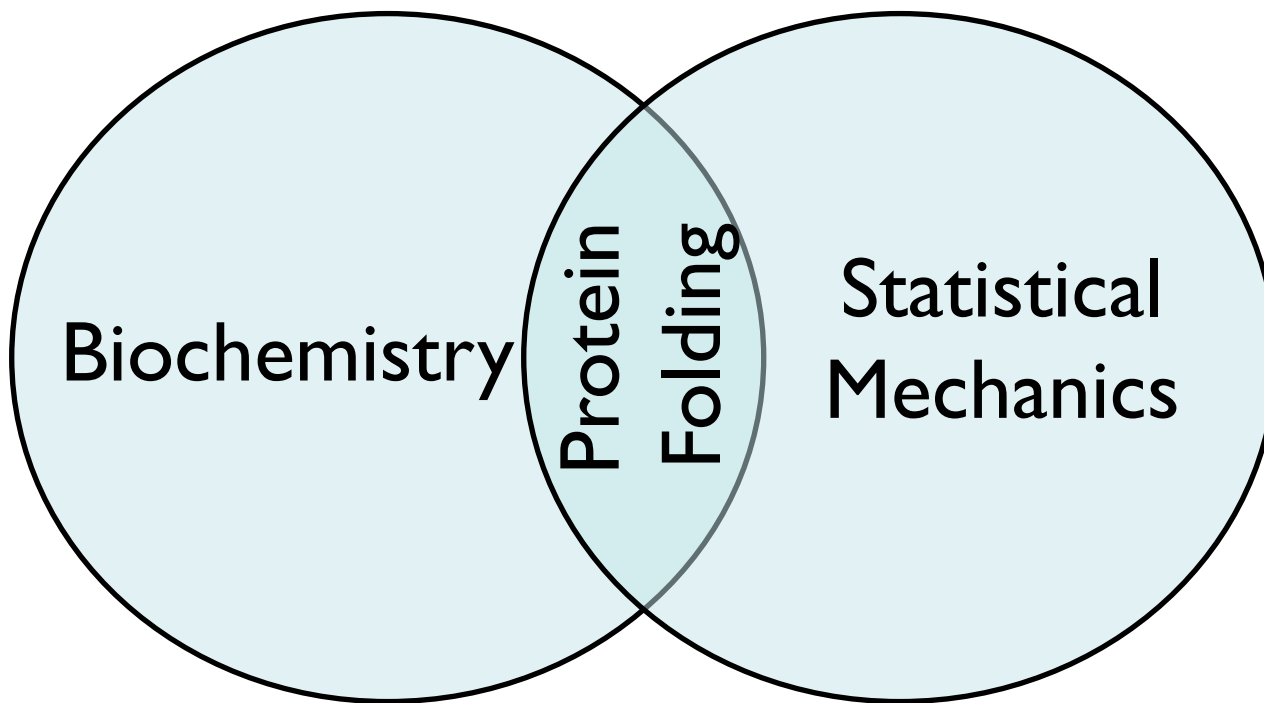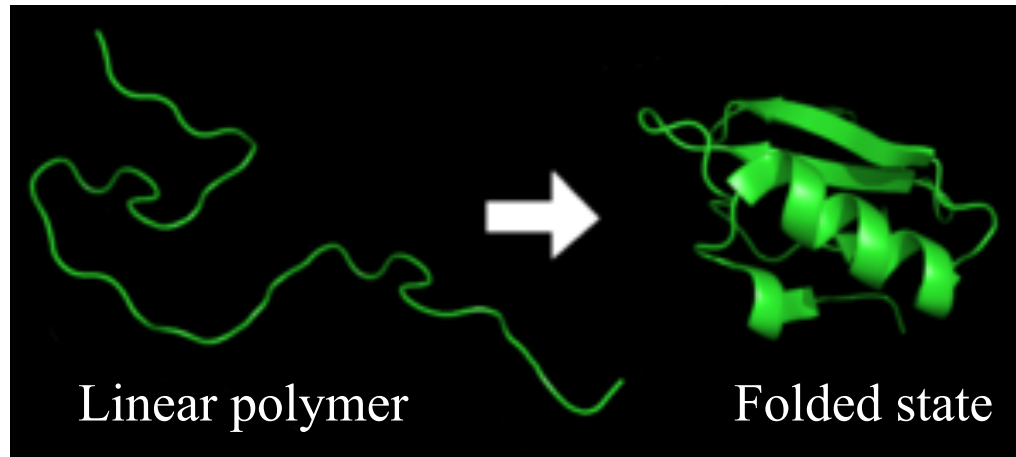# Bioinformatics: Practical Application of Simulation and Data Mining

# Protein Folding I

Prof. Corey O'Hern
Department of Mechanical Engineering & Materials Science
Department of Physics
Department of Applied Physics
Program in Computation Biology & Bioinformatics
Integrated Graduate Program in Physical & Engineering Biology
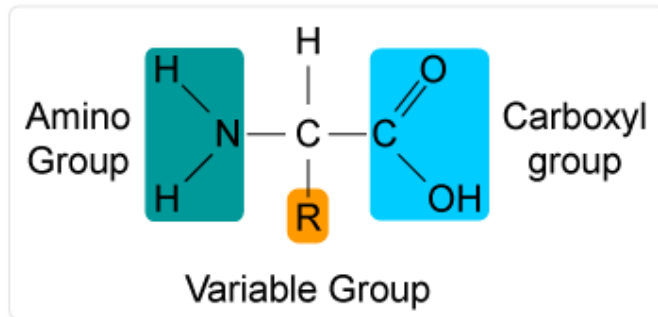Yale University

1

# What are proteins?



Linear polymer → Folded state

- Proteins are important; e.g. for catalyzing and regulating biochemical reactions, transporting molecules, …
- Linear polymer chain composed of tens (peptides) to thousands (proteins) of monomers
- Monomers are 20 naturally occurring amino acids
- Different proteins have different amino acid sequences
- *Structureless*, extended unfolded state
- Compact, 'unique' native folded state (with secondary and tertiary structure) required for biological function
- Sequence determines protein structure (or lack thereof)
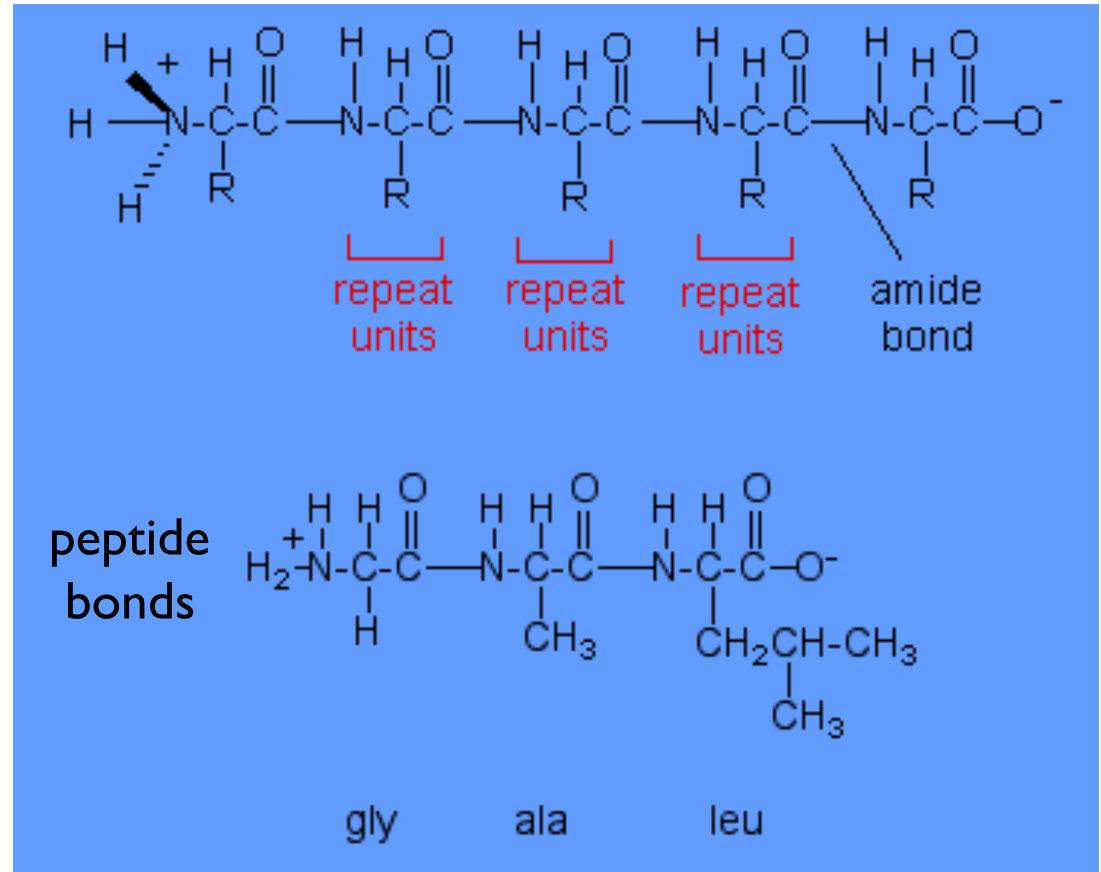- Proteins unfold or denature with increasing temperature or chemical denaturants

3
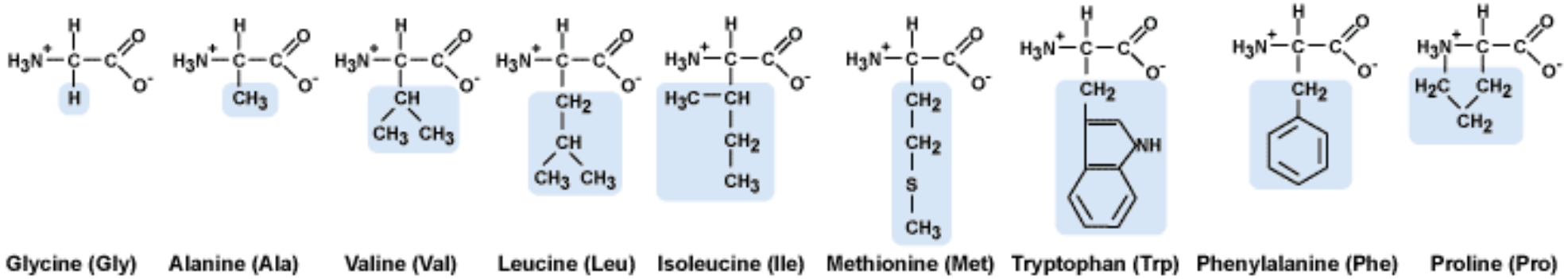
# Amino Acids I

## General structure of Amino Acids

Amino Group

Carboxyl group

Variable Group

N-terminal    $C_\alpha$    C-terminal

R
variable
side chain

peptide
bonds

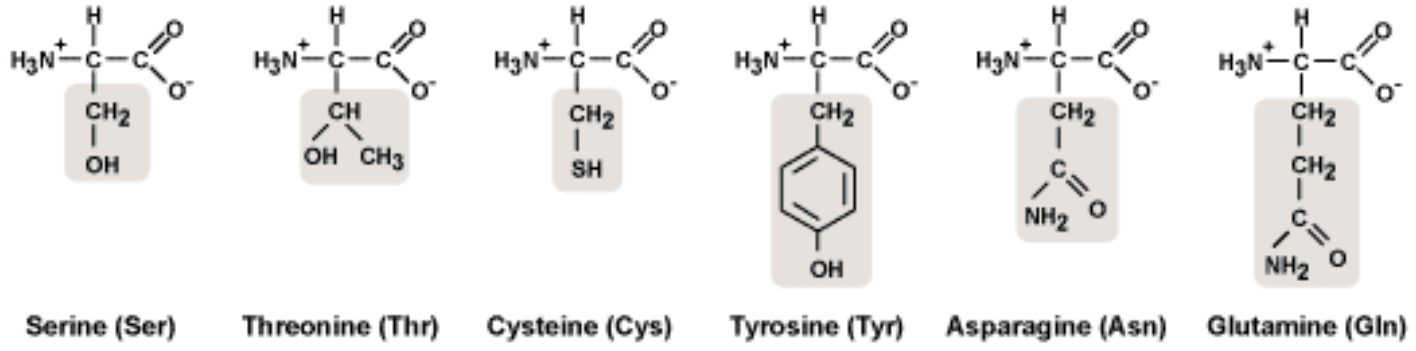repeat units    repeat units    repeat units    amide bond

gly    ala    leu

- Side chains differentiate amino acid repeat units
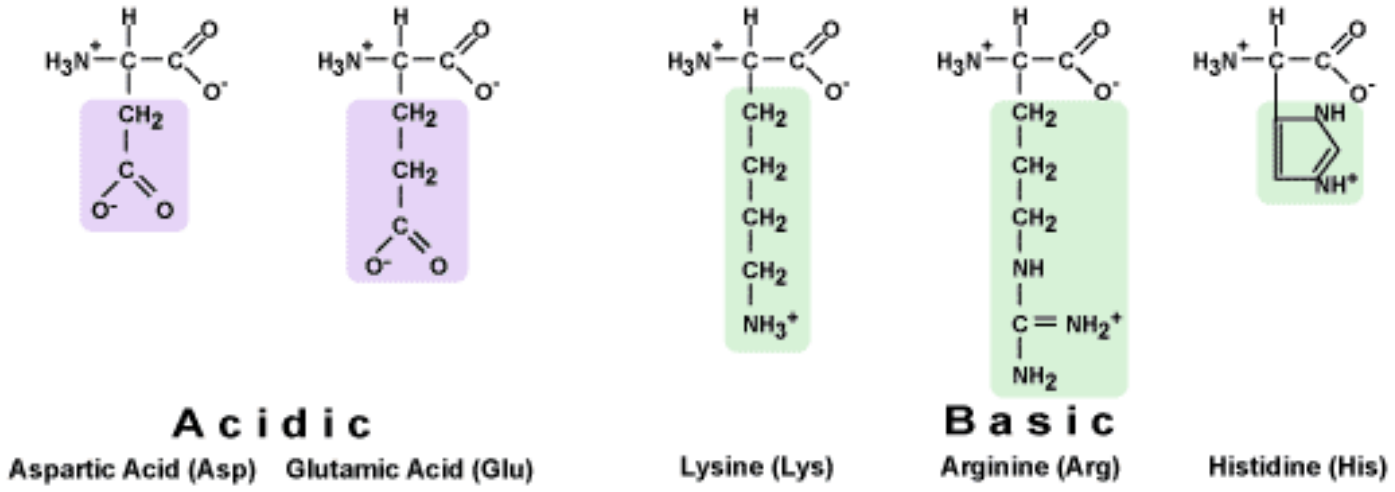- Peptide bonds link residues into polypeptides

4

# Amino Acids II



**NONPOLAR**

Glycine (Gly)  Alanine (Ala)  Valine (Val)  Leucine (Leu)  Isoleucine (Ile)  Methionine (Met)  Tryptophan (Trp)  Phenylalanine (Phe)  Proline (Pro)

**POLAR**

Serine (Ser)  Threonine (Thr)  Cysteine (Cys)  Tyrosine (Tyr)  Asparagine (Asn)  Glutamine (Gln)

**Electrically Charged**

**Acidic**

Aspartic Acid (Asp)  Glutamic Acid (Glu)

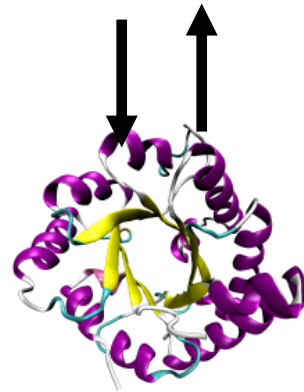**Basic**

Lysine (Lys)  Arginine (Arg)  Histidine (His)

5

# The Protein Folding Problem:

What is 'unique' folded 3D structure of a protein based on its amino acid sequence?                    Sequence → Structure

Lys–Asn–Val–Arg–Ser–Lys–Val–Gly–Ser–Thr–Glu–Asn–Ile–Lys– His–Gln–Pro– Gly–Gly–Gly–...

# Driving Forces

- Folding: hydrophobicity, hydrogen bonding, van der Waals interactions, …
- Unfolding: increase in conformational entropy, electric charge…
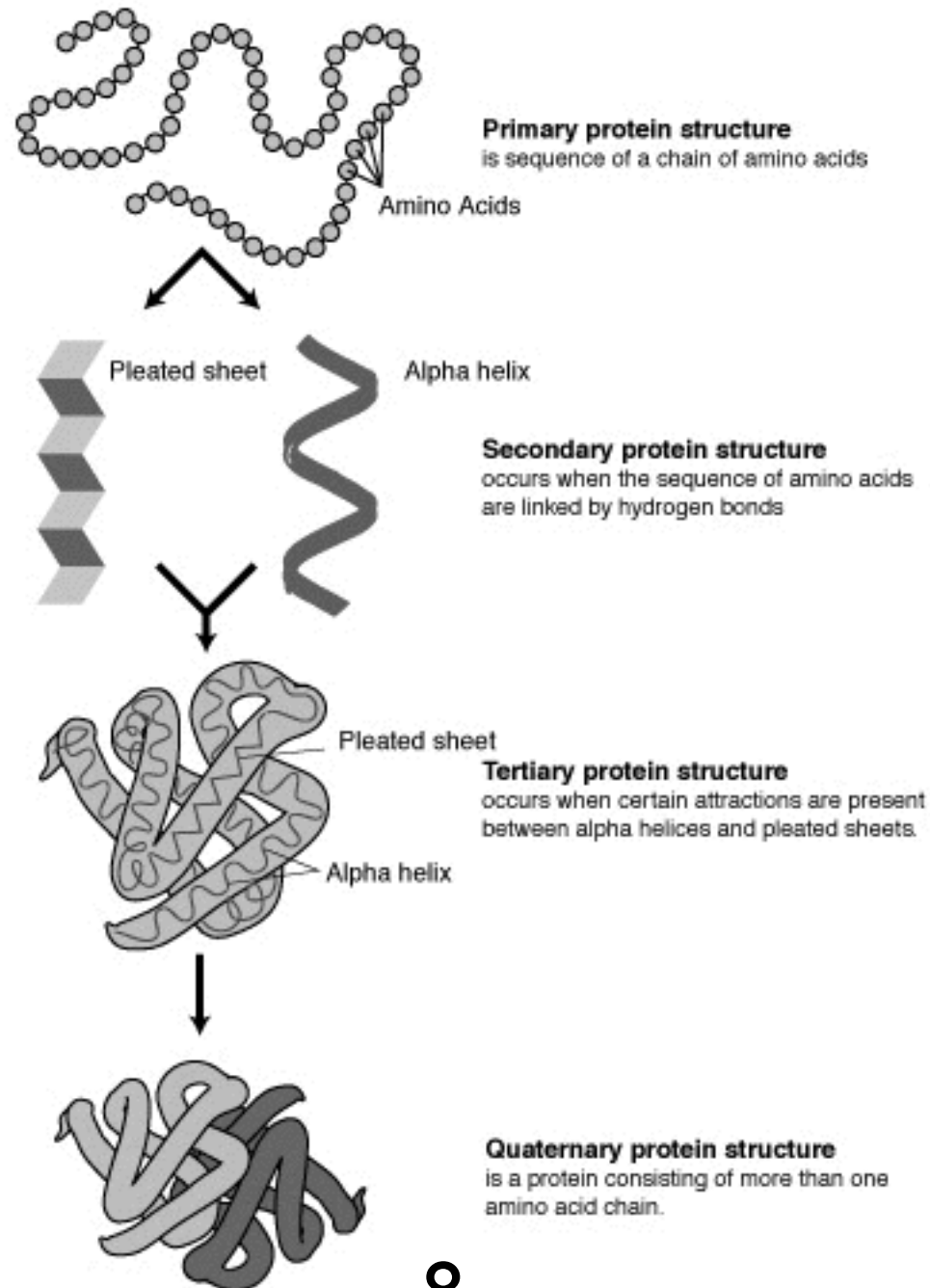
inside     H (hydrophobic)

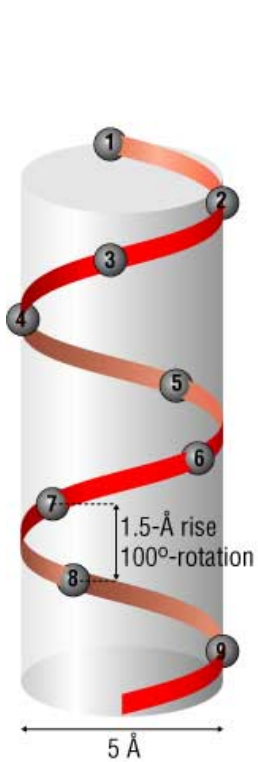outside    P (polar)

### Hydrophobicity index

| At pH 2[A] | | At pH 7[B] | |
|---|---|---|---|
| **Very Hydrophobic** | | | |
| Leu | 100 | Phe | 100 |
| Ile | 100 | Ile | 99 |
| Phe | 92 | Trp | 97 |
| Trp | 84 | Leu | 97 |
| Val | 79 | Val | 76 |
| Met | 74 | Met | 74 |
| **Hydrophobic** | | | |
| Cys | 52 | Tyr | 63 |
| Tyr | 49 | Cys | 49 |
| Ala | 47 | Ala | 41 |
| **Neutral** | | | |
| Thr | 13 | Thr | 13 |
| Glu | 8 | His | 8 |
| Gly | 0 | Gly | 0 |
| Ser | -7 | Ser | -5 |
| Gln | -18 | Gln | -10 |
| Asp | -18 | | |
| **Hydrophilic** | | | |
| Arg | -26 | Arg | -14 |
| Lys | -37 | Lys | -23 |
| Asn | -41 | Asn | -28 |
| His | -42 | Glu | -31 |
| Pro | -46 | Pro | -46 (used pH 2) |
| | | Asp | -55 |

[A]pH 2 values: Normalized from Sereda et al., J. Chrom. 676: 139-153 (1994).

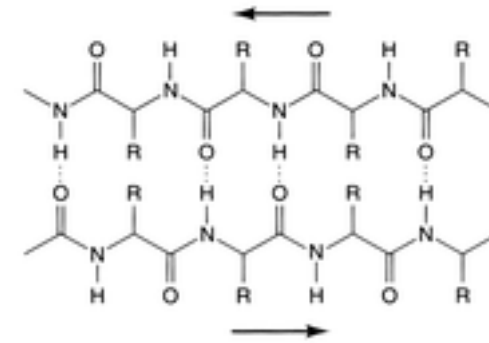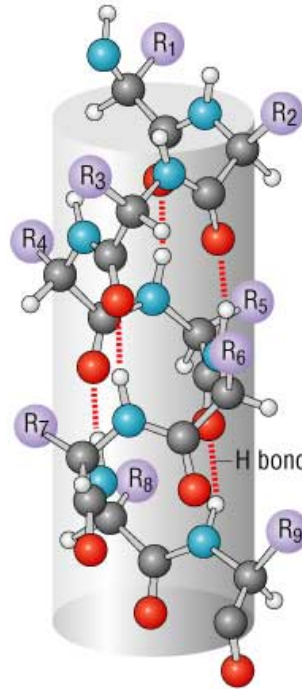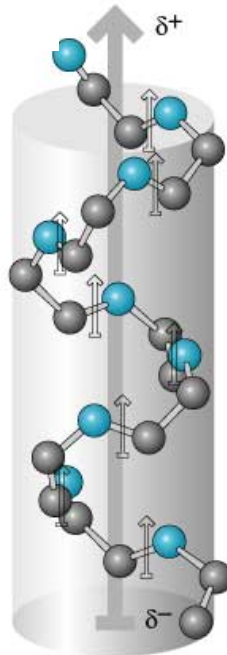[B]pH 7 values: Monera et al., J. Pept. Sci. 1: 319-329 (1995).
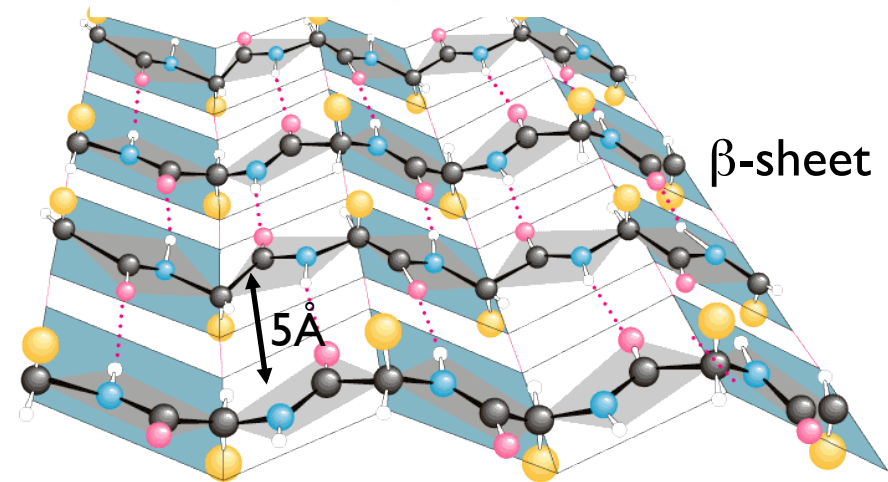
7

# Higher-order Structure

**Primary protein structure**
is sequence of a chain of amino acids

Amino Acids

Pleated sheet          Alpha helix

**Secondary protein structure**
occurs when the sequence of amino acids
are linked by hydrogen bonds

Pleated sheet

**Tertiary protein structure**
occurs when certain attractions are present
between alpha helices and pleated sheets.

Alpha helix

**Quaternary protein structure**
is a protein consisting of more than one
amino acid chain.

8

# Secondary Structure: Loops, α-helices, β-strands/sheets
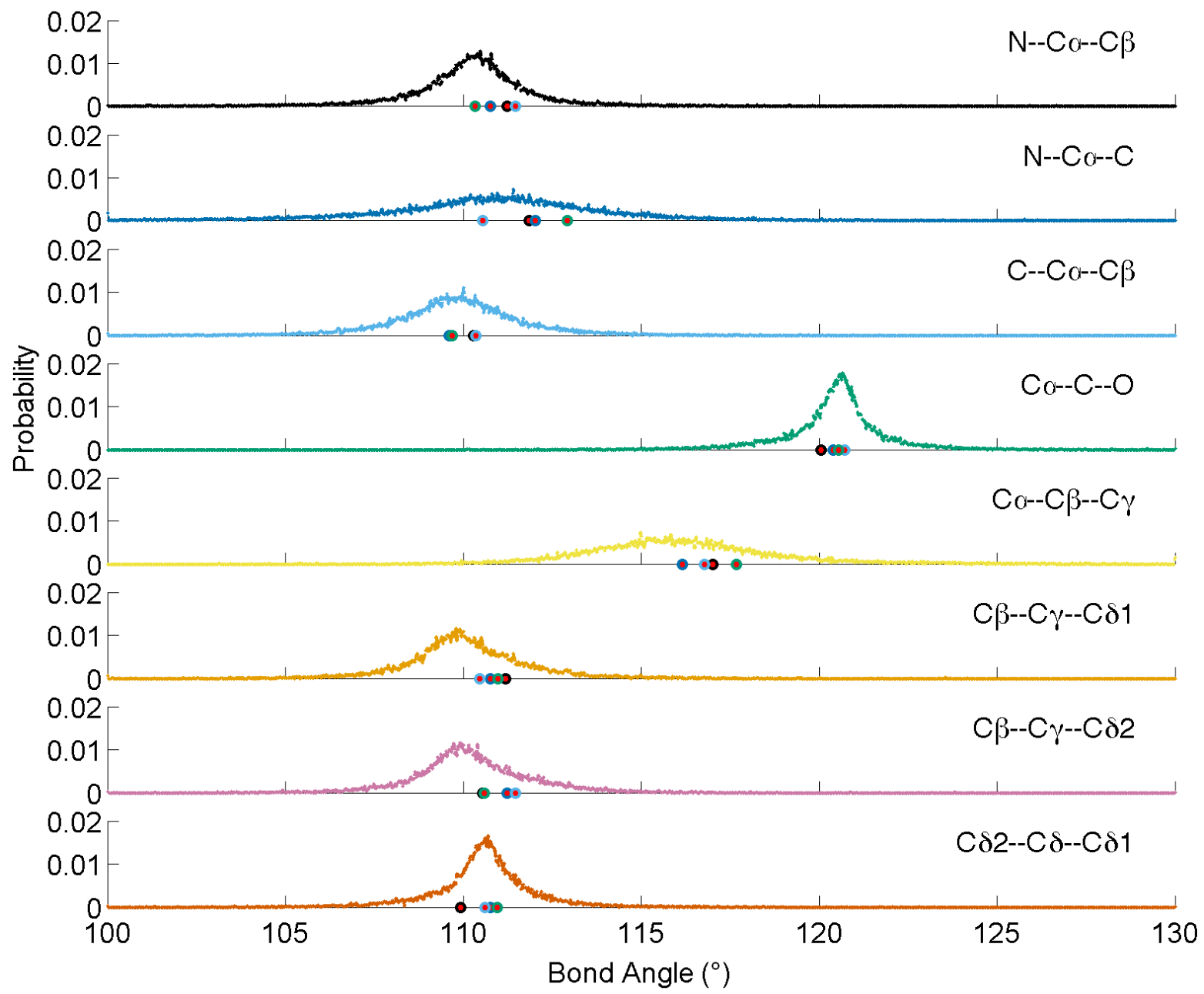
## α-helix



β-strand

β-sheet

•Right-handed; three turns
•Vertical hydrogen bonds between $NH_2$ (teal/white) backbone group and C=O (grey/red) backbone group four residues earlier in sequence
•Side chains (R) on outside; point upwards toward $NH_2$
•Each amino acid corresponds to 100°, 1.5Å, 3.6 amino acids per turn
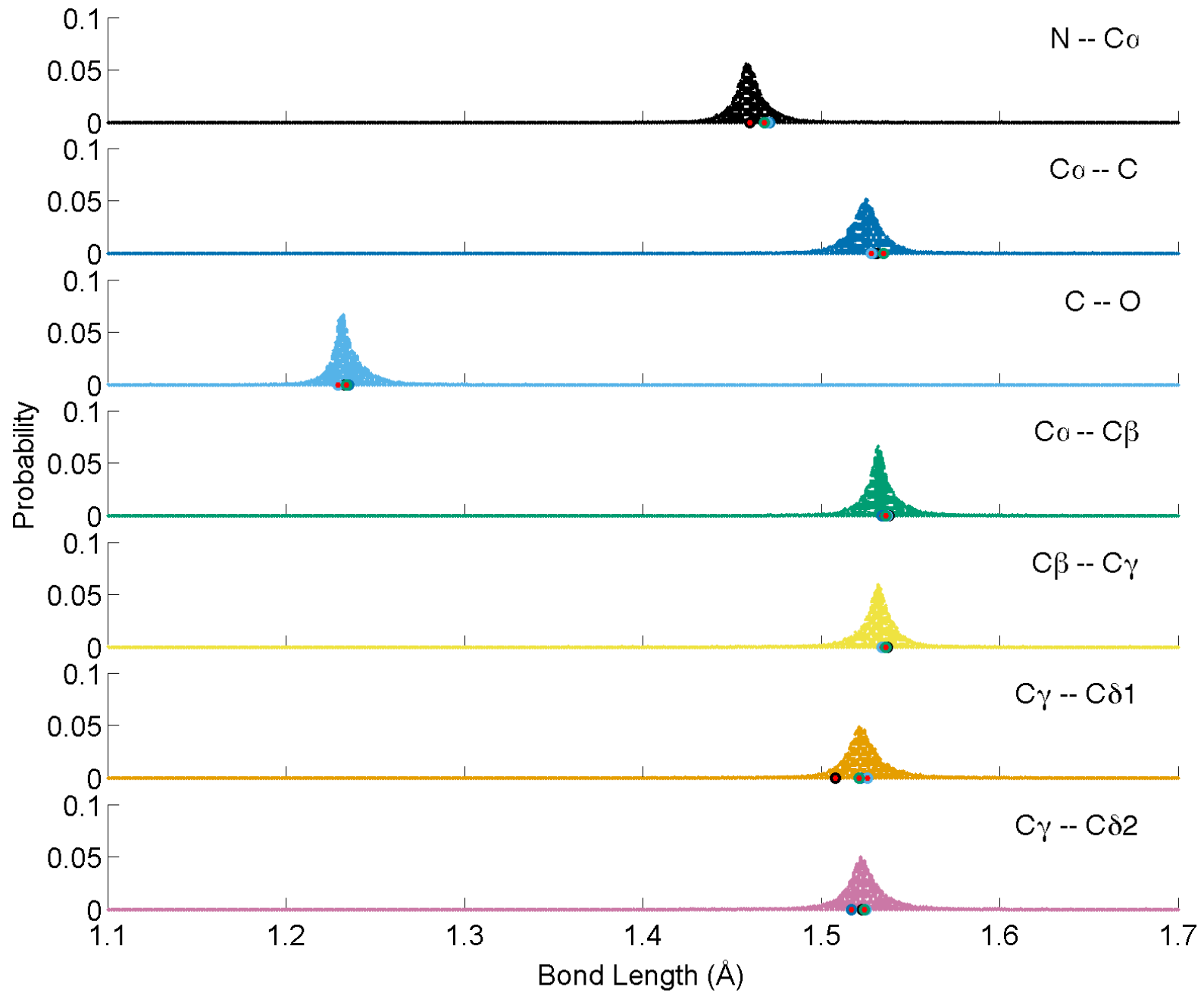•$(\phi,\psi)=(-60°,-45°)$
•α-helix propensities: Met, Ala, Leu, Glu

•5-10 residues; peptide backbones fully extended
•NH (blue/white) of one strand hydrogen-bonded to C=O (black/red) of another strand
•$C_\alpha$, side chains (yellow) on adjacent strands aligned; side chains along single strand alternate up and down
•$(\phi,\psi)=(-135°,135°)$
•β-strand propensities: Val, Thr, Tyr, Trp, Phe, Ile

9

# Backbonde Dihedral Angles



$$\cos\theta = \hat{\pi}_1 \bullet \hat{\pi}_2$$

10

N -- Cα

Cα -- C

C -- O

Cα -- Cβ
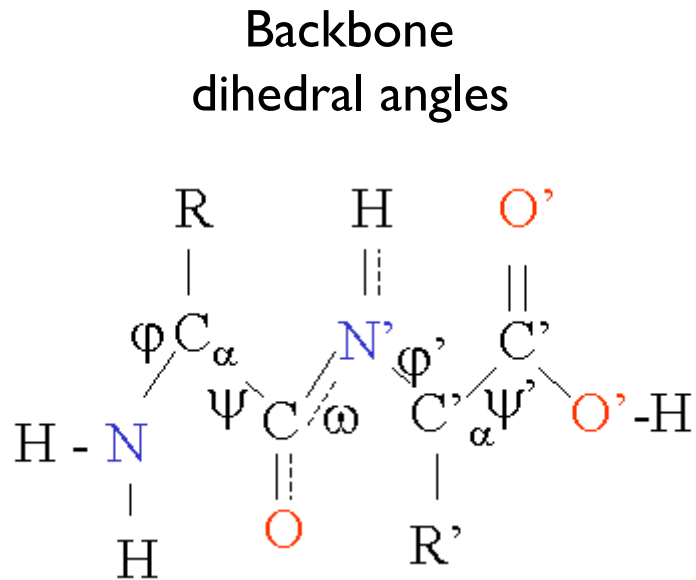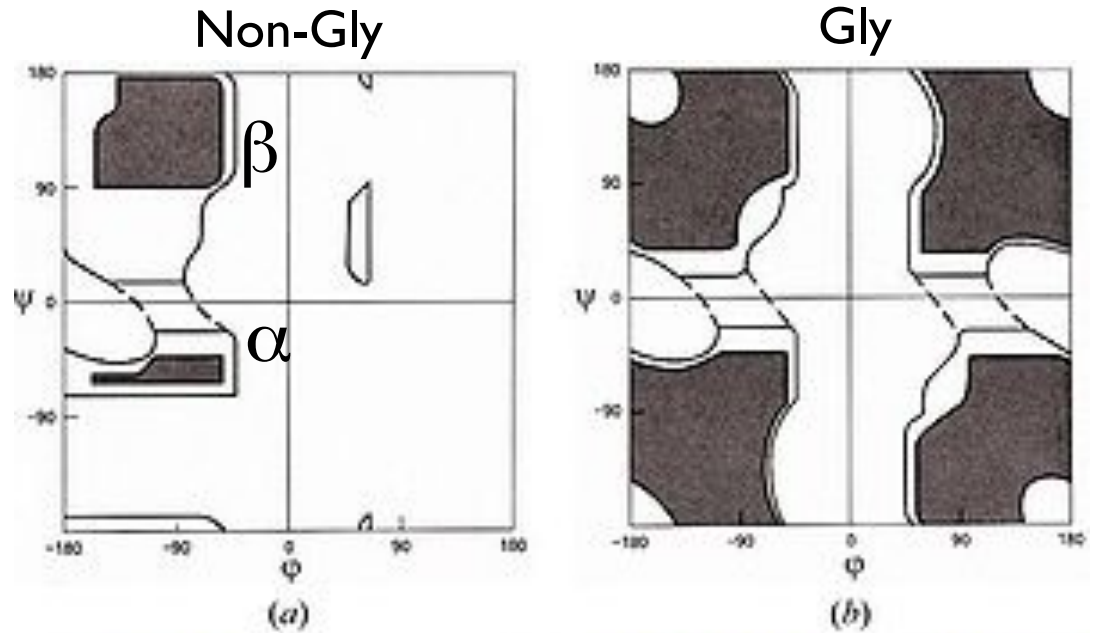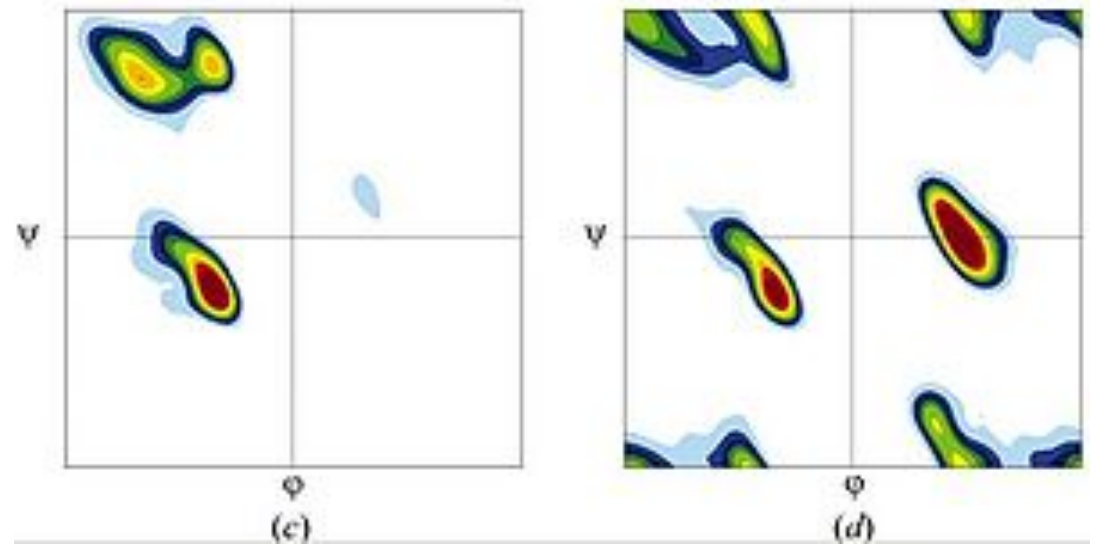
Cβ -- Cγ

Cγ -- Cδ1

Cγ -- Cδ2

Probability

Bond Length (Å)

# Ramachandran Plot: Determining Steric Clashes

Backbone dihedral angles



4 atoms define dihedral angle:

$CC_\alpha NC$     $\phi$

$C_\alpha N\,CC_\alpha$     $\omega=0,180°$
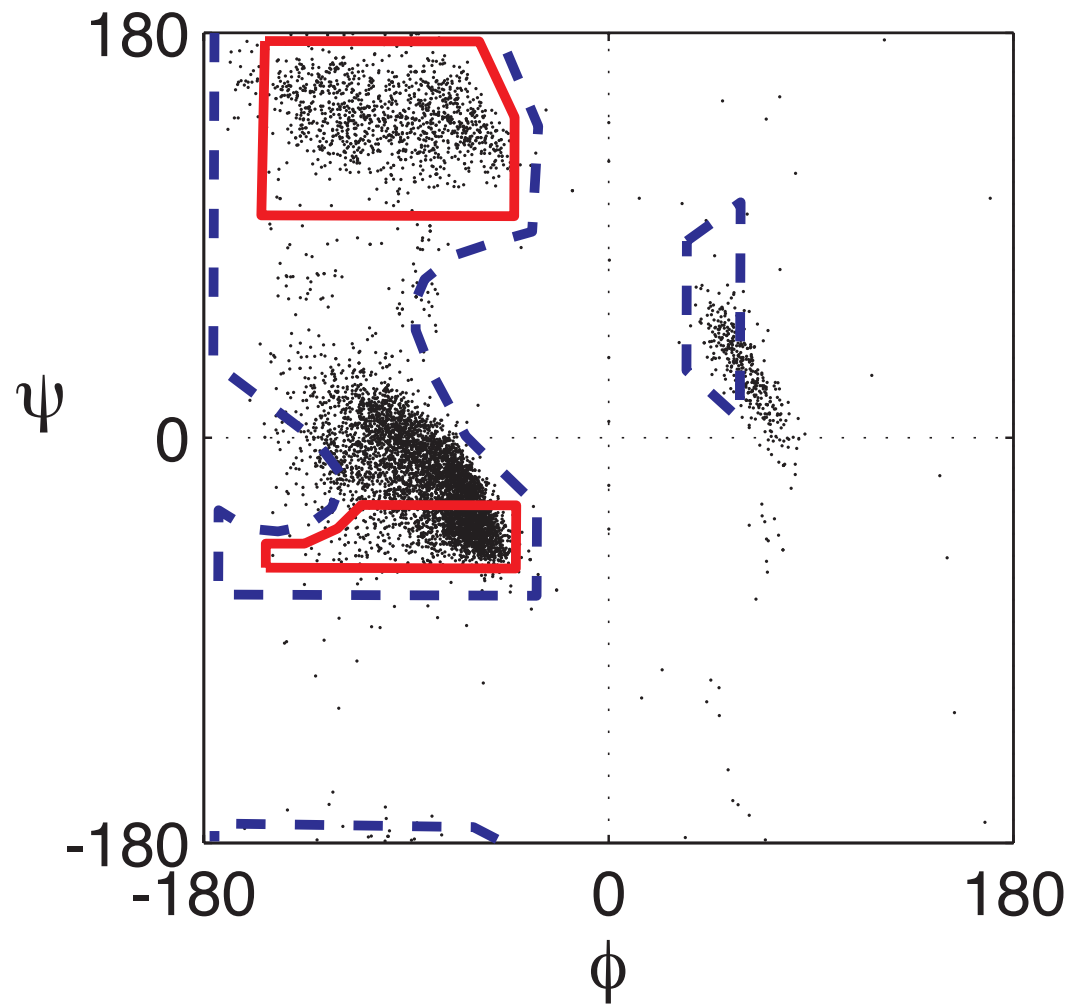
$N\,CC_\alpha N$     $\psi$

Non-Gly        Gly



theory

PDB

14

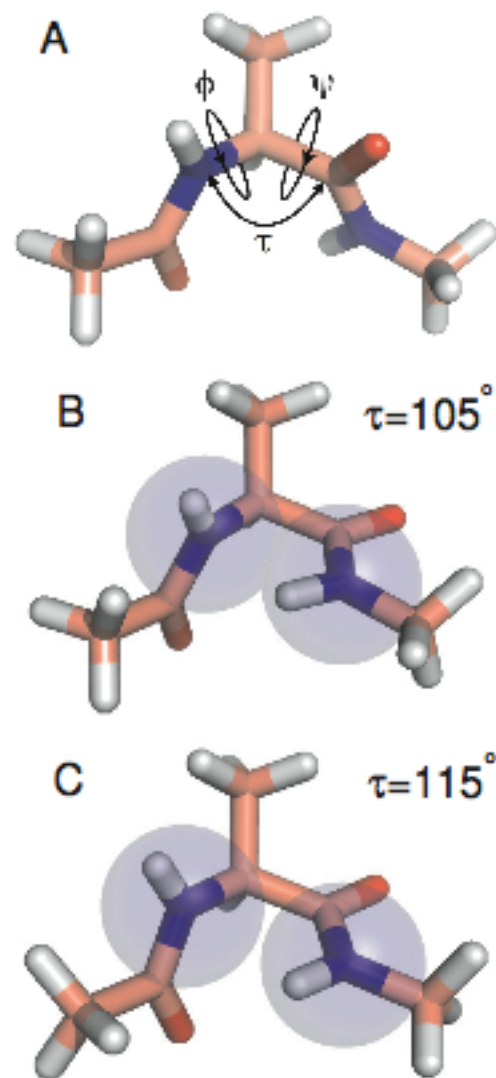⬜ vdW radii

— < vdW radii

--- backbone flexibility

# Backbone dihedral angles from PDB

Figure 1: Stick representation of an alanyl dipeptide mimetic. Atom types are color-coded: carbon=pink, nitrogen=blue, oxygen=red, hydrogen=white. **A**: The backbone dihedral angles $\phi$ and $\psi$ and the bond angle $\tau$ are indicated. **B**: $\tau = 105°$, $\phi = -90°$, $\psi = 0°$ (i.e. bridge region values of $\phi$ and $\psi$). Blue-shaded spheres indicate steric ove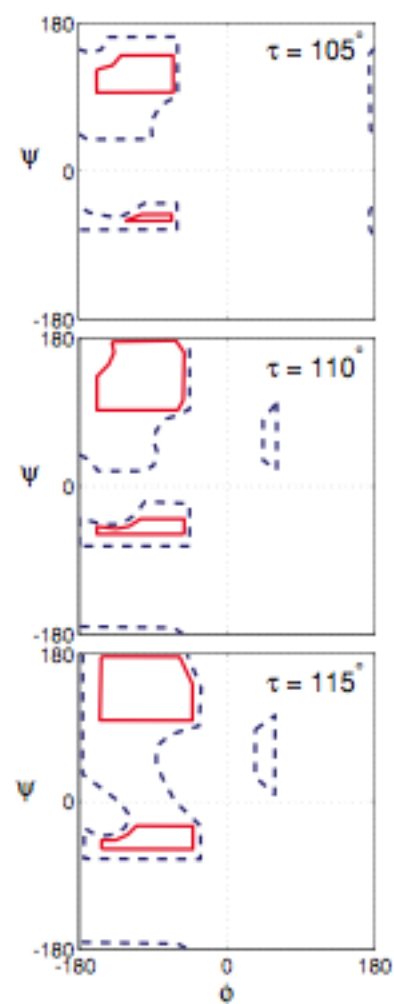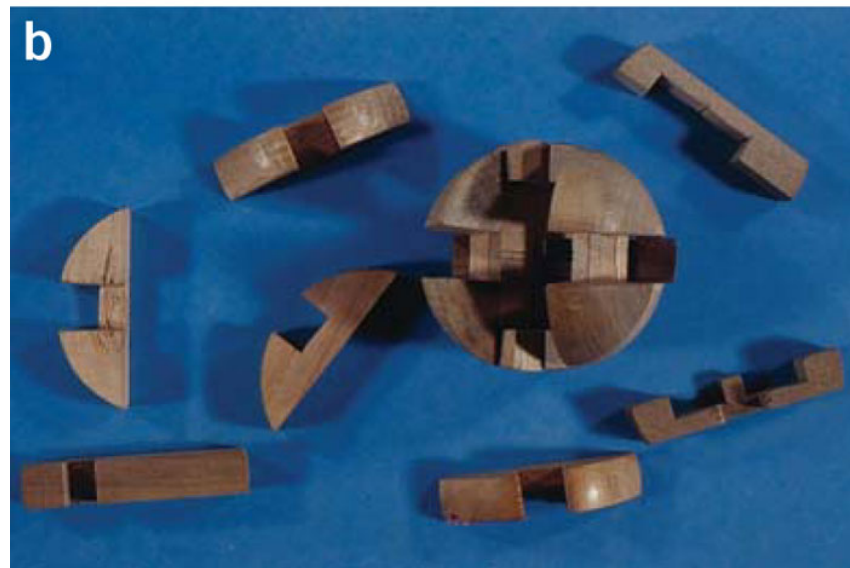rlap between main-chain nitrogens for this value of $\tau$. **C**: $\tau = 115°$, $\phi = -90°$, $\psi = 0°$ (i.e. bridge region values of $\phi$ and $\psi$). Blue-shaded spheres indicate no steric overlap between main-chain nitrogens for this value of $\tau$.
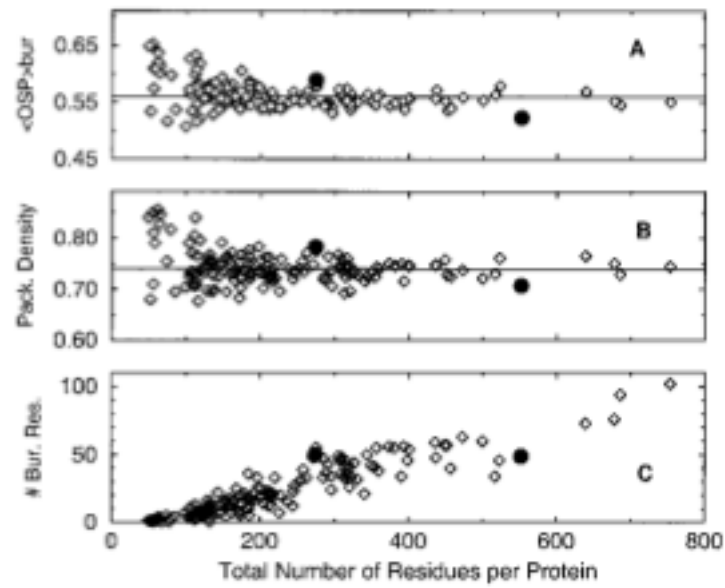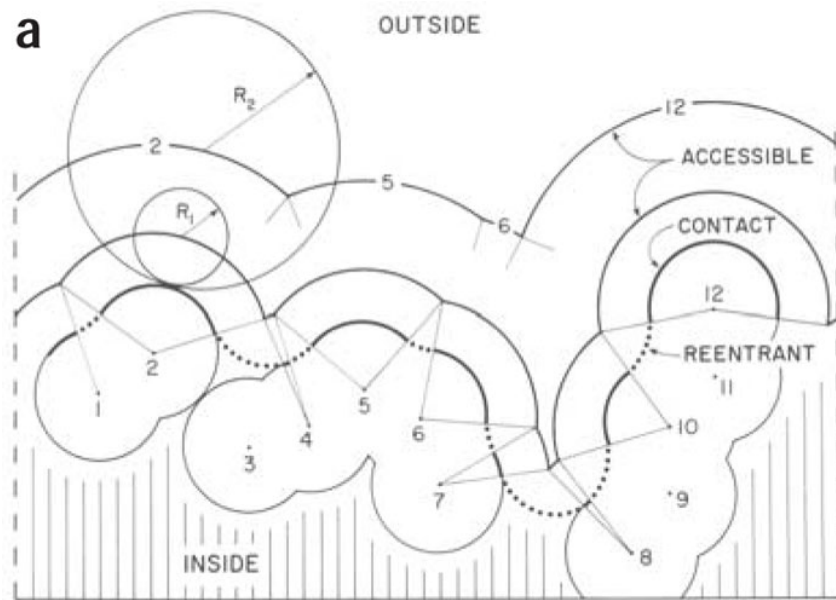
Figure 2: Ramachandran plots of allowed $\phi/\psi$ combinations for 3 values of $\tau$ [2]. The solid red lines enclose the 'normally allowed' $\phi/\psi$ combinations and the dashed blue line indicates the 'outer limit'.
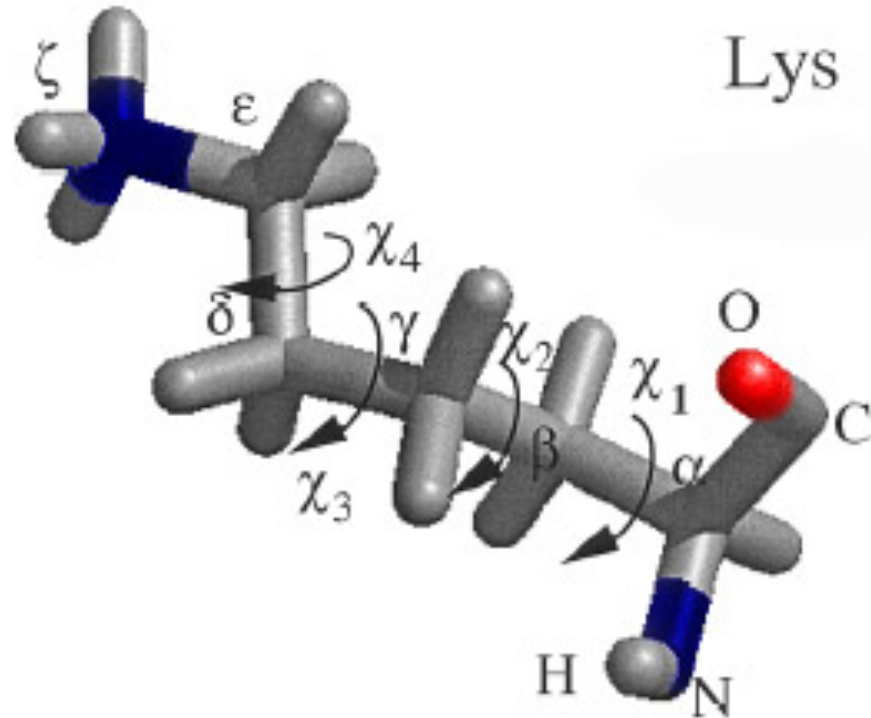
16

Prof. Fred Richards, Yale



a

OUTSIDE

ACCESSIBLE

CONTACT

REENTRANT

INSIDE

b

# Side-Chain Dihedral Angles

$\chi_4$: Lys, Arg

$\chi_5$: Arg

Side chain: $C_\alpha$-$CH_2$-$CH_2$-$CH_2$-$CH_2$-$NH_3$

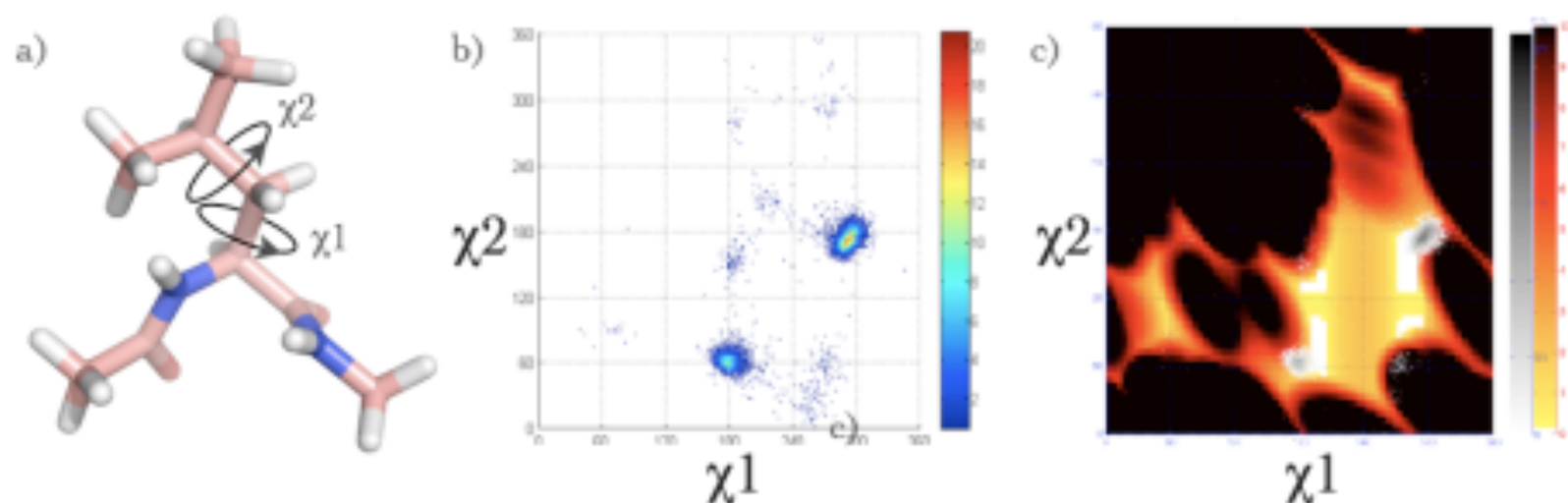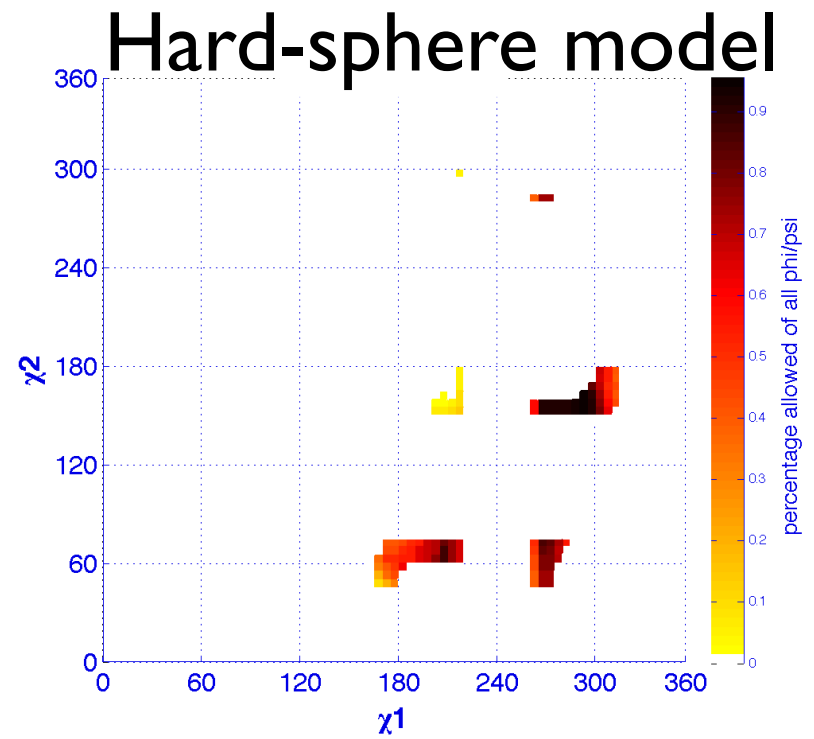Use $NC_\alpha C_\beta C_\gamma C_\delta C_\varepsilon N_\zeta$ to define $\chi_1, \chi_2, \chi_3, \chi_4$
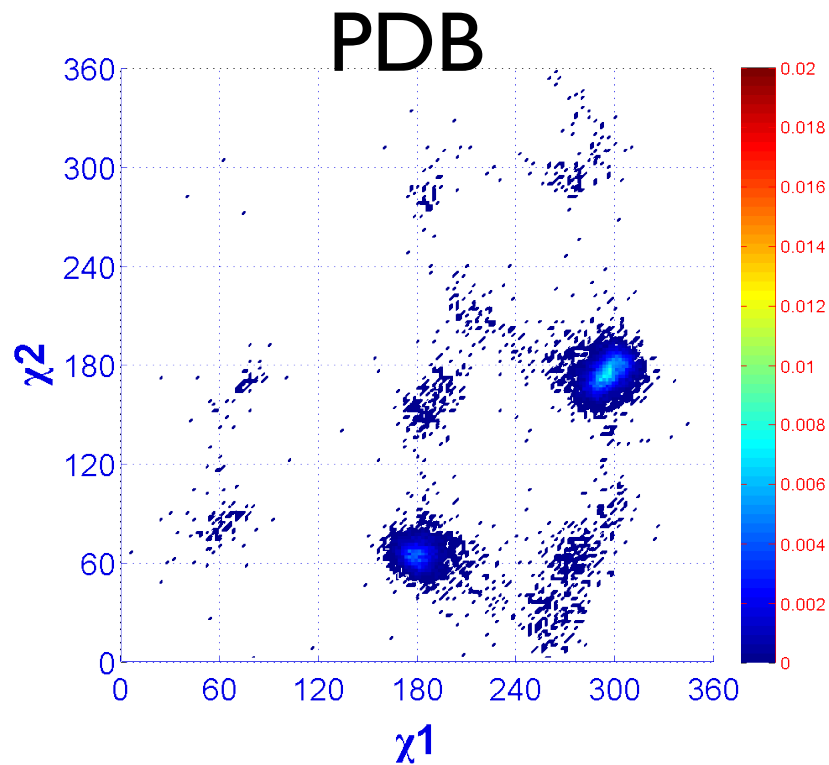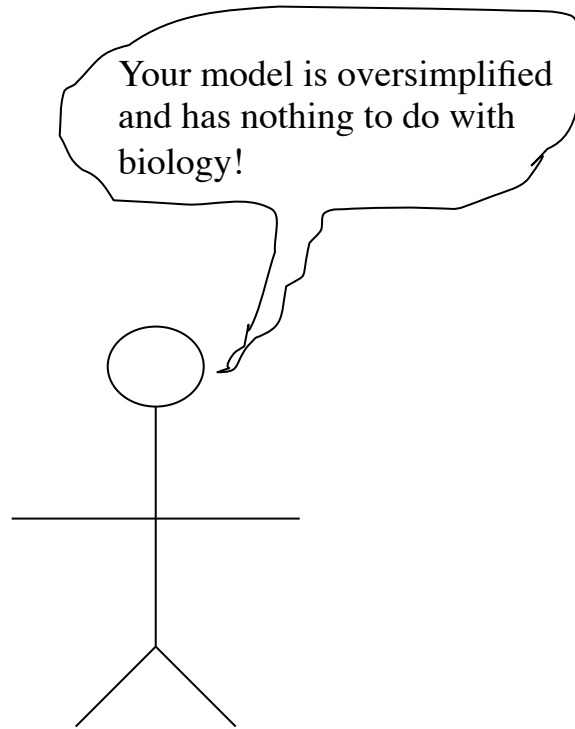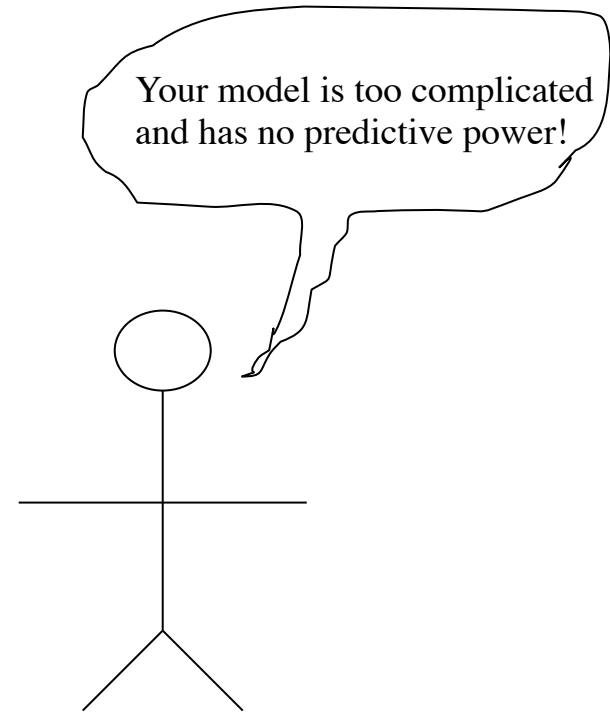
18

Figure 2: (a) Stick representation of a Leu dipeptide showing the side-chain dihedrals chi1 and chi2 (carbon=pink, nitrogen=blue, oxygen=red, hydrogen=white). (b) Density plot of chi1/chi2 value for every Leu in the Dunbrack database [5]. (c) My calculated energy landscape for the Leu dipeptide using the repulsive Lennard-Jones interaction potential overlaid on the Dunbrack probability distribution (grey scale). White regions correspond to low-energy minima with energy increasing from yellow to black.

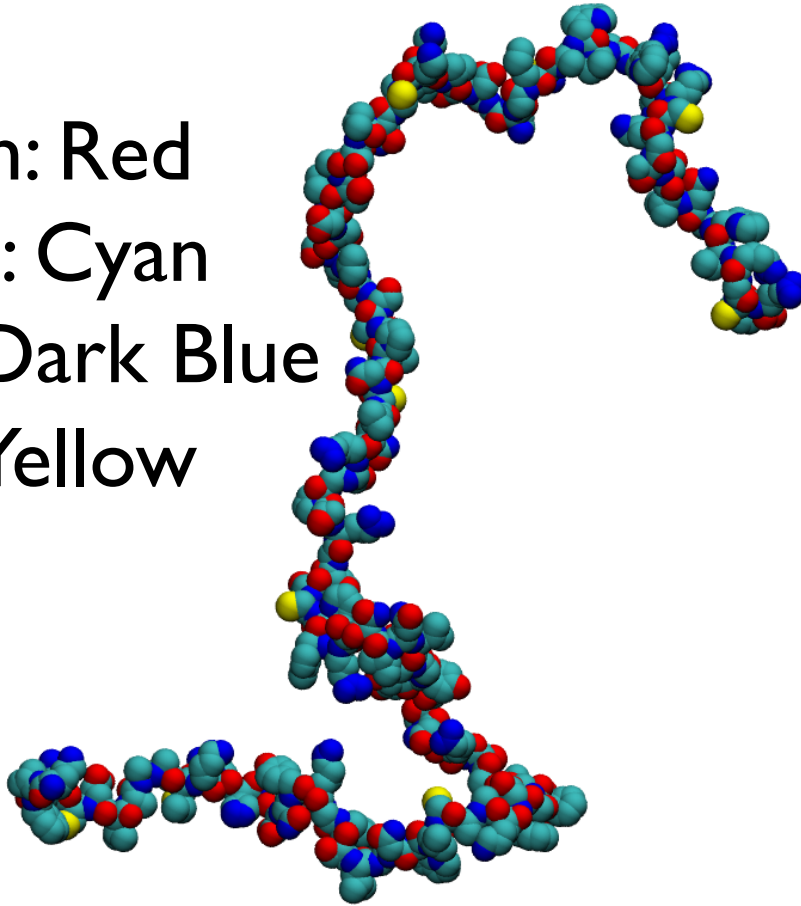# Sidechain Dihedral Angle Distributions for Leu

## PDB



## Hard-sphere model
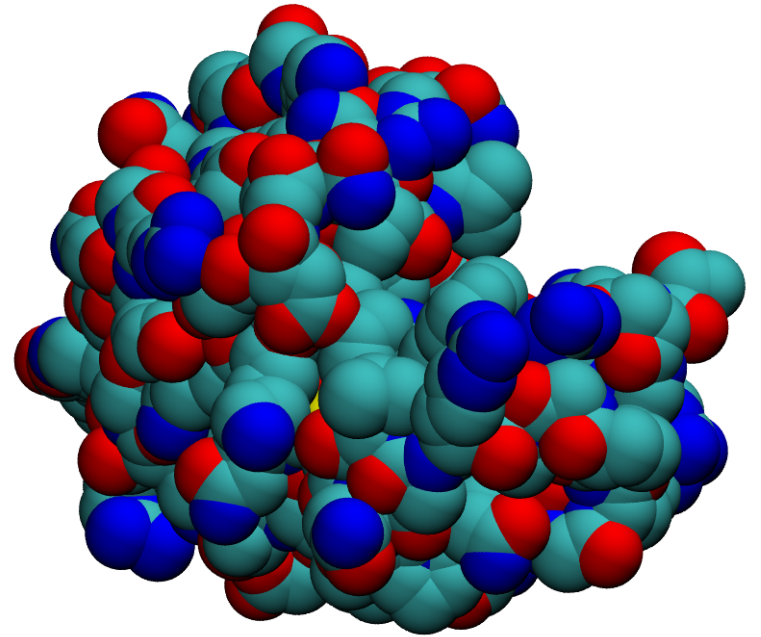
Molecular biologist          Biological Physicist

# Folding Transition

Oxygen: Red
Carbon: Cyan
Nitrogen: Dark Blue
Sulfur: Yellow

$T > T_m$

$T < T_m$

# Possible Strategies for Understanding Protein Folding

- For all possible conformations, compute free energy from atomic interactions within protein and protein-solvent interactions; find conformation with lowest free energy…e.g using all-atom molecular dynamics simulations

<span style="color:red">Not possible?, limited time resolution</span>

- Use coarse-grained models with effective interactions between residues and residues and solvent

<span style="color:red">General, but qualitative</span>

# Why do proteins fold (correctly & rapidly)??

Levinthal's paradox:

For a protein with N amino acids, number of backbone conformations/minima

$$N_c \sim \mu^{2N}$$

$\mu$ = # allowed dihedral angles

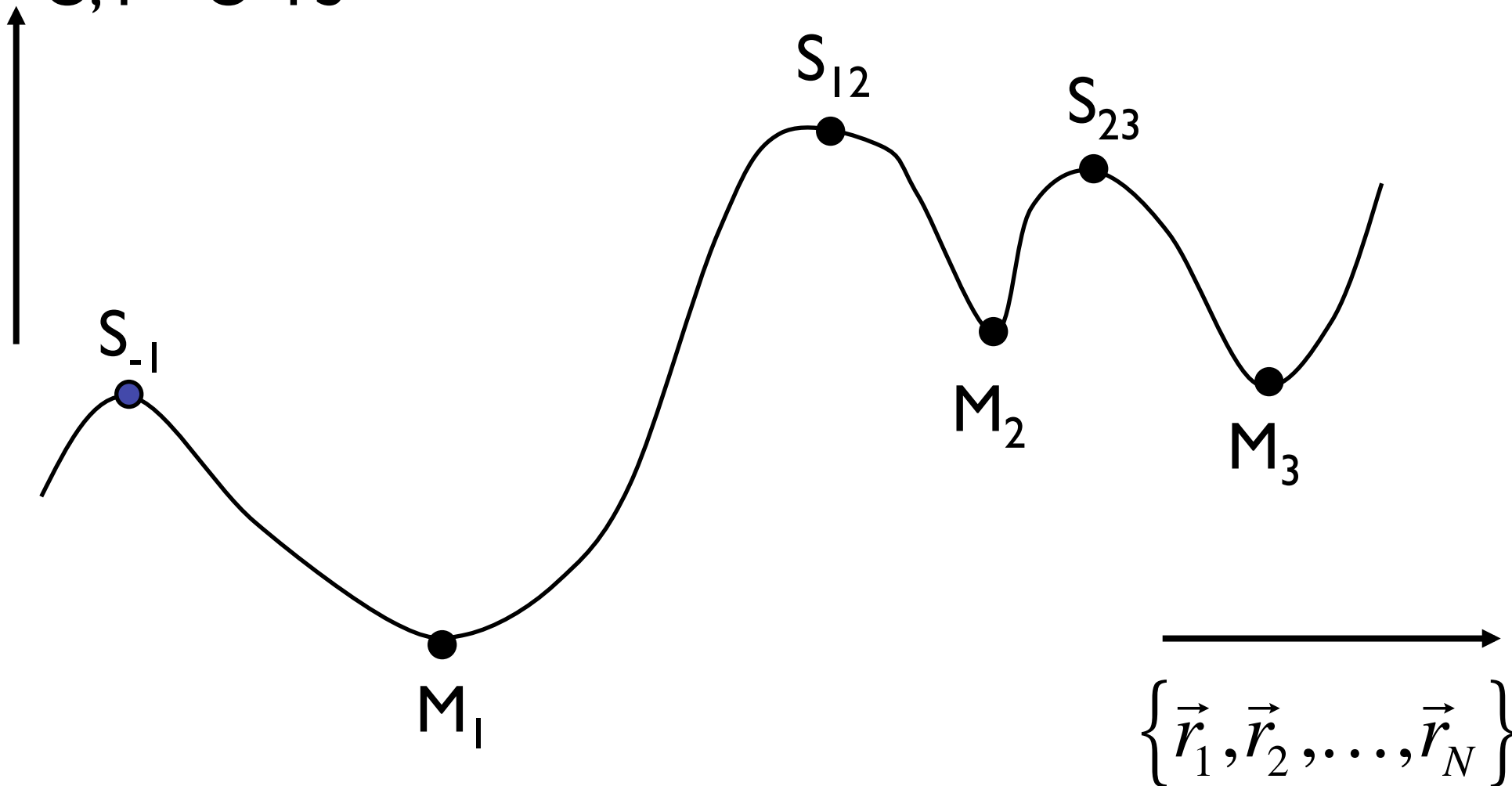How does a protein find the global optimum w/o global search?  Proteins fold much faster.

$$N_c \sim 3^{200} \sim 10^{95}$$

$$\tau_{fold} \sim N_c \, \tau_{sample} \sim 10^{83} \text{ s} \quad \text{vs} \quad \tau_{fold} \sim 10^{-6}\text{-}10^{-3} \text{ s}$$

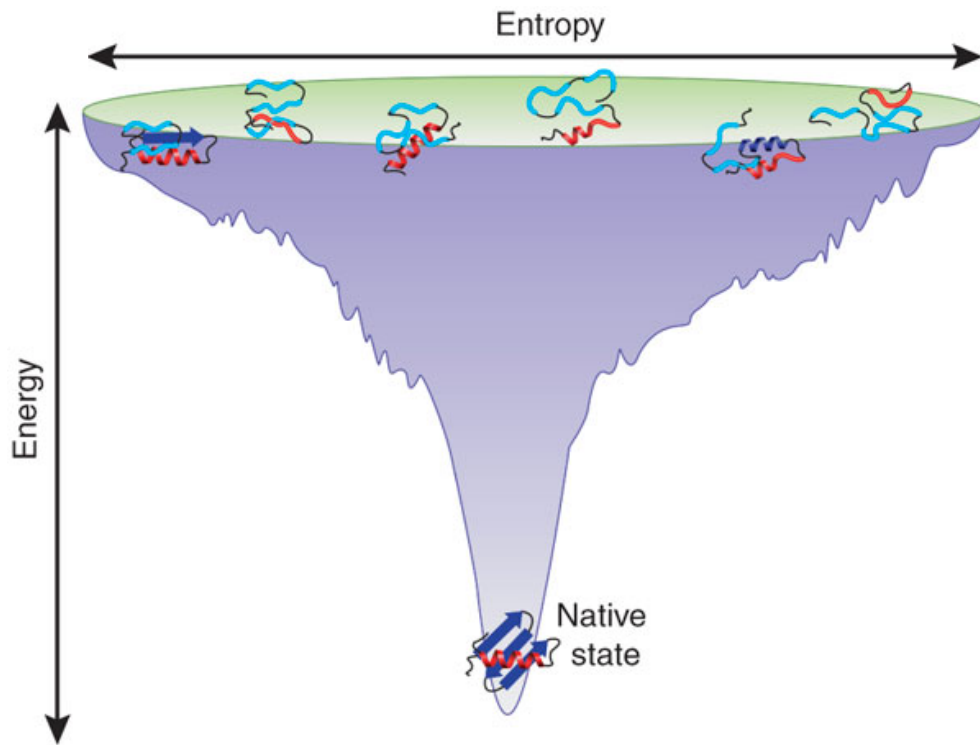$$\tau_{universe} \sim 10^{17} \text{ s}$$

24

# Energy Landscape

$U, F = U - TS$

$S_{12}$

$S_{23}$

$S_{-1}$

$M_2$

$M_3$

$M_1$

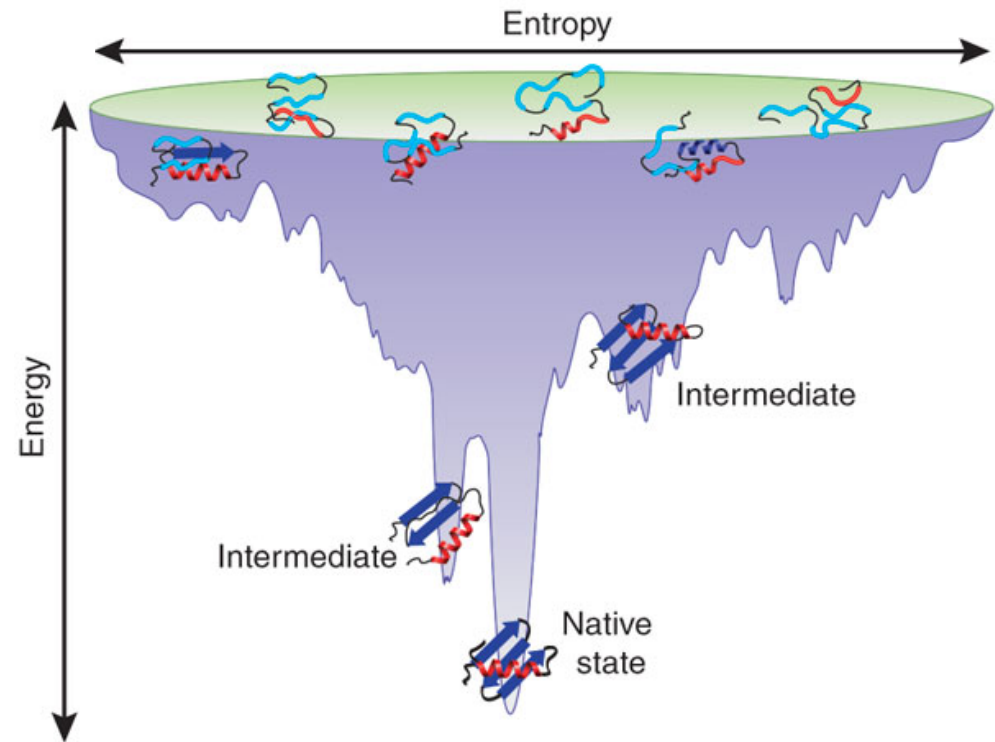$$\{\vec{r}_1, \vec{r}_2, \ldots, \vec{r}_N\}$$

all atomic coordinates; dihedral angles

$$\vec{\nabla} U = 0 \quad \begin{cases} \nabla^2 U > 0 & \text{minimum} \\ \nabla^2 U = 0 & \text{saddle point} \\ \nabla^2 U < 0 & \text{maximum} \end{cases}$$

25

# Roughness of Energy Landscape



smooth, funneled

(Wolynes et. al. 1997)

rough

26

# Folding Pathways



Collapsed structures

dead end

Native fold

Few paths

Many paths

similarity to native state

$n_g$    $n$

27
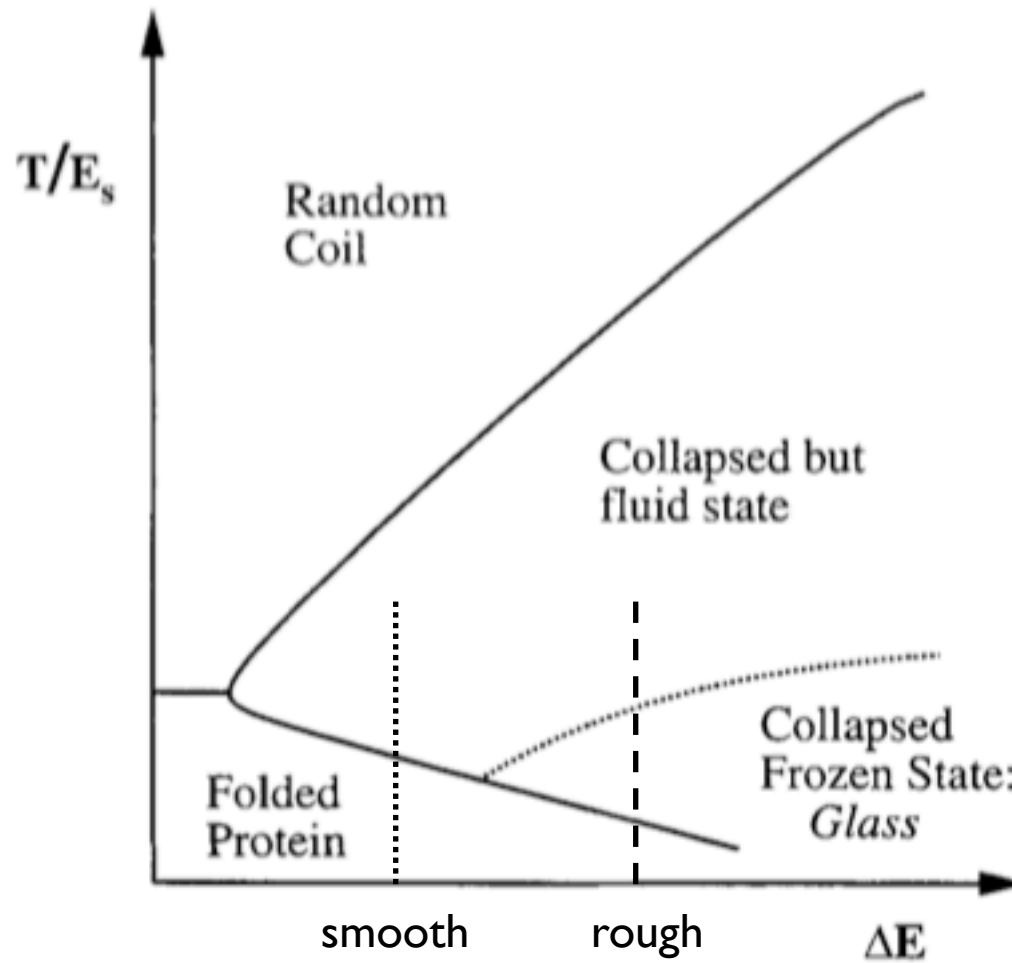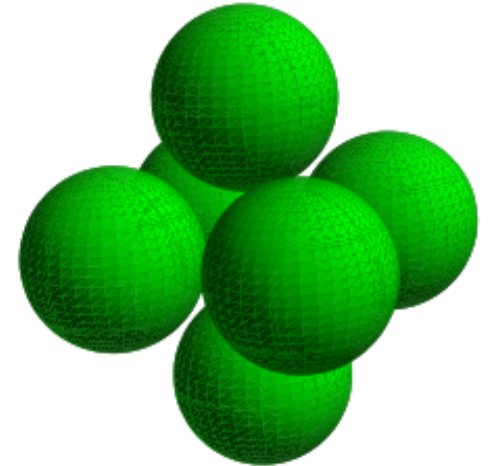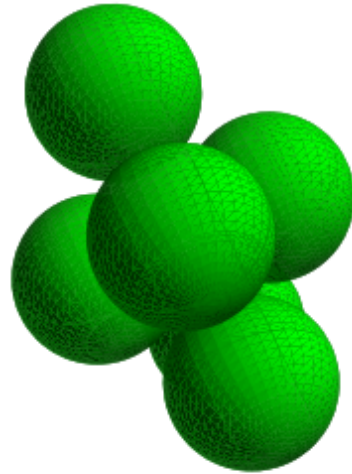
# Folding Phase Diagram



28

# Open Questions

• What differentiates the native state from other low-lying energy minima?

• How many low-lying energy minima are there?  Can we calculate landscape roughness from sequence?

• What determines whether protein will fold to the native state or become trapped in another minimum?

• What are the pathways in the energy landscape that a given protein follows to its native state?
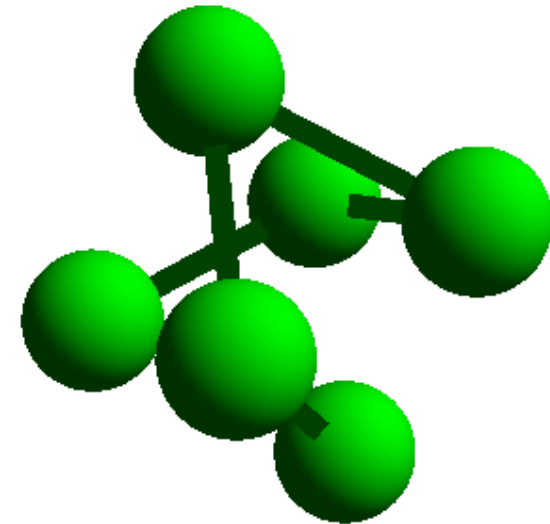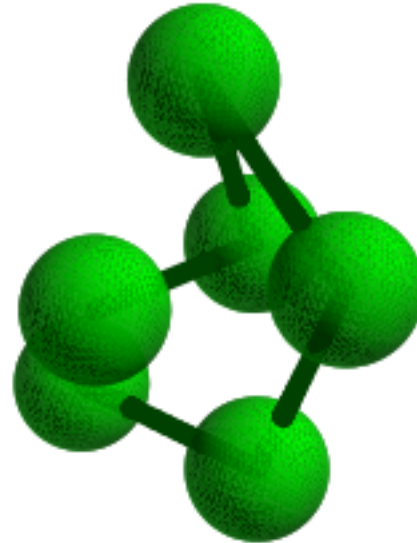
## NP Hard Problem!

29

# Digression---Number of Energy Minima for Sticky Spheres

| $N_m$ | $N_s$ | $N_p$ |
|-------|-------|-------|
| 4 | 1 | 1 |
| 5 | 1 | 6 |
| 6 | 2 | 50 |
| 7 | 5 | 486 |
| 8 | 13 | 5500 |
| 9 | 52 | 49029 |
| 10 | - | - |

$N_s \sim \exp(aN_m)$;
$N_p \sim \exp(bN_m)$ with $b > a$

sphere packings

polymer packings

30