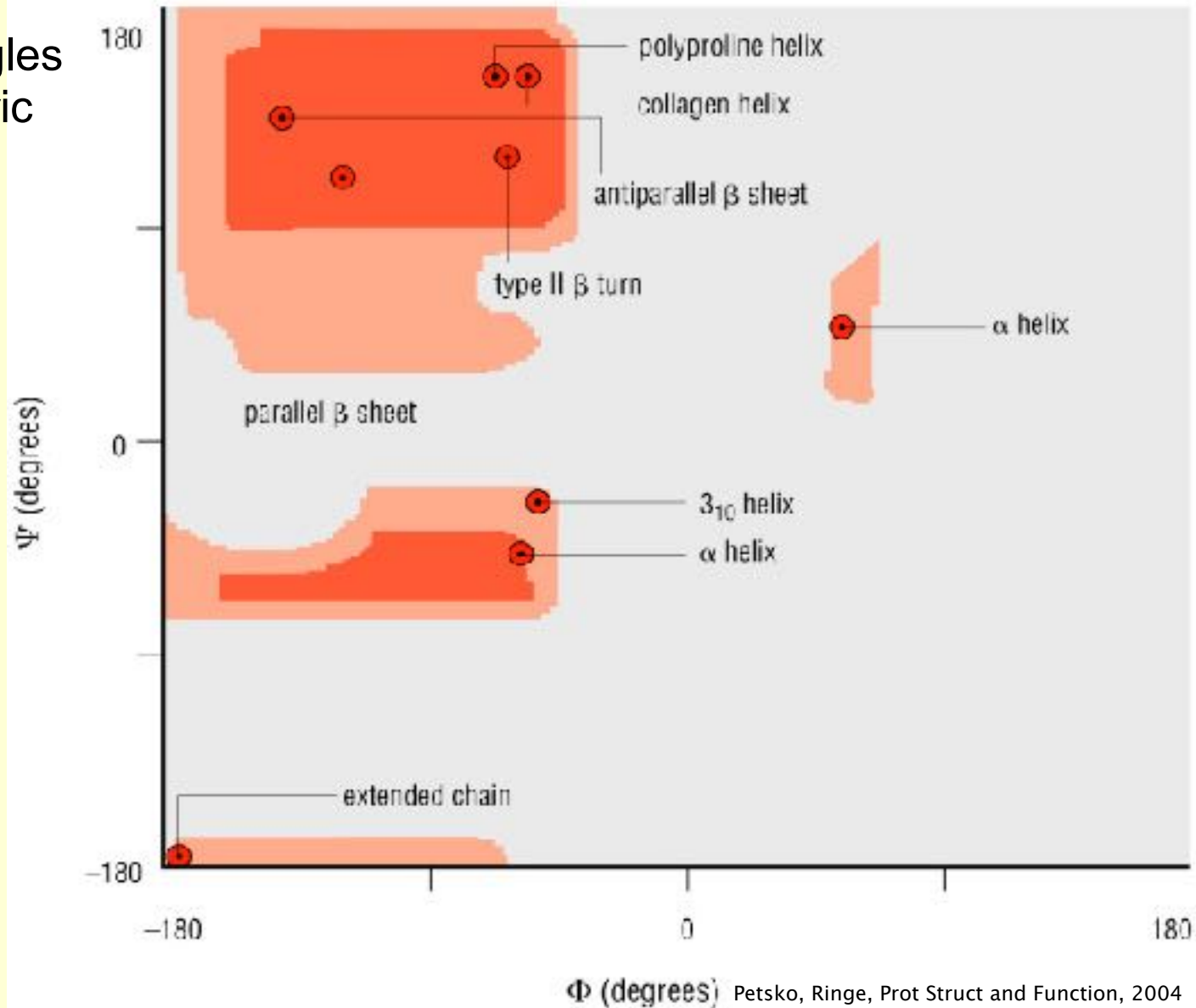
Quaternary Structure

# Phi/Psi Histogram

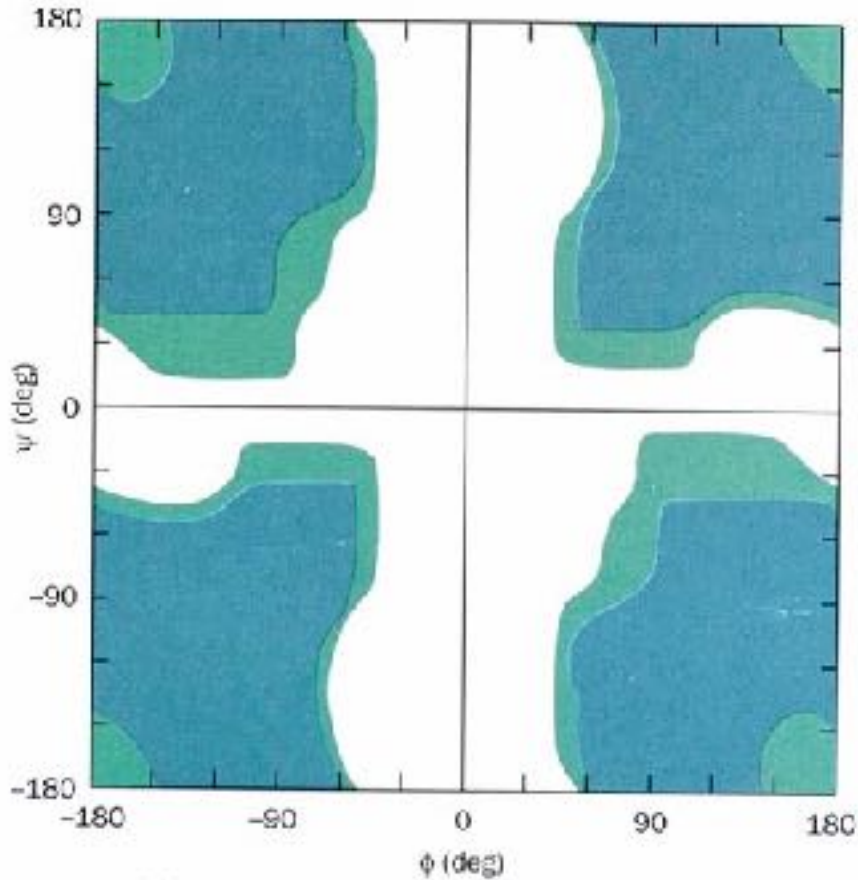


# Ramachandran Plot

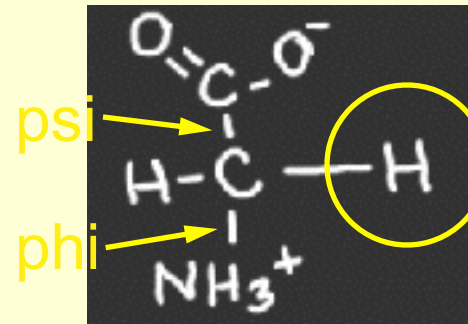
Allowed phi/psi angles do not result in steric interference



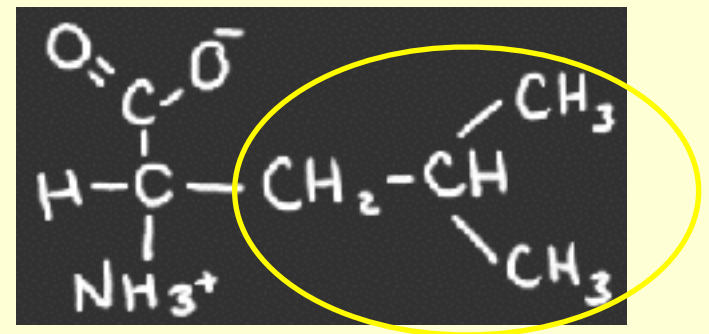
# Phi / Psi Histogram for Glycine



## Glycine

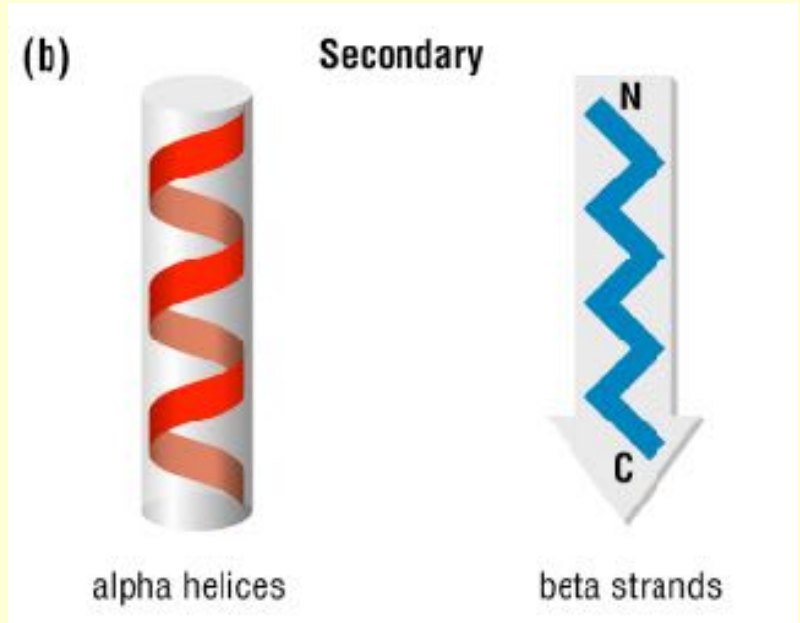


## Leucine



# Secondary Structure

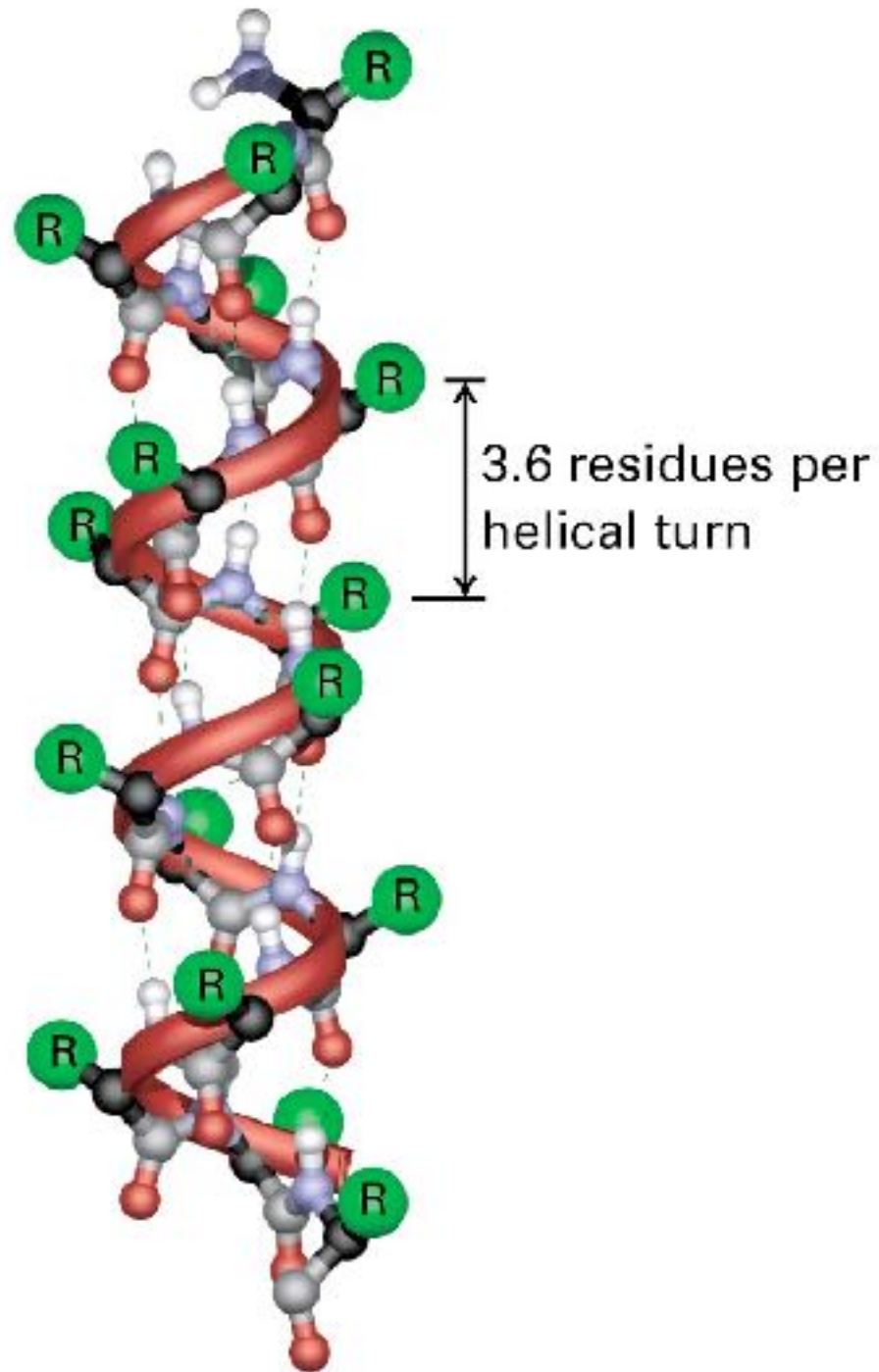
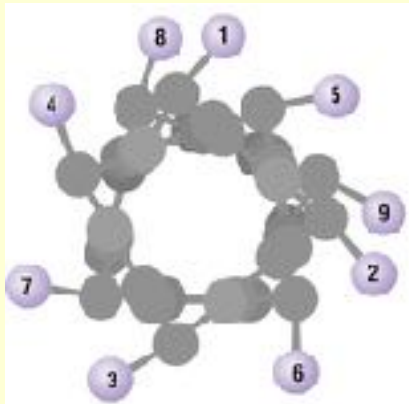
No true gold standard definition  
May be dynamic - change with  
changing protein state



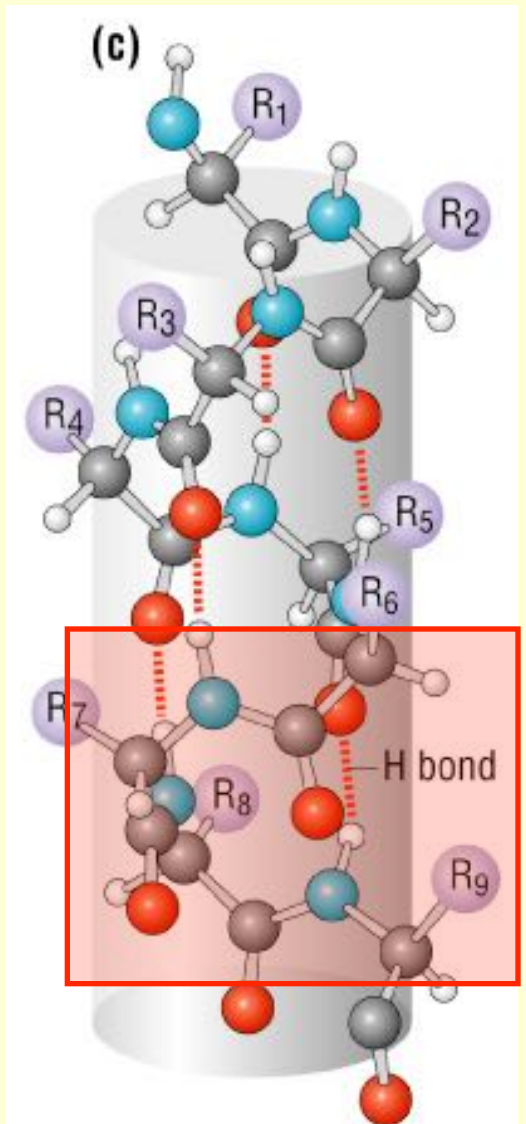
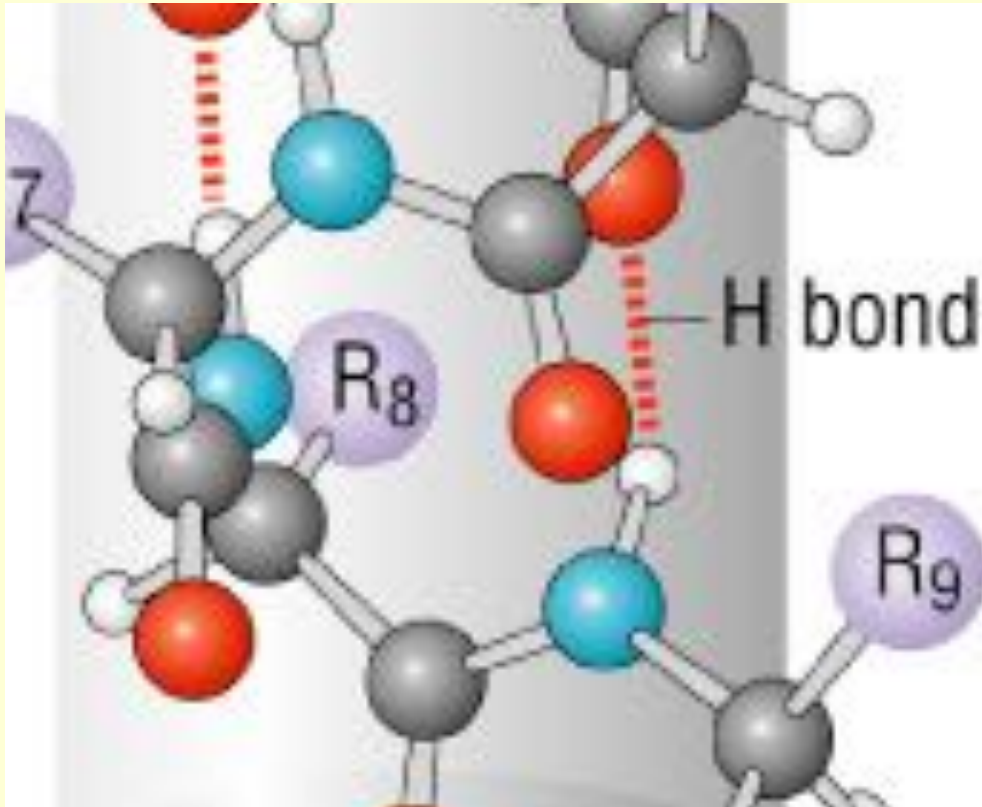


# The Alpha Helix

- Most common secondary structure type
- Hydrogen bonding between carbonyl oxygen atom of residue  $n$  and amide nitrogen of residue  $n+4$ .
- 3.6 Residues per turn
- Cylindrical structure with hydrogen-bonded wall and outside studded with side chains

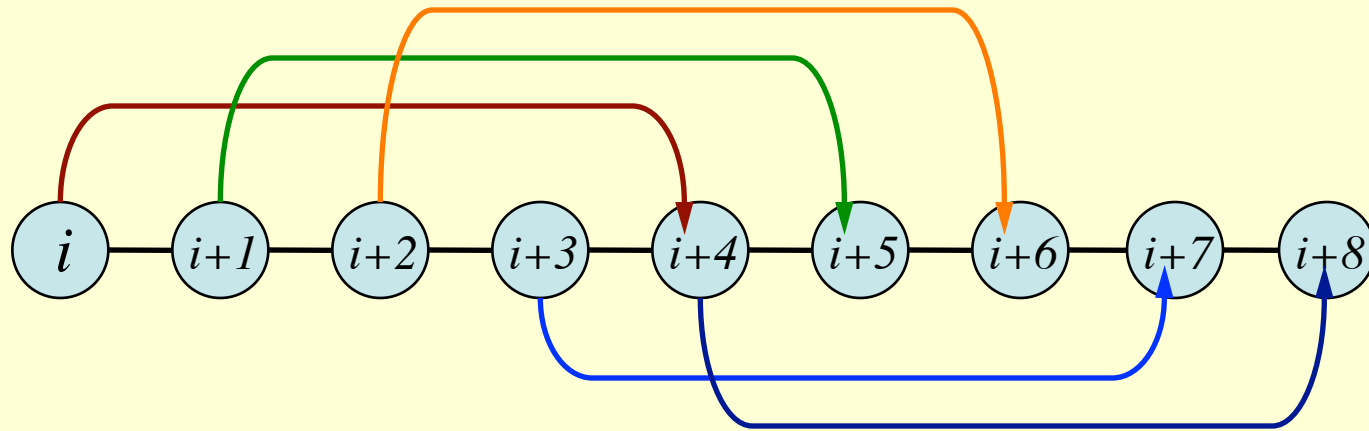


# The Alpha Helix



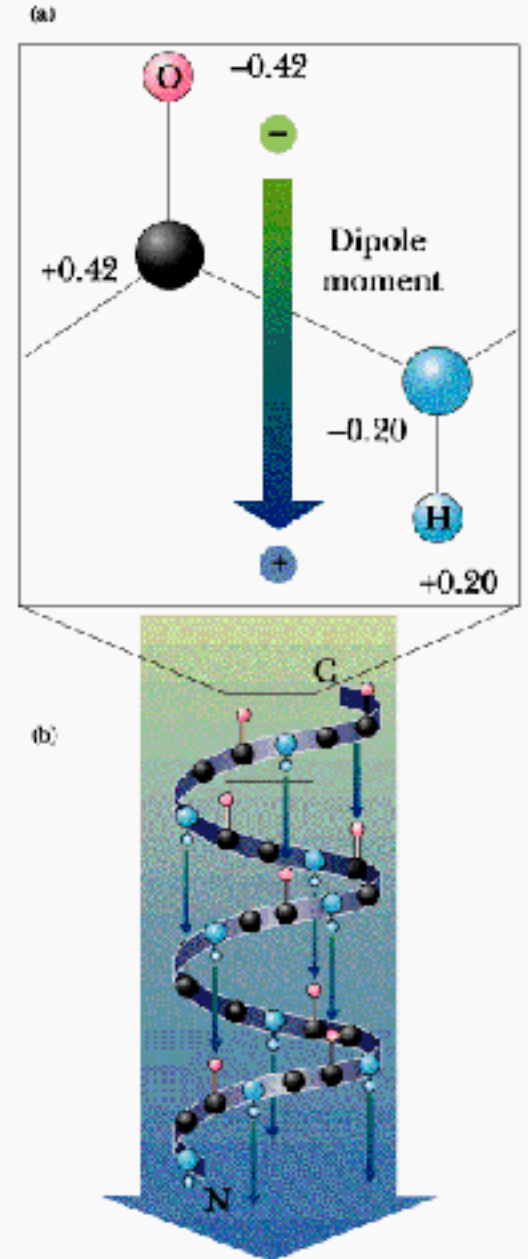
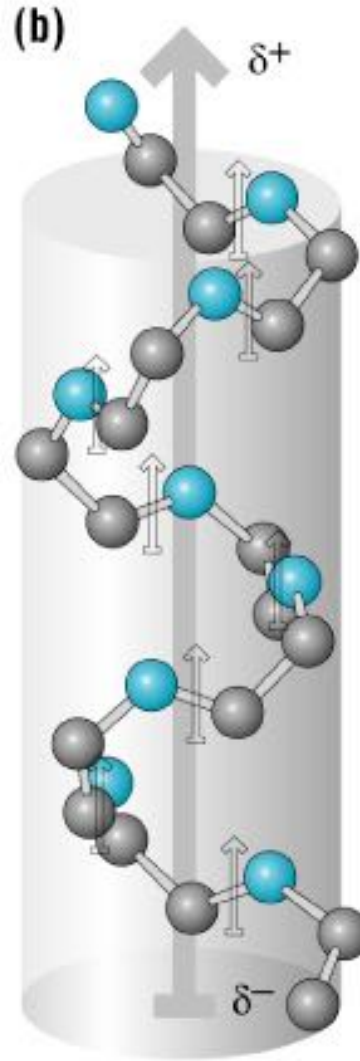
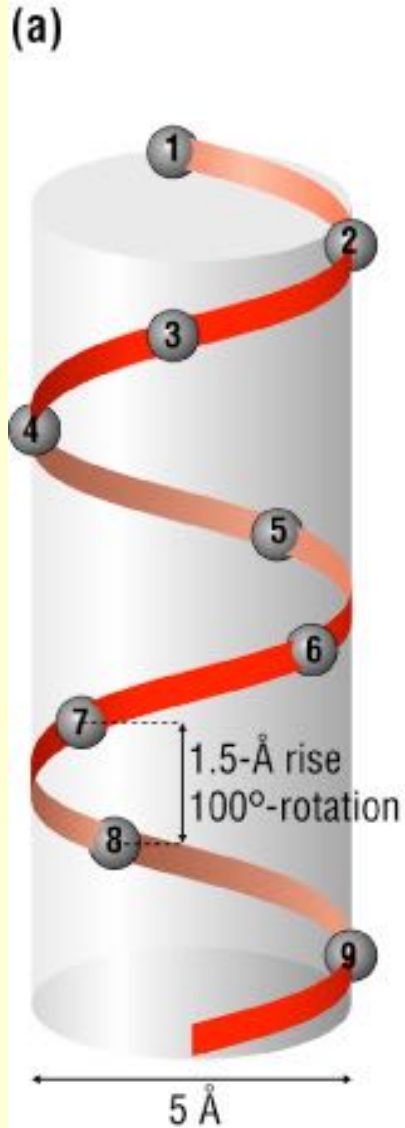
# Hydrogen Bonding Pattern

---

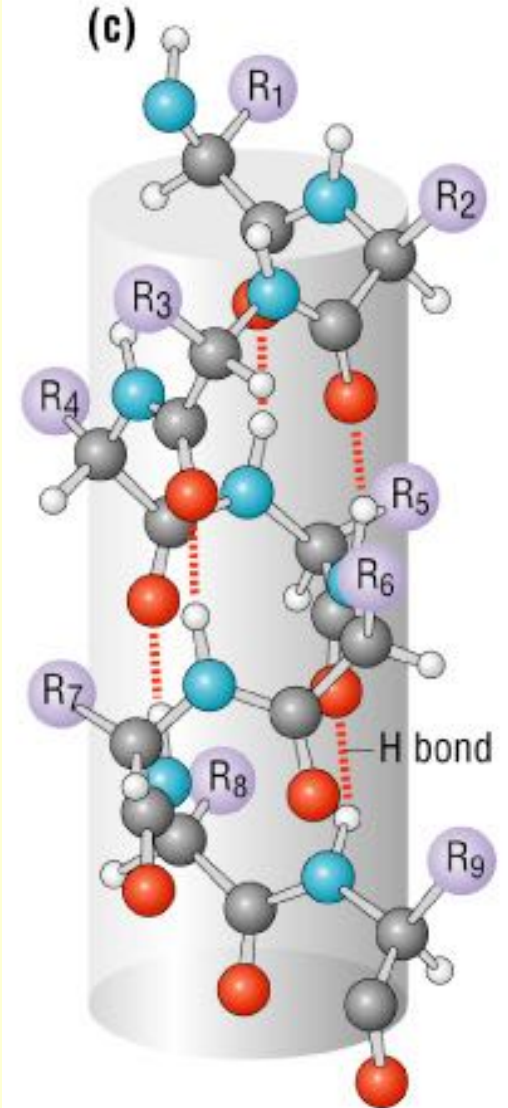
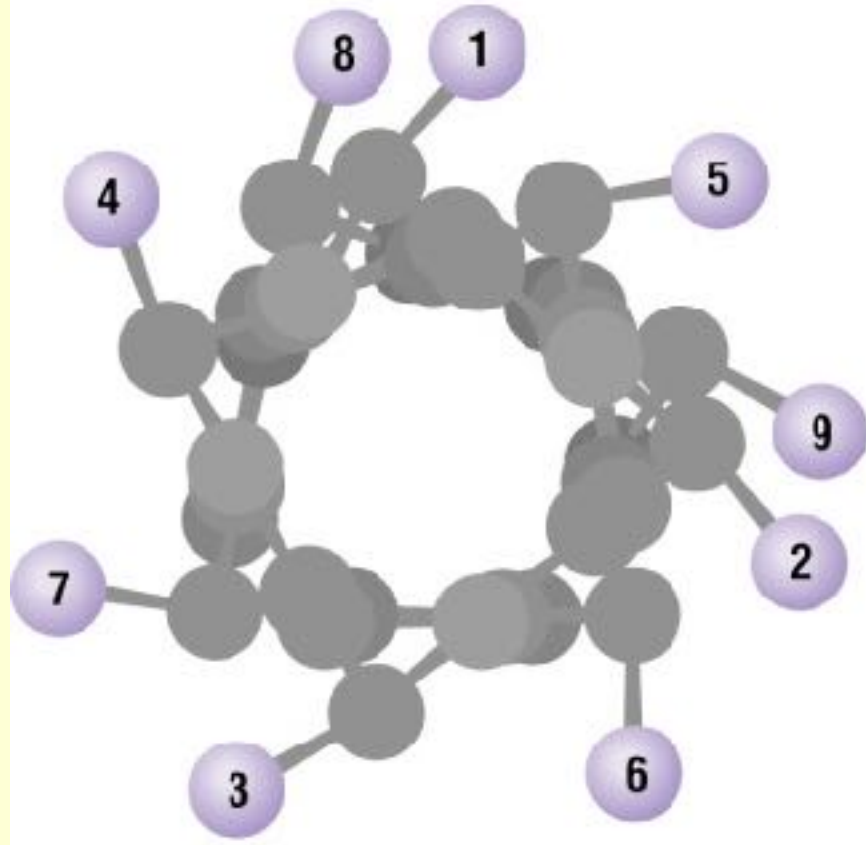


# The Alpha Helix

Garrett & Grisham: Biochemistry, 2/e  
Figure 6.8



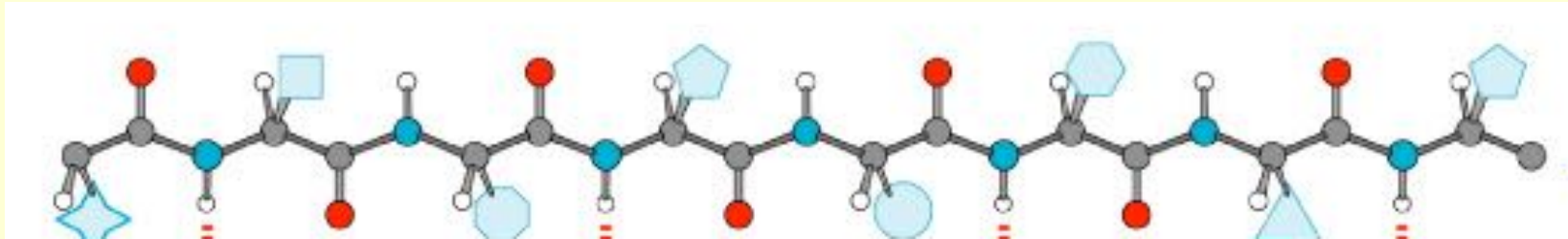
# The Alpha Helix



# Beta Strands

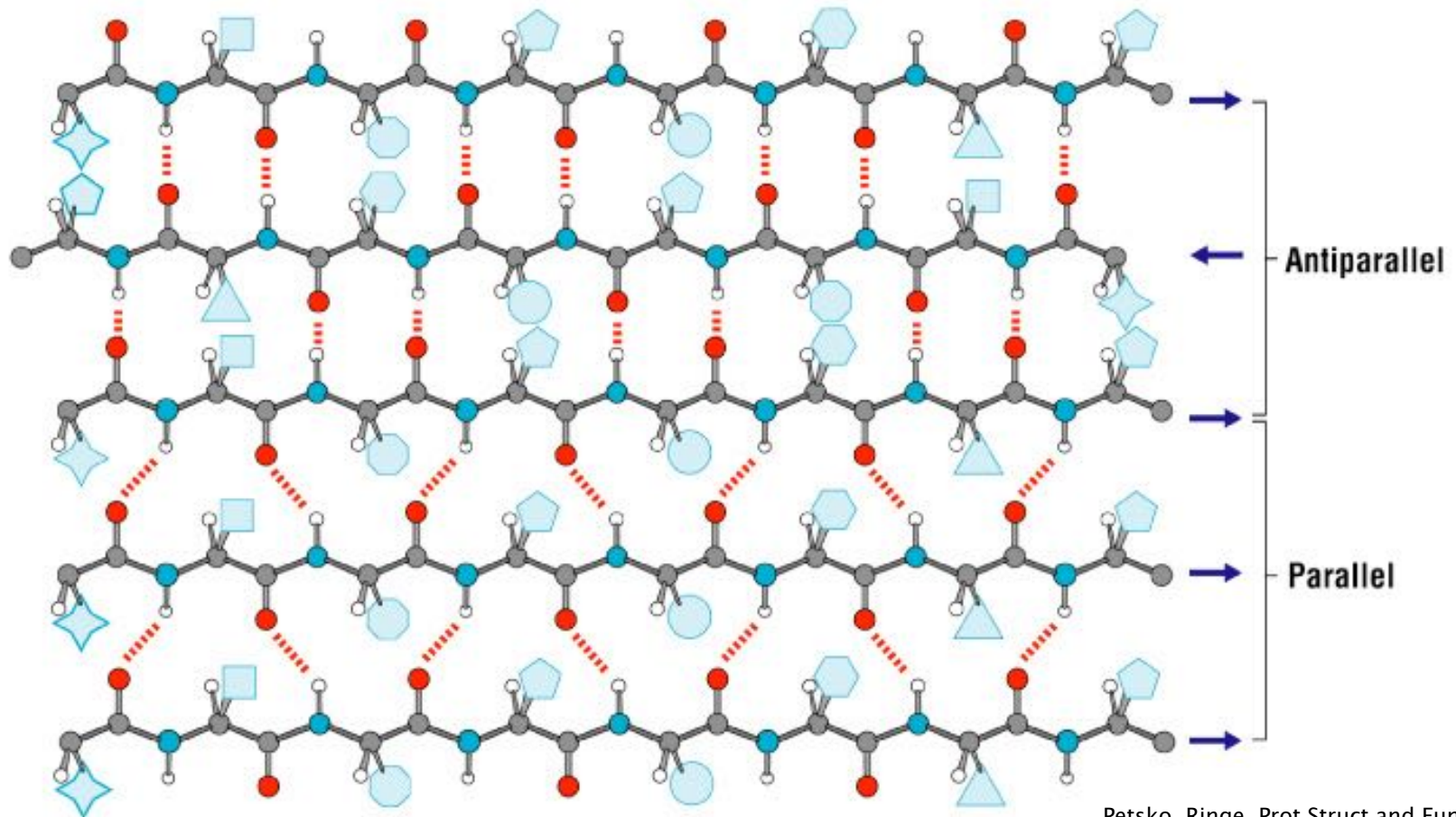
---

- Side-chains protrude in opposite directions
- Relatively linear backbone

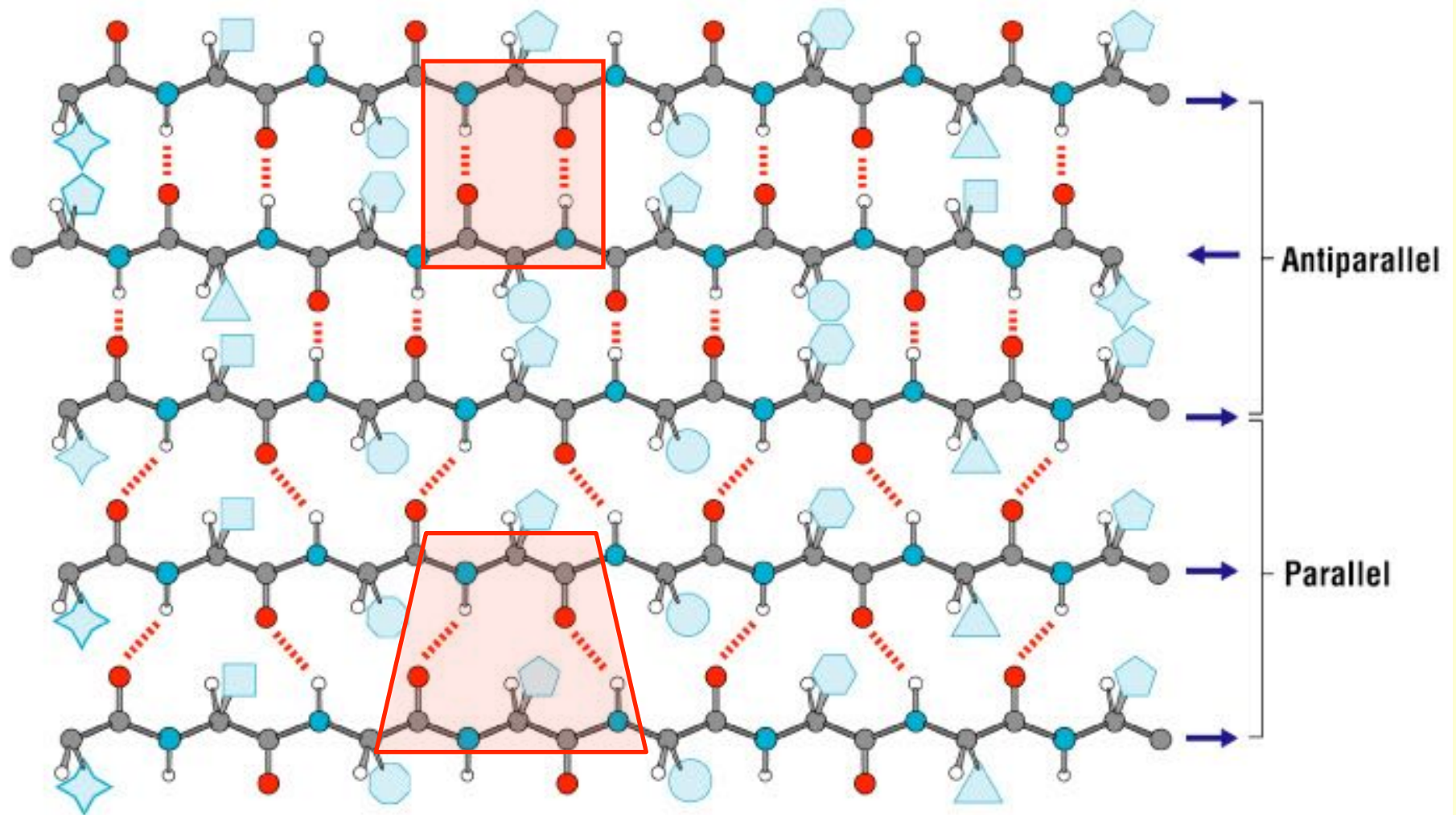


# Beta Sheets

- Hydrogen bonding between backbone of strands
- Strands of a sheet may be separated by arbitrary number of amino acids



# Hydrogen Bonding

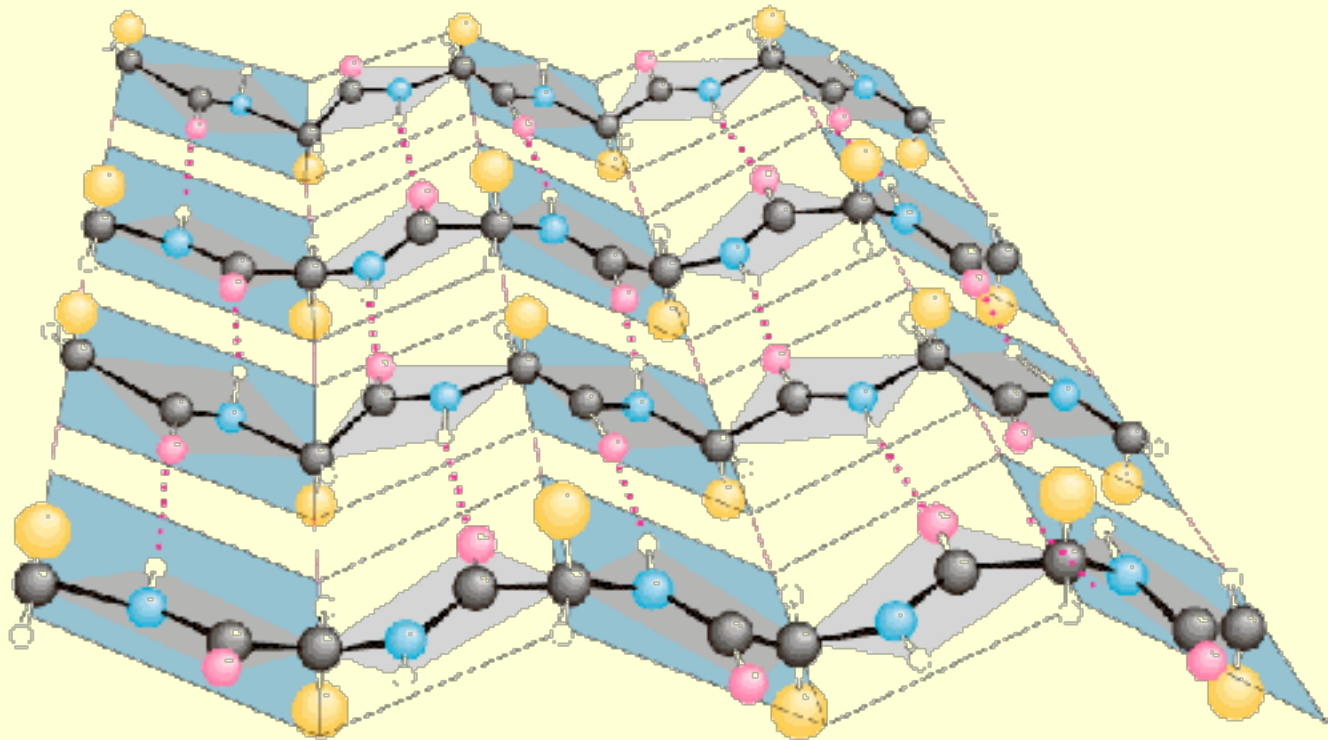




# Beta Sheets

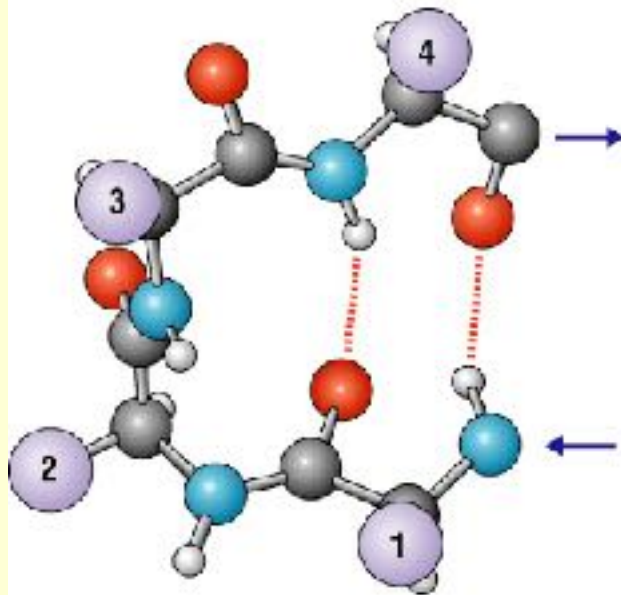
---

- Hydrogen bonding between backbone of strands
- Strands of a sheet may be separated by arbitrary number of amino acids

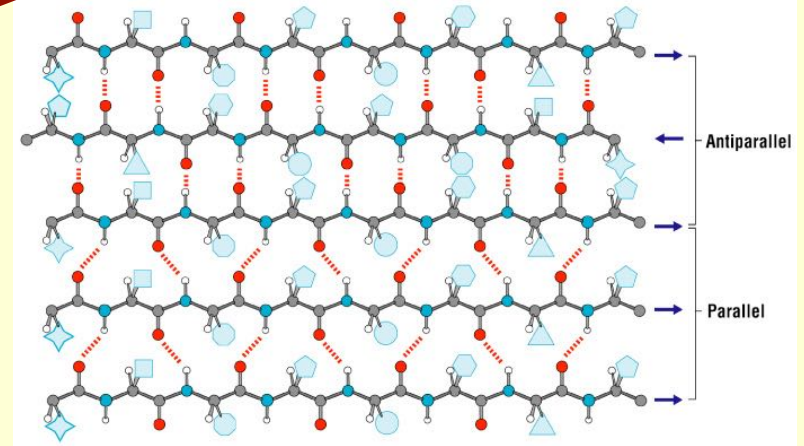


# AntiParallel Beta Sheets

- Antiparallel sheets most commonly have **beta-turns** (aka. hairpin turn) connecting strands



Hydrogen bonds between  
residues 1 and 4



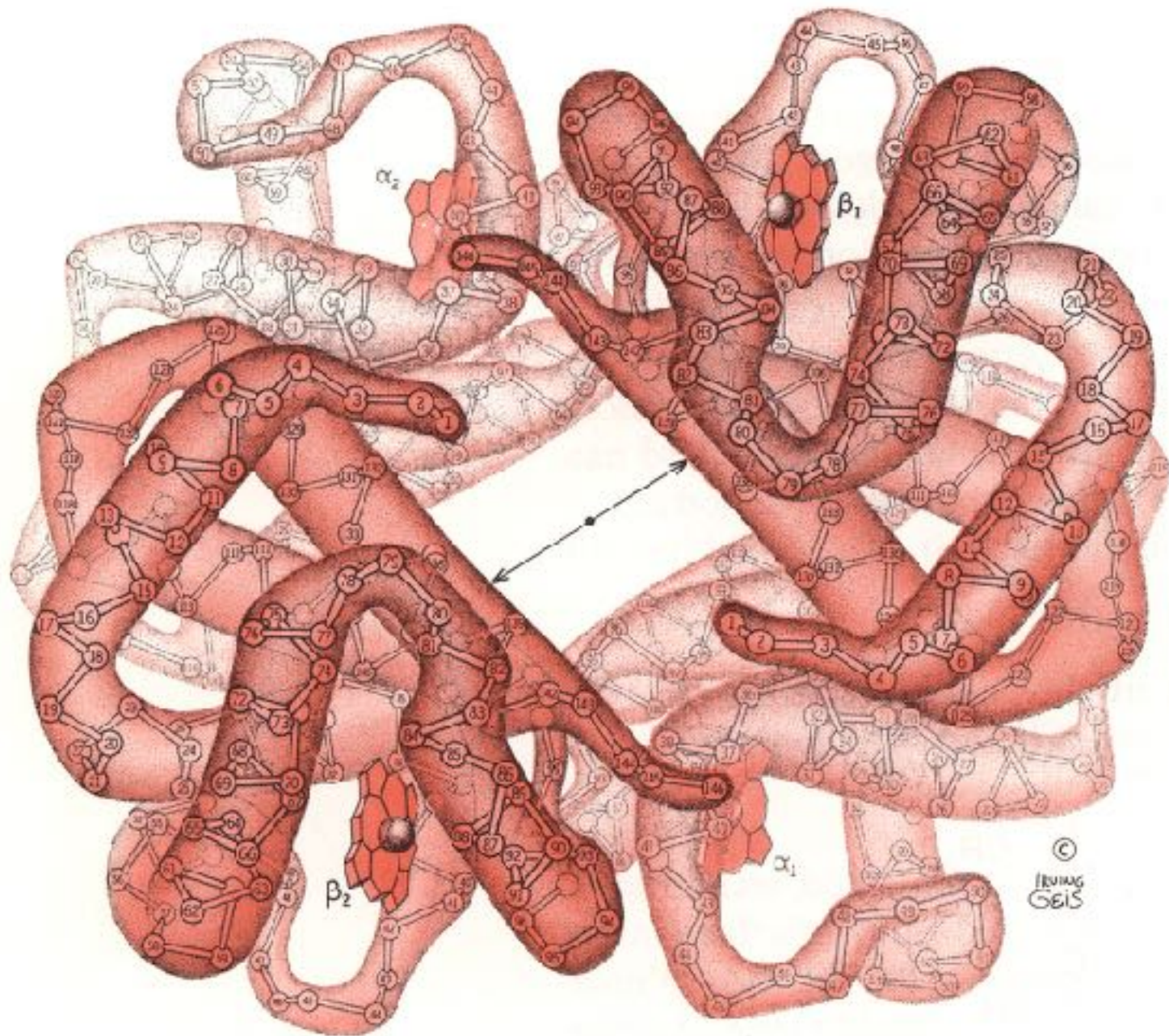
# Parallel Beta Sheets

---

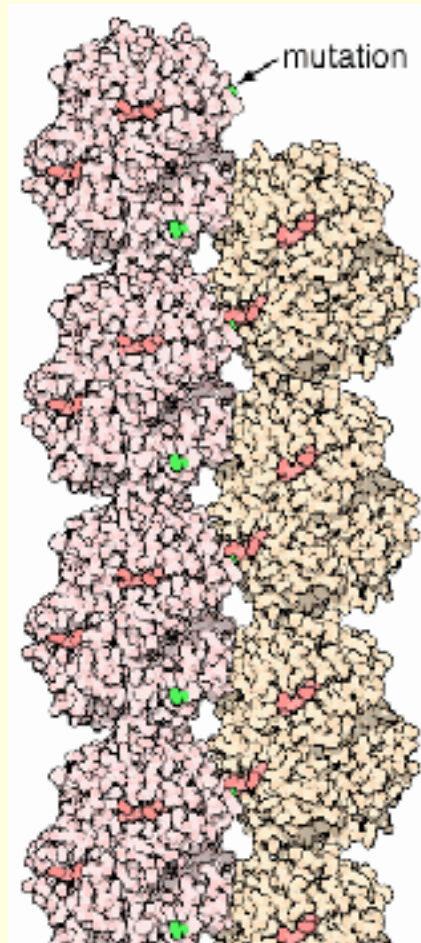
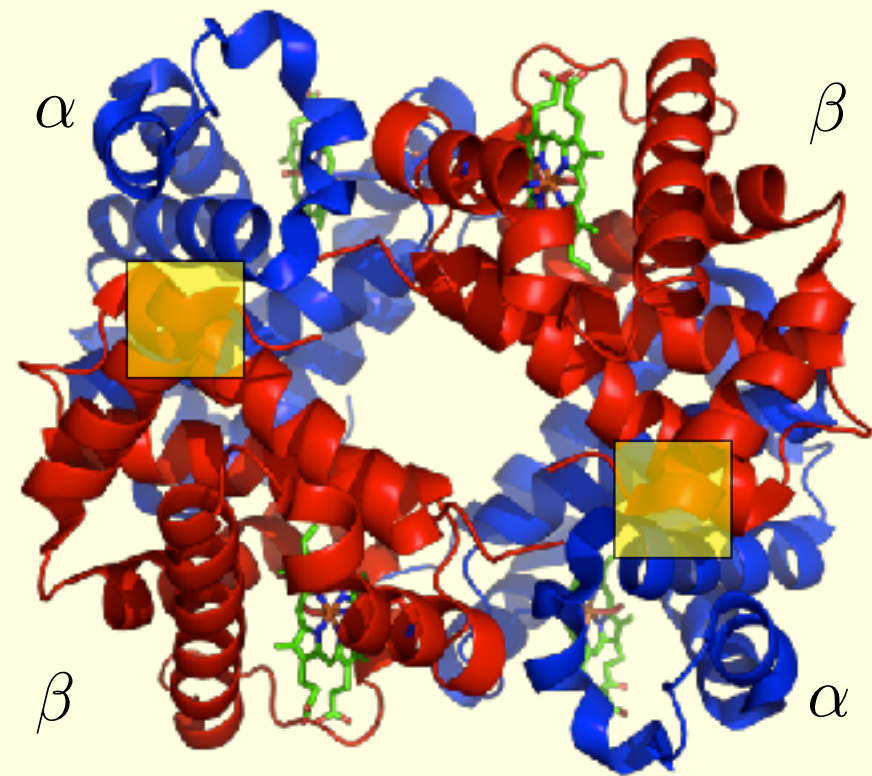
- Discontiguous by necessity  
Often connected by alpha-helix
- Less twisted than antiparallel sheets



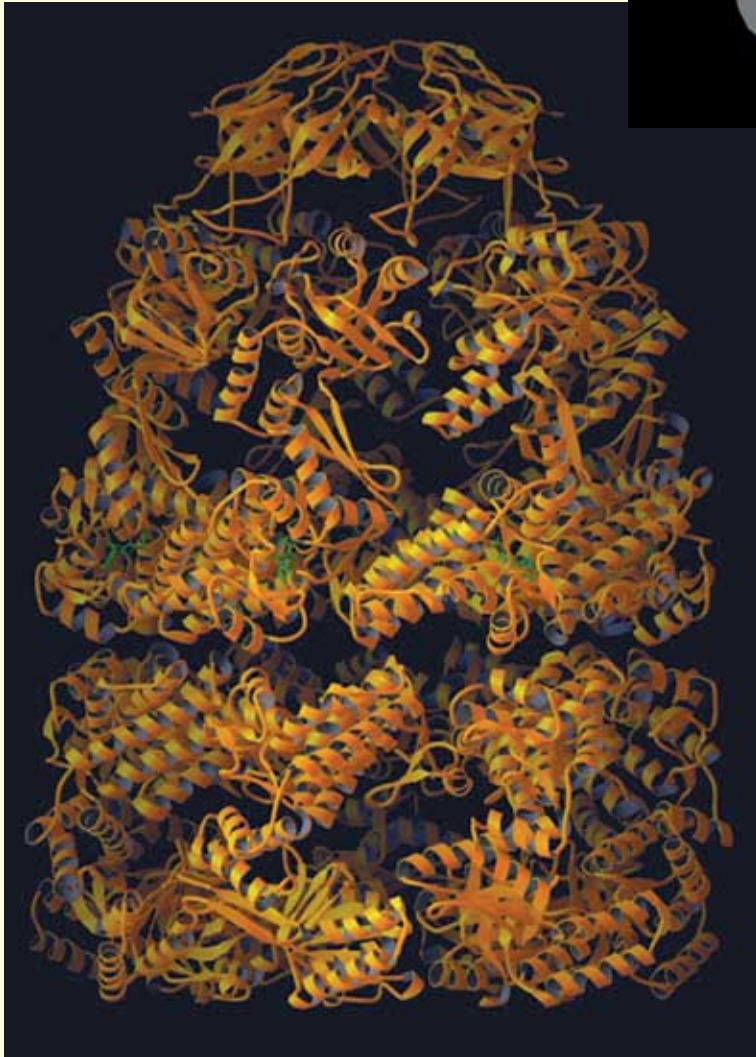
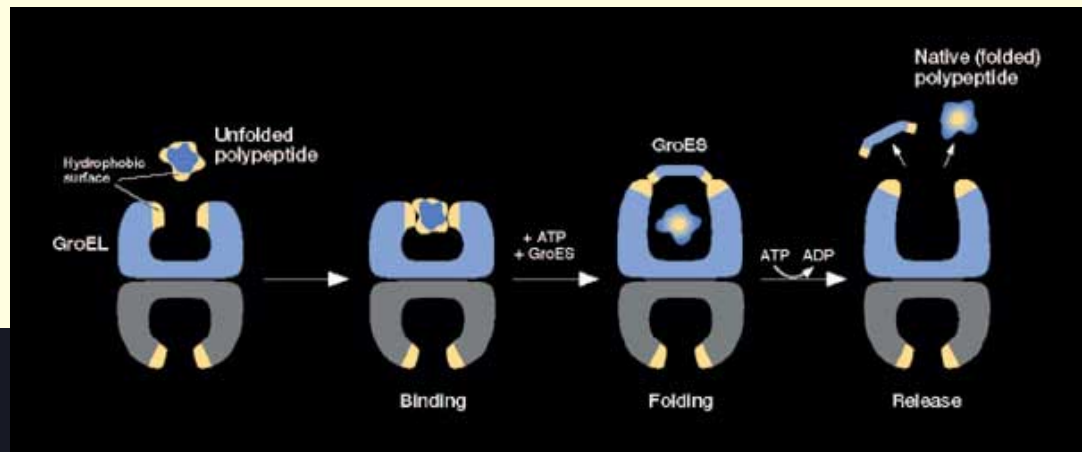
$\beta$ - $\alpha$ - $\beta$  Loop



Deoxyhemoglobin



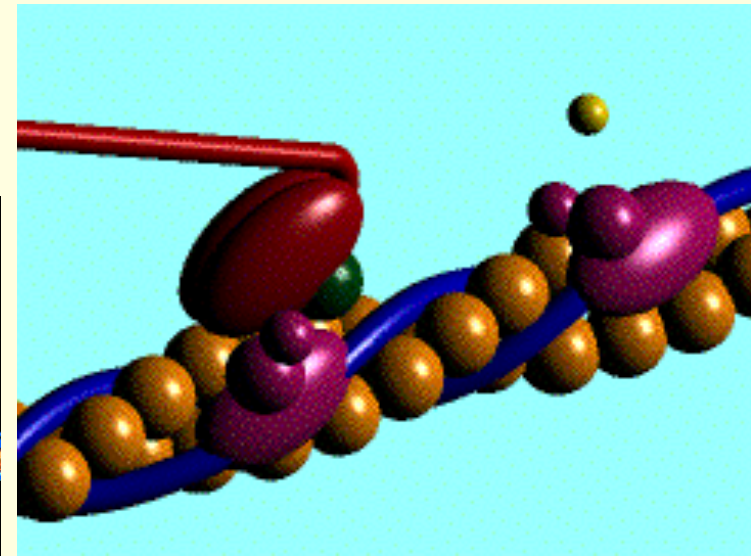
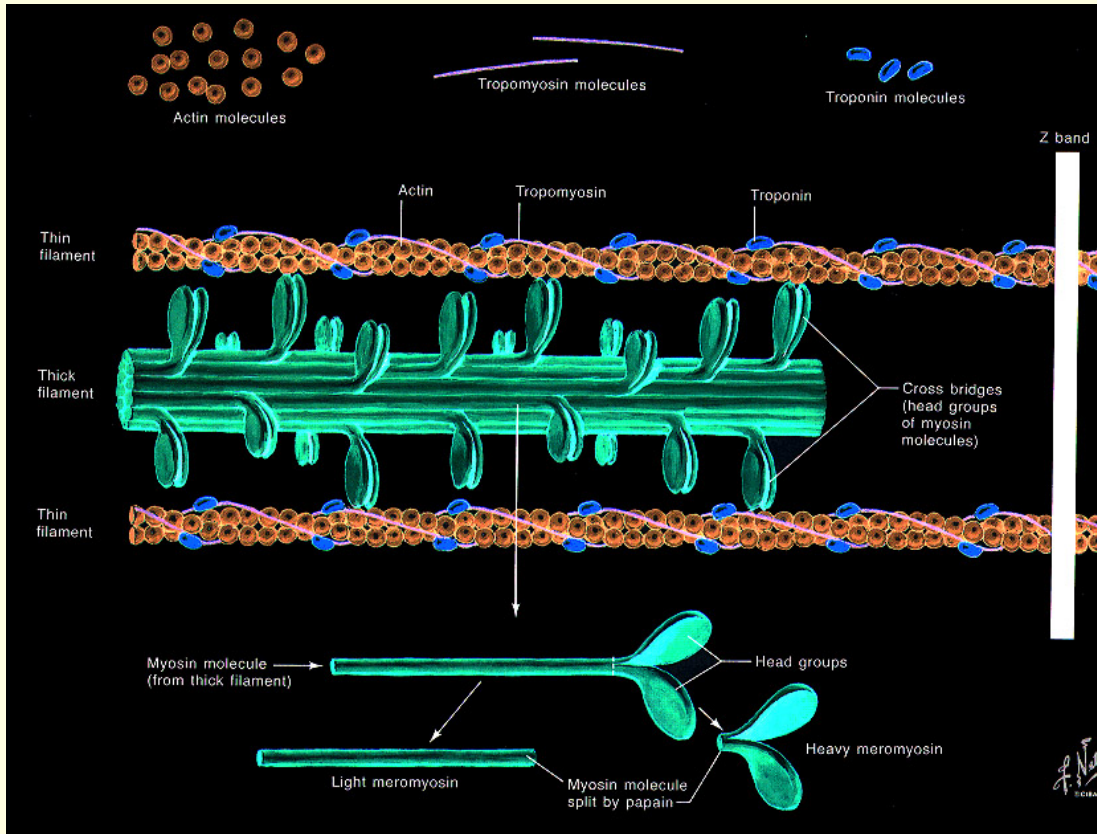
A GLU to VAL mutation at 6th amino acid in the  $\beta$ -subchains causes hemoglobin to aggregate, resulting in sickle-cell anemia.



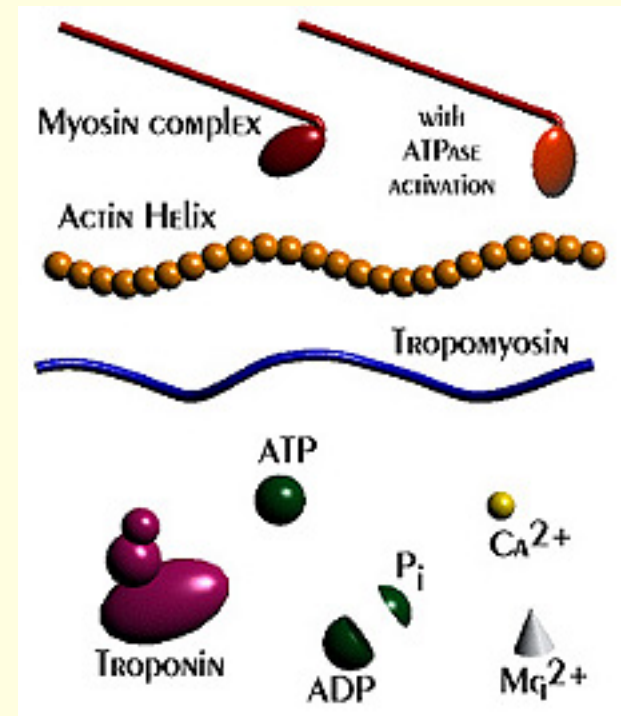
Some proteins need “help” during the folding process from “chaperonins”.

GroEL-GroES complex (Horwich *et al.*)

# Actin and Myosin



[www.sci.sdsu.edu/movies/actin\\_myosin.html](http://www.sci.sdsu.edu/movies/actin_myosin.html)



Titin is the largest known protein, with 38,138 residues (4200 kDA).

# “Computing” Protein Structure

- Molecular Dynamics Simulations
- *De novo* Prediction
- Homology Modeling
- X-ray diffraction
- Nuclear Magnetic Resonance Spectroscopy



# Molecular Dynamics

- We can also use folding techniques to study proteins of known structure.
- By simulating physics/chemistry, we can study normal modes, binding, interaction, etc.

<http://www.stanford.edu/group/pandegroup/folding/villin/>

Formation of a 30-residue  $\alpha$ -helix

**Experimental Methods**

XRC / NMR

**Predictive Methods**

Homology modeling

**Redesign**

**Comparison & Analysis**

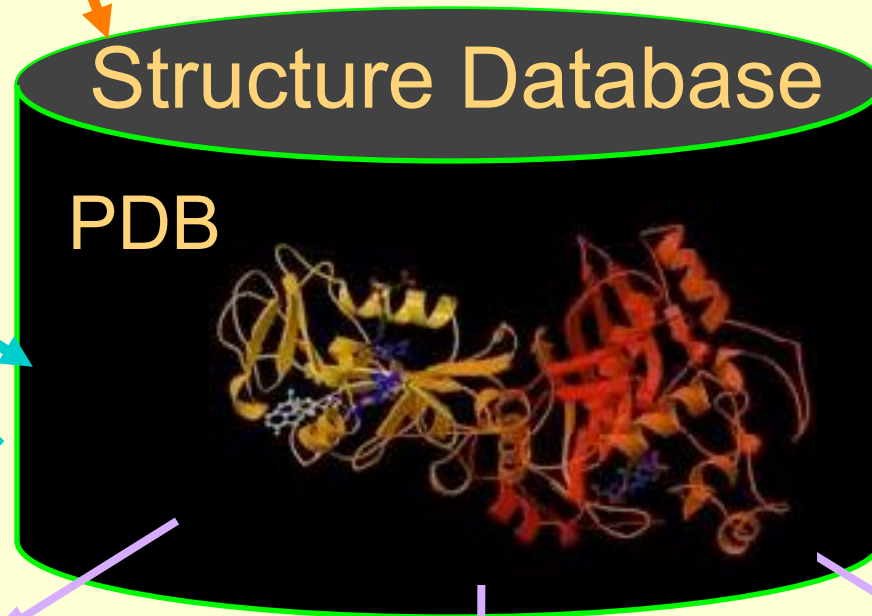
Mutation prediction,  
Structural Homology

**Function**

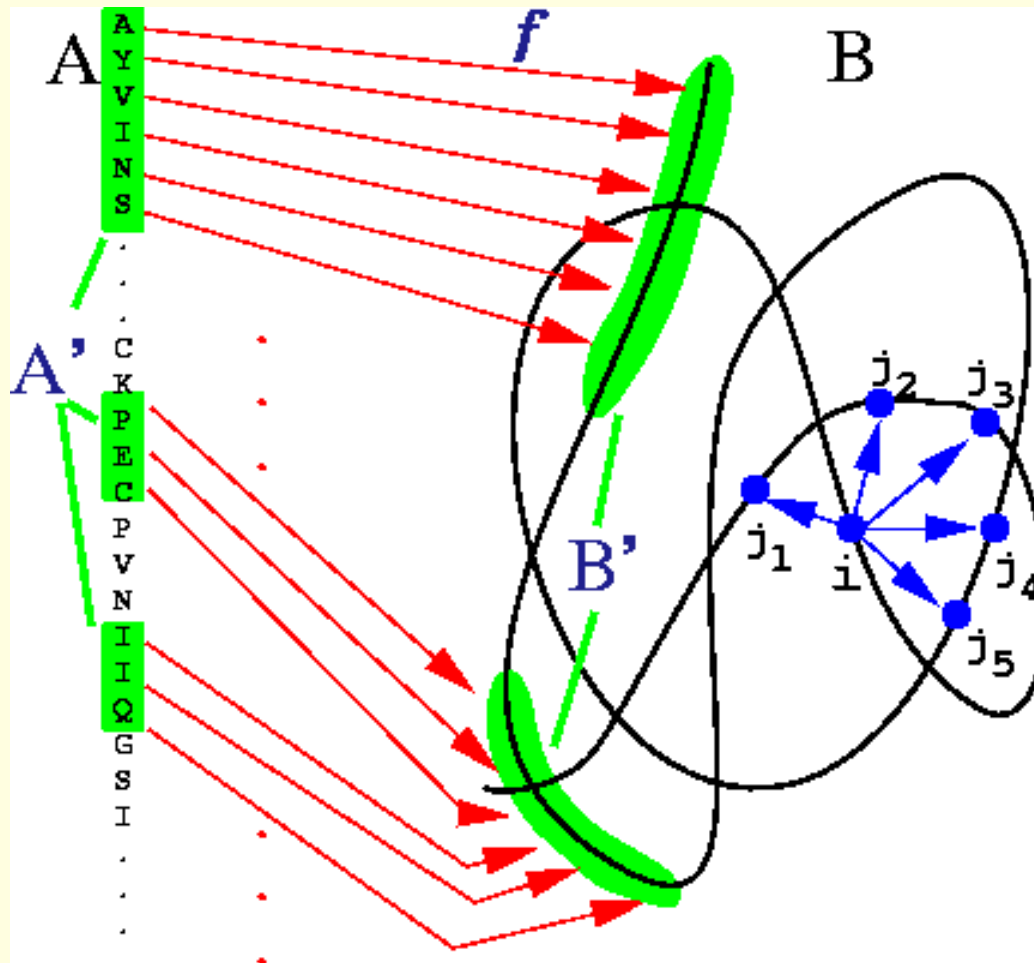
Catalysis, Binding,  
Dynamics

**Drug Design**

Docking, de novo,  
Database Search



# Leveraging Existing Structures



primary sequence

known structure

Input:

- 1) Primary sequence P for unknown structure.
- 2) Known homologous structure.

Output:

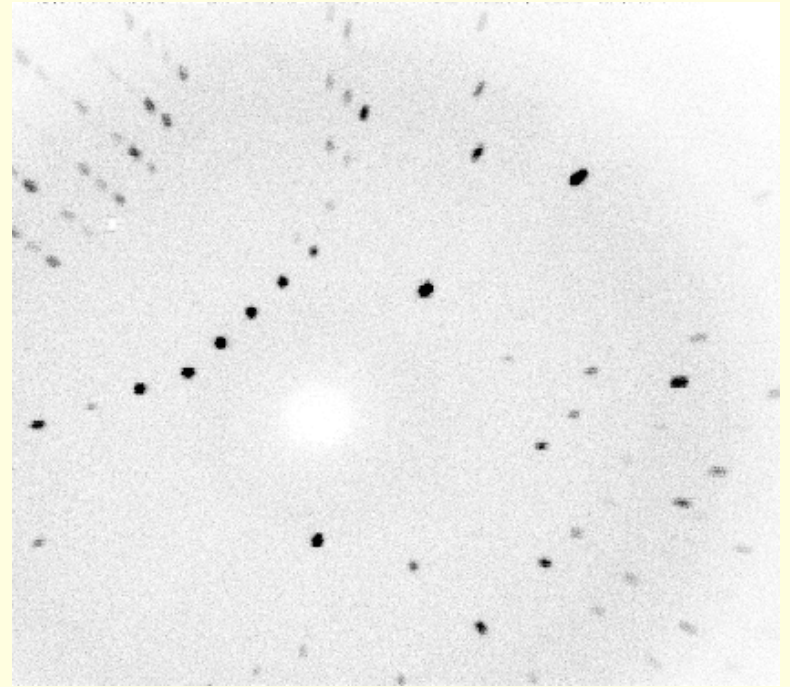
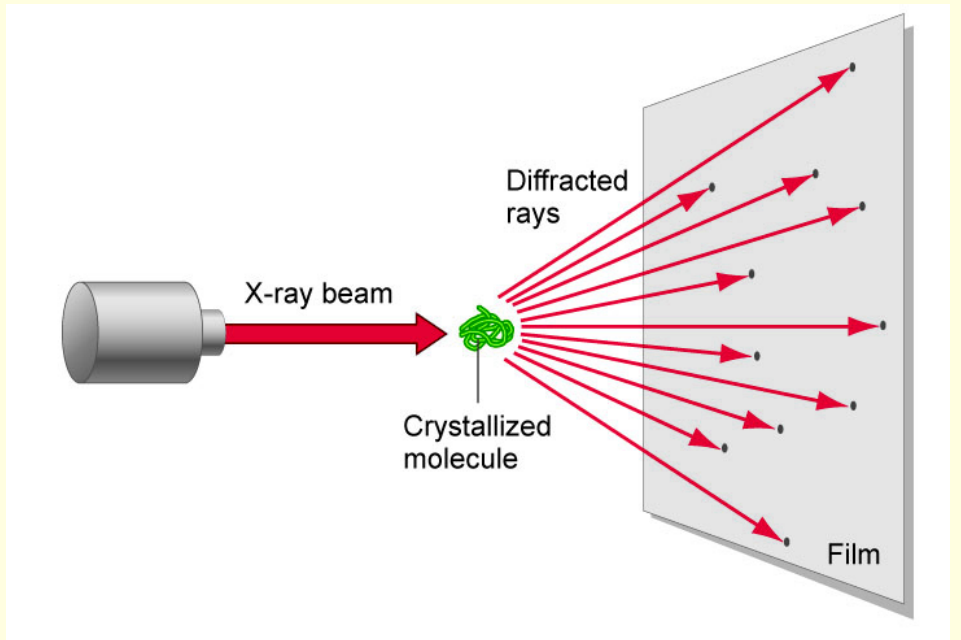
3D structure of A.

Goal:

Minimize energy of A.

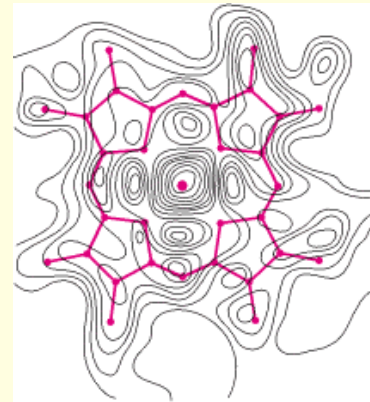
**Threading is NP-hard!**

# X-ray Crystallography



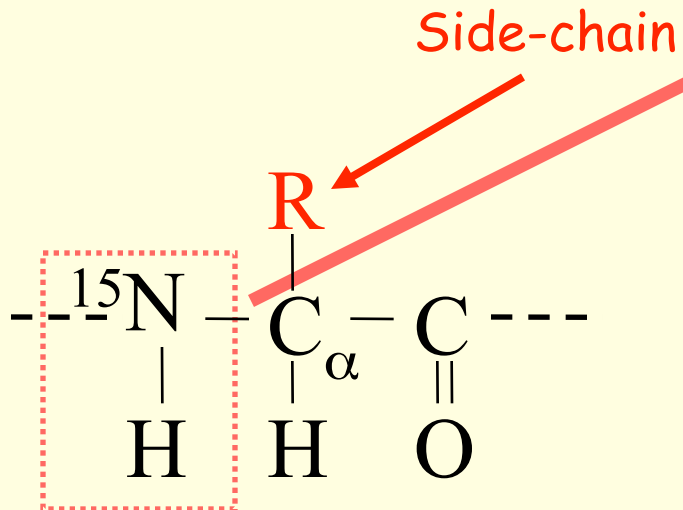
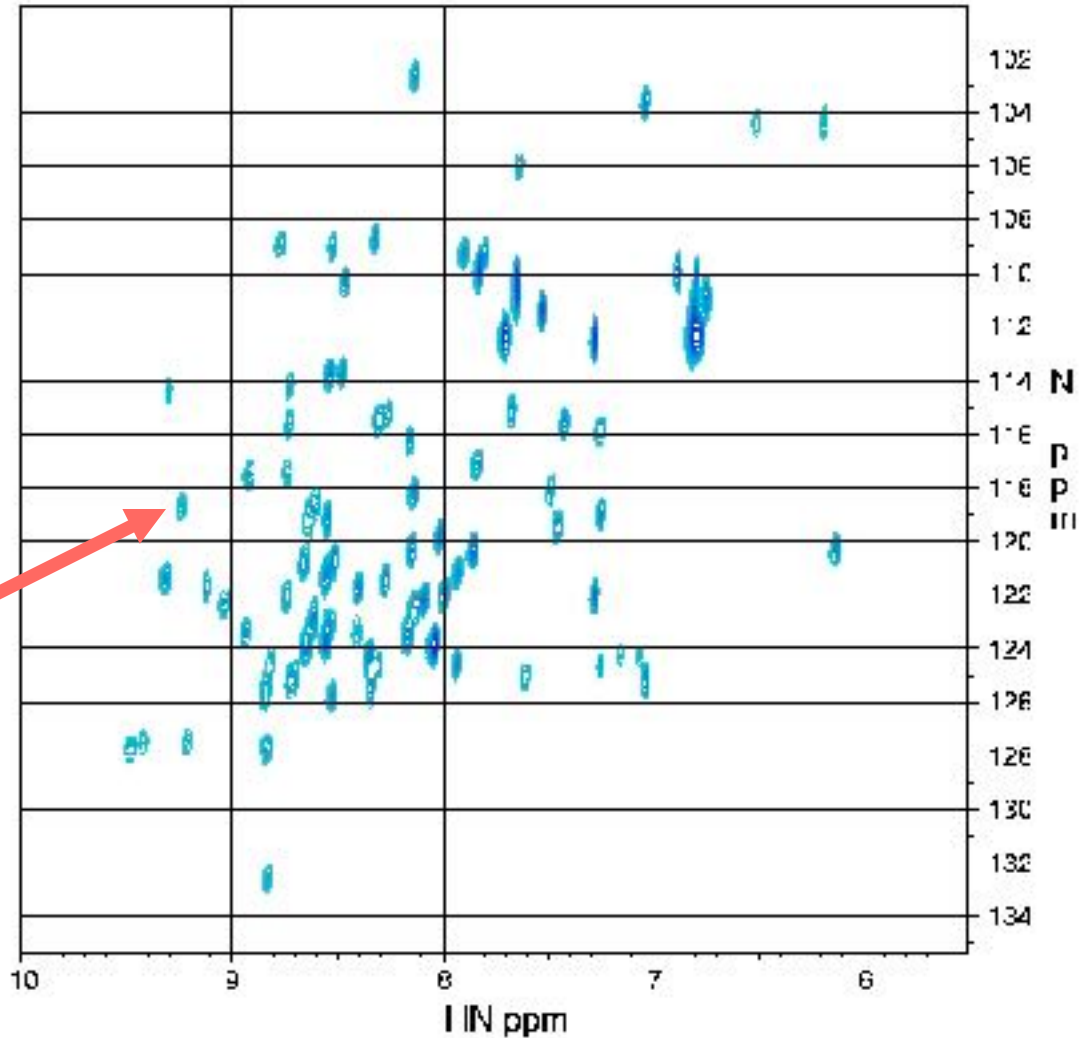
The diffraction pattern is the Fourier transform of electron density. To compute electron density, we can “invert” the diffraction pattern.

But we cannot measure the **phase** of diffracted wave! We must simulate it...



Myoglobin structure was solved by X-ray crystallography (Perutz and Kendrew 1958)

# NMR Spectroscopy

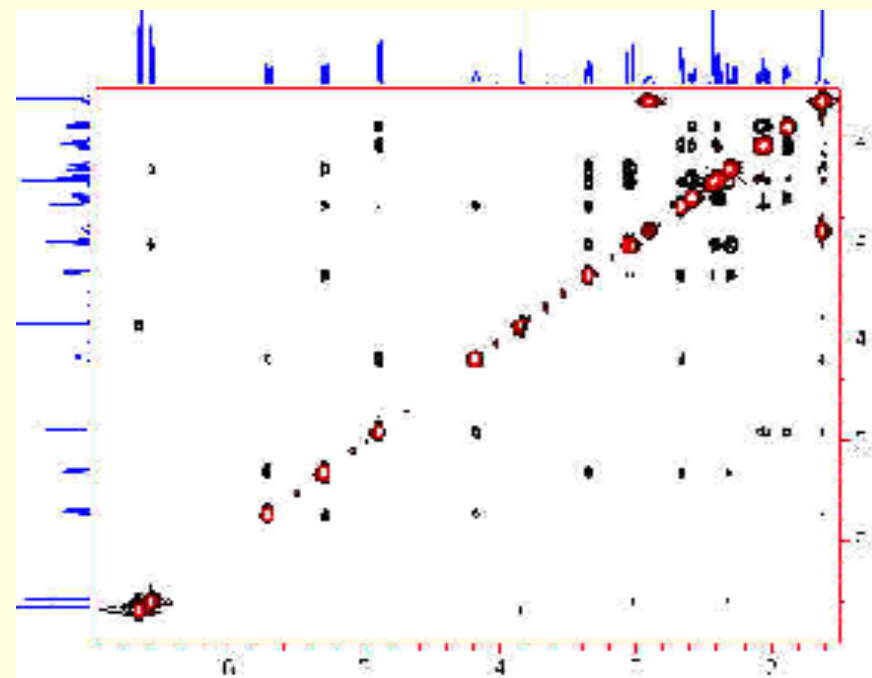


Amino Acid

Ubiquitin:  $\text{H}^{\text{N}} - ^{15}\text{N}$  HSQC

# NMR Structure Determination

Distance-based methods use *assigned* NOE restraints.



2D NOE spectrum



|       | $a_1$ | $a_2$ | $a_3$ | ... | $a_n$ |
|-------|-------|-------|-------|-----|-------|
| $a_1$ | 0     | ?     | <6    | ... |       |
| $a_2$ | ?     | 0     | <6    | ... |       |
| $a_3$ | ?     | ?     | 0     | ... |       |
| ...   | ...   | ...   | ...   | ... | ...   |
| $a_n$ | ...   | ...   | ...   | ... | ...   |

# NMR Structure Determination

Structure determination from NOE restraints is **NP-hard!** [Saxe '79; Hendrickson '92, '95]

|          | $a_1$    | $a_2$    | $a_3$    | $\dots$  | $a_n$ |
|----------|----------|----------|----------|----------|-------|
| $a_1$    | 0        | ?        | <6       | $\dots$  |       |
| $a_2$    | ?        | 0        | <6       | $\dots$  |       |
| $a_3$    | ?        | ?        | 0        | $\dots$  |       |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |       |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |       |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |       |
| $a_n$    |          |          |          |          |       |

Distance Geometry Method  
[Crippen&Havel'88]

Exponential Time

