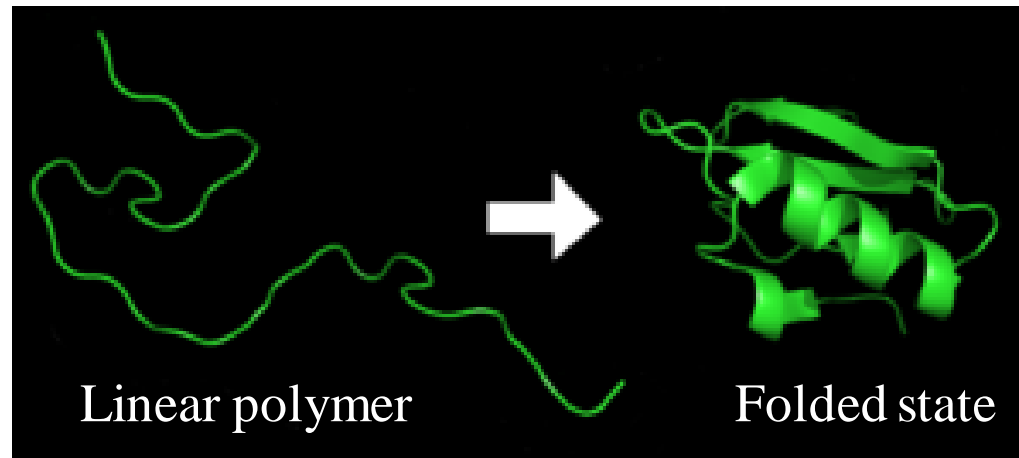


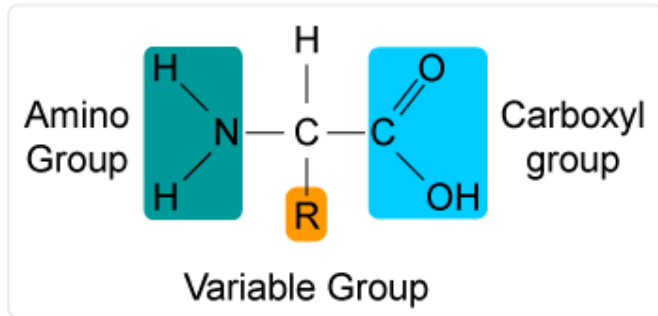
# What are proteins?



- 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

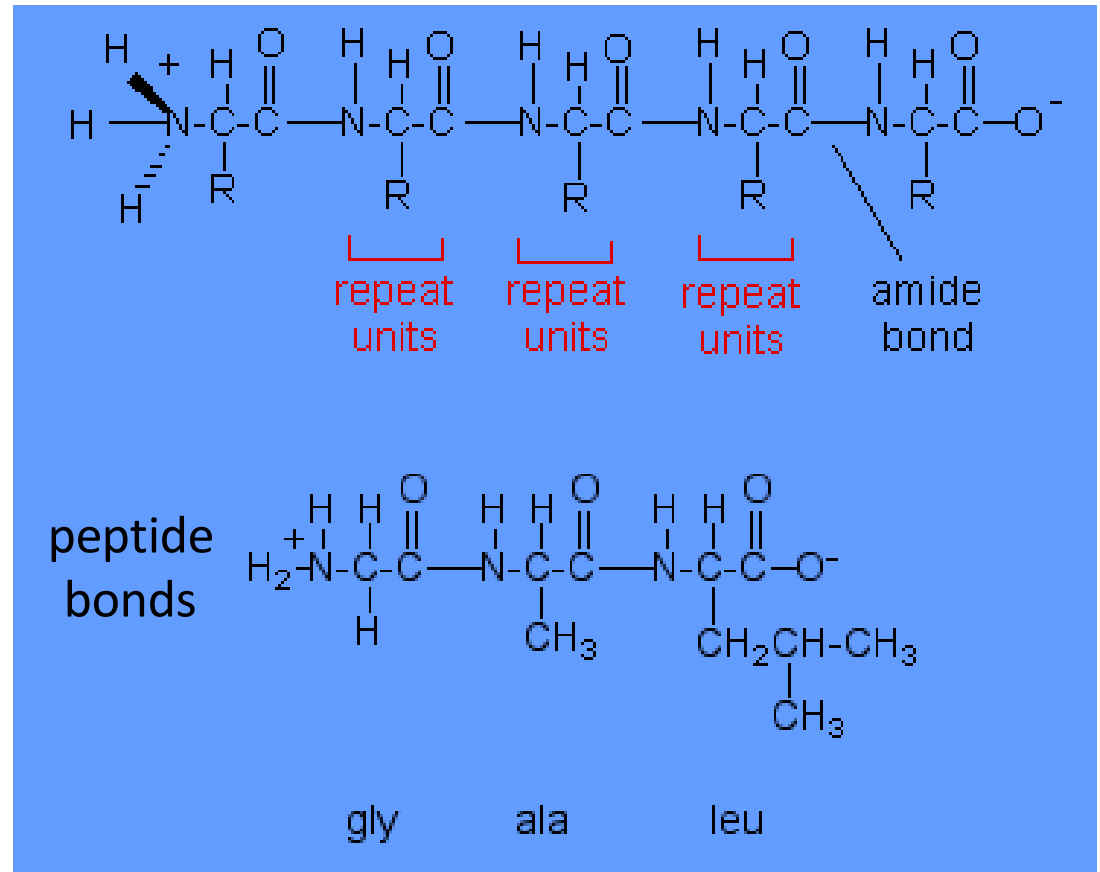
# Amino Acids I

General structure of Amino Acids



N-terminal       $C_{\alpha}$       C-terminal

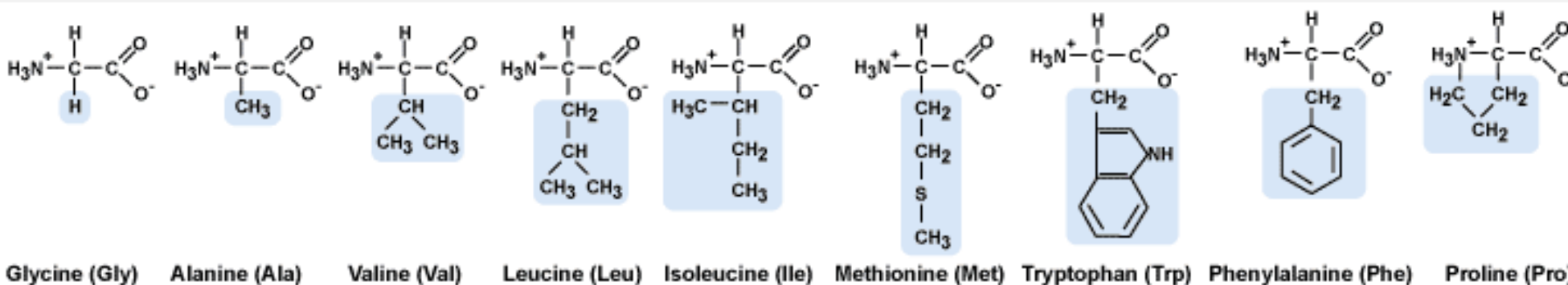
R  
variable  
side chain



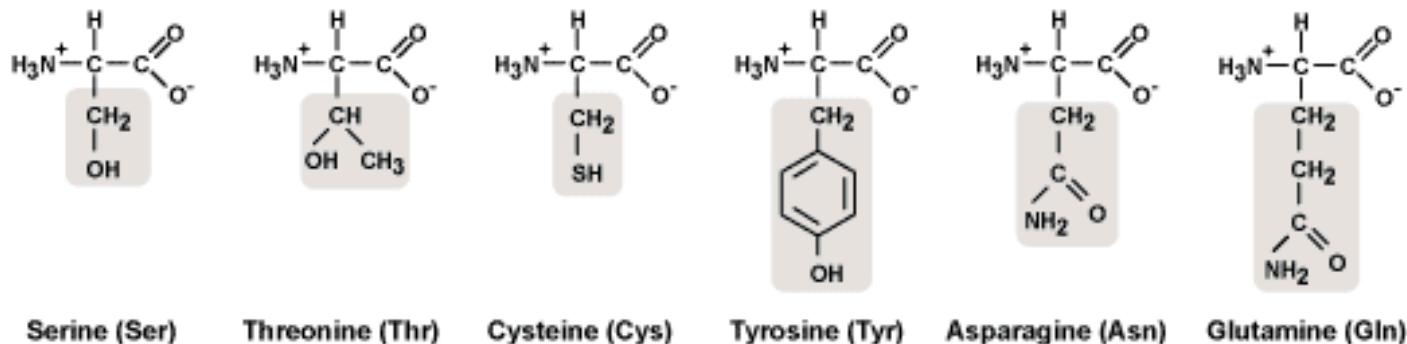
- Side chains differentiate amino acid repeat units
- Peptide bonds link residues into polypeptides

# Amino Acids II

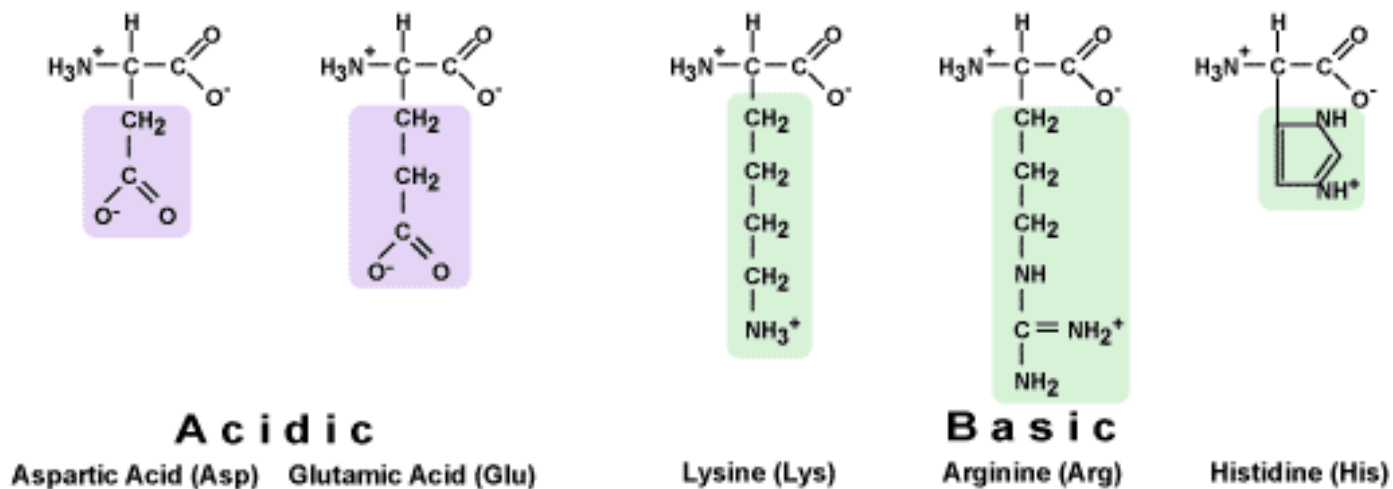
NONPOLAR



POLAR



Electrically Charged



(-)

3

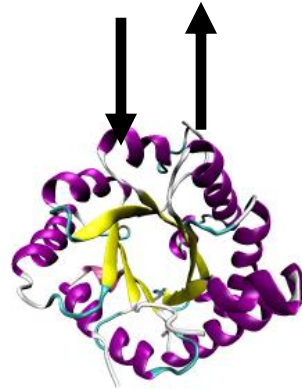
(+)

# 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-...



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

Levinthal's paradox:

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

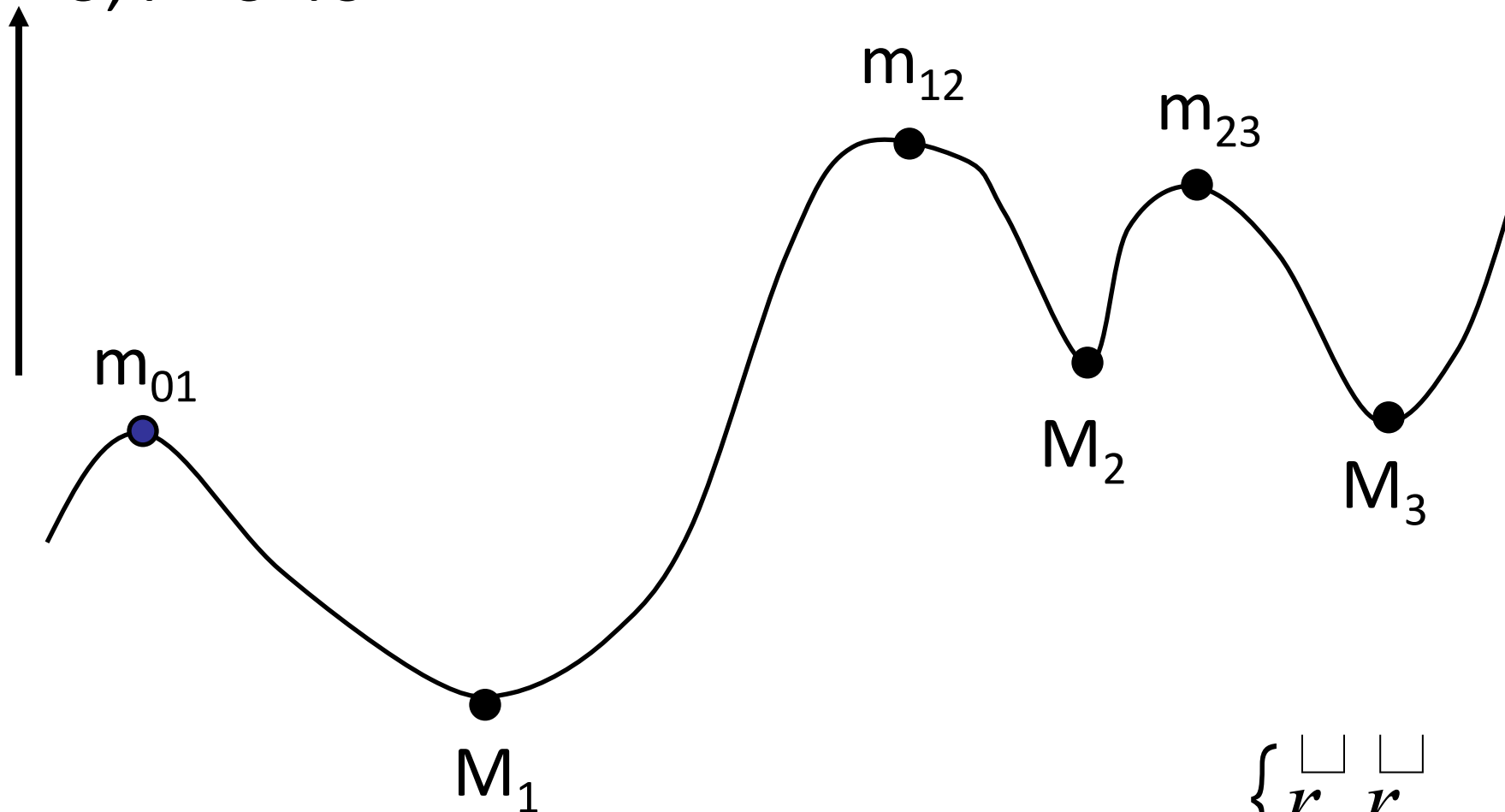
$$N_c \sim \mu^{2N} \quad \mu = \# \text{ allowed dihedral angles}$$

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

$$\begin{aligned} N_c &\sim 3^{200} \sim 10^{95} \\ \tau_{\text{fold}} &\sim N_c \tau_{\text{sample}} \sim 10^{83} \text{ s} \quad \text{vs} \quad \tau_{\text{fold}} \sim 10^{-6} - 10^{-3} \text{ s} \\ \tau_{\text{universe}} &\sim 10^{17} \text{ s} \end{aligned}$$

# Energy Landscape

$U, F = U - TS$



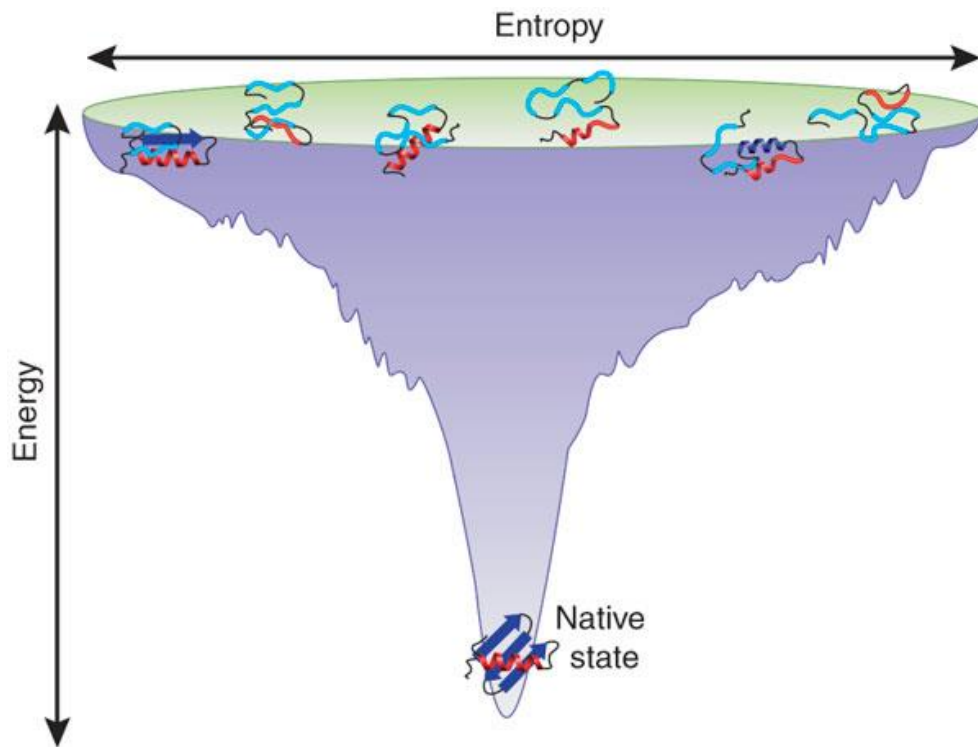
$$\vec{\nabla} U = \mathbf{0}$$

- $\nabla^2 U > 0$  Minimum (M)
- $\nabla^2 U = 0$  saddle point
- $\nabla^2 U < 0$  Maximum (m) 6

$$\left\{ \begin{array}{c} \square \quad \square \quad \square \\ r_1, r_2, \dots, r_N \end{array} \right\}$$

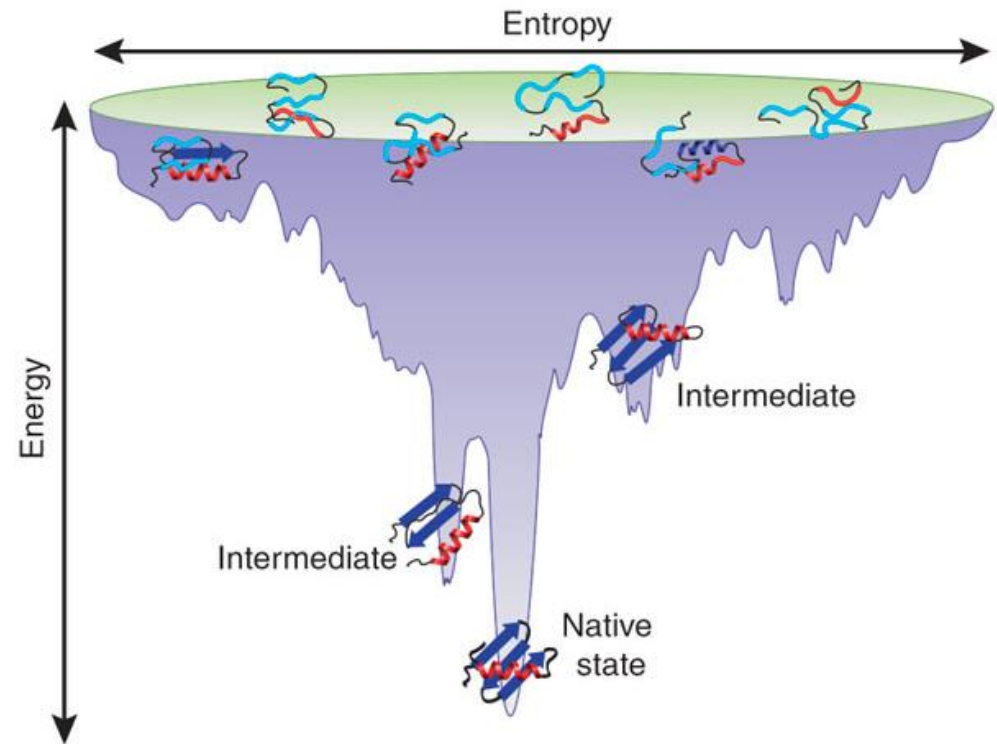
all atomic  
coordinates;  
dihedral angles

# Roughness of Energy Landscape



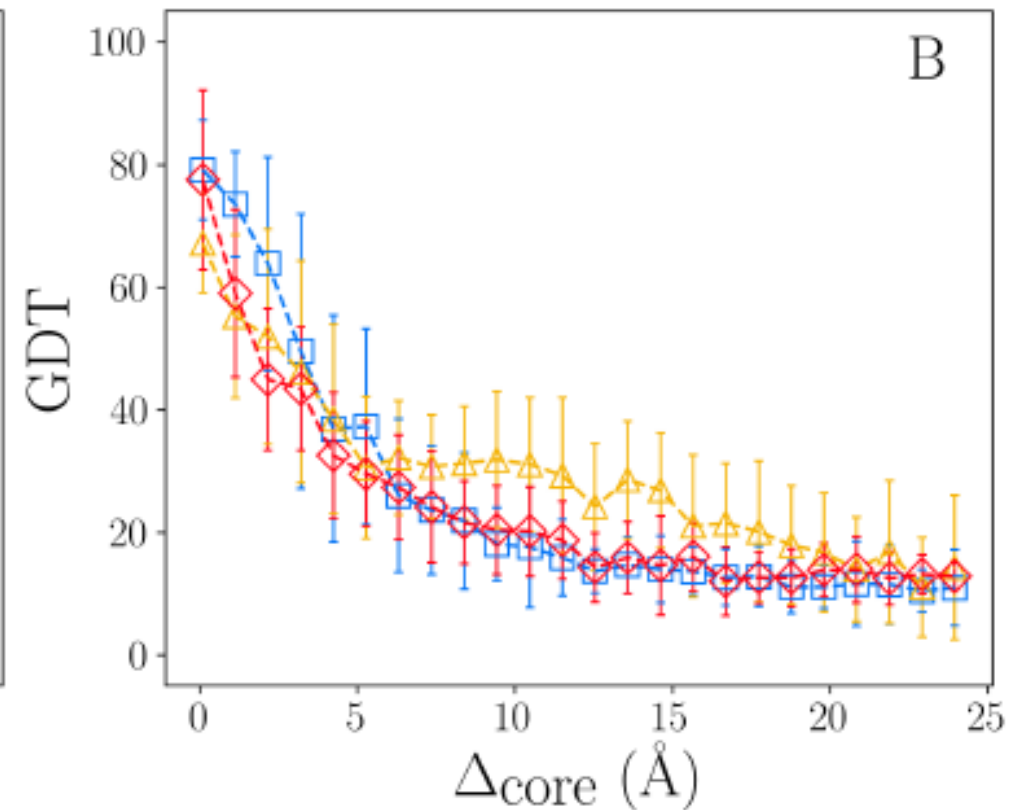
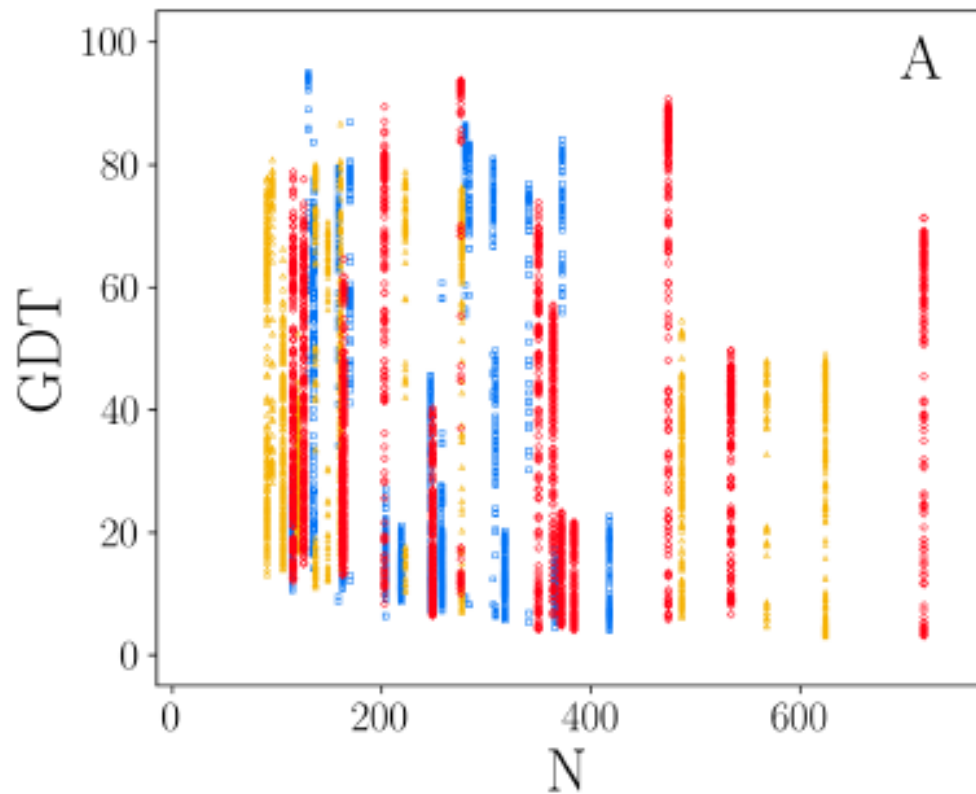
smooth, funneled

(Wolynes et. al. 1997)

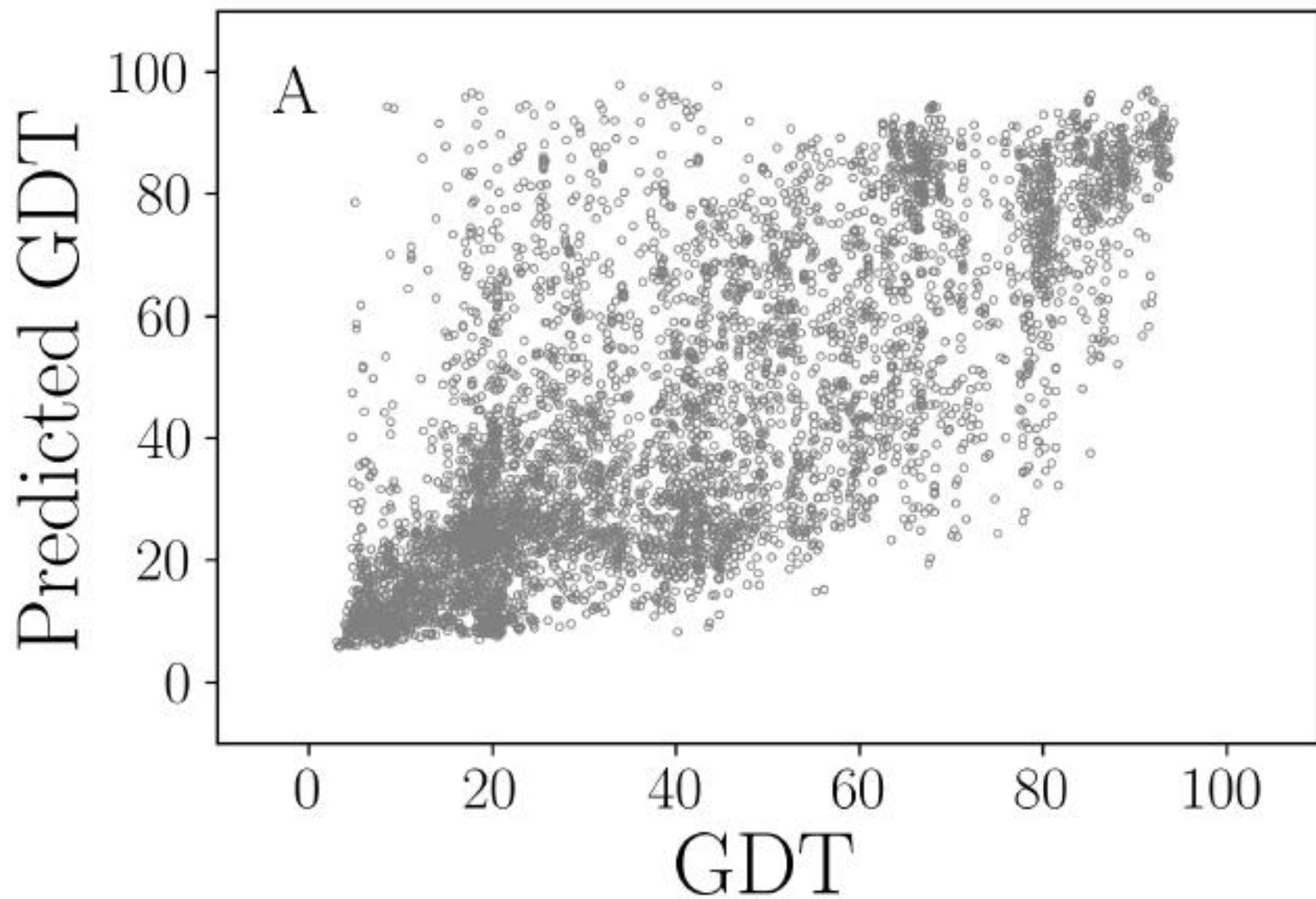


rough

# Critical Assessment of Structure Prediction (CASP)







# Driving Forces

- Folding: hydrophobicity, hydrogen bonding, van der Waals interactions, ...
- Unfolding: increase in conformational entropy, electric charge...

inside

H (hydrophobic)

outside

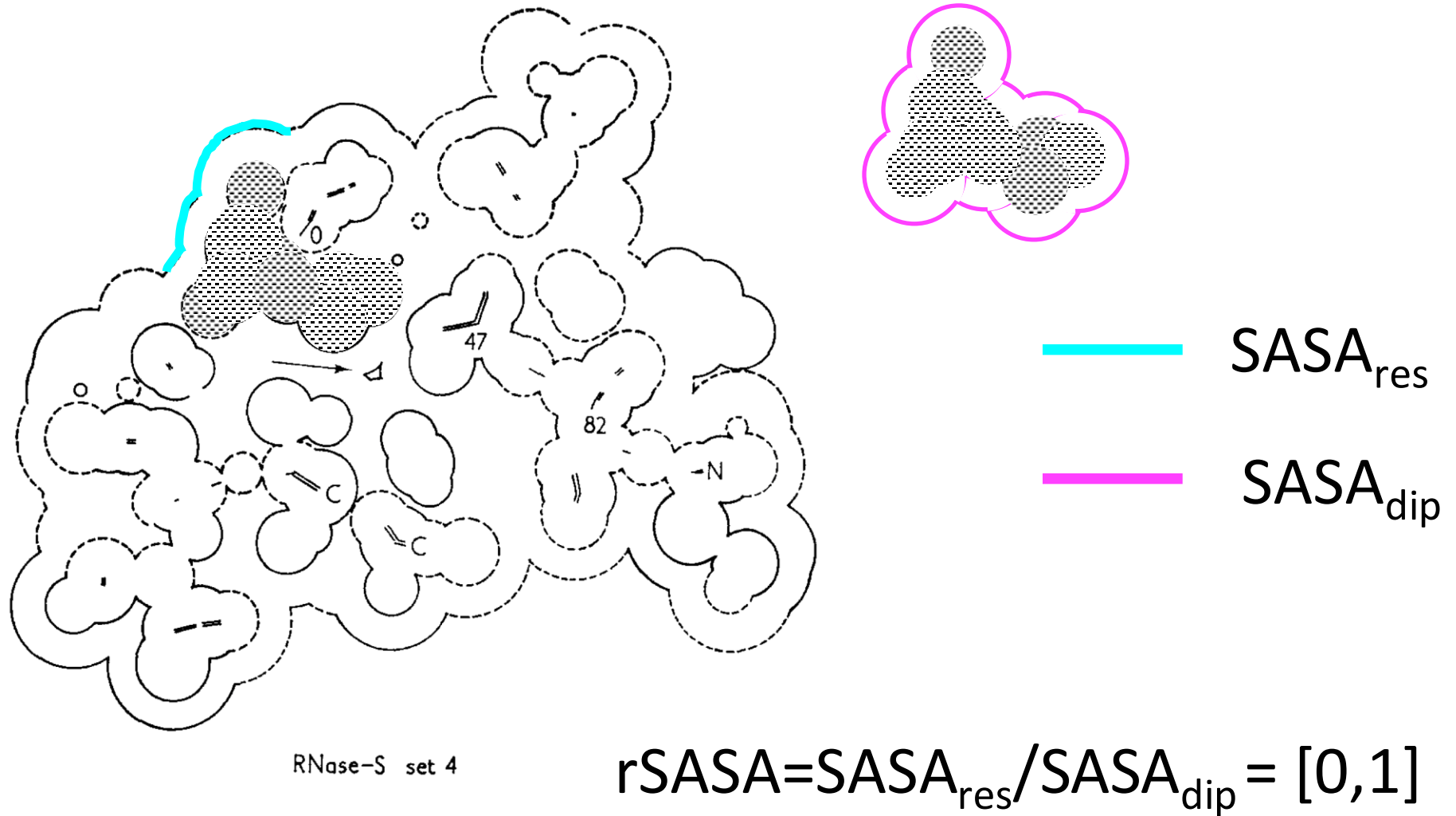
P (polar)

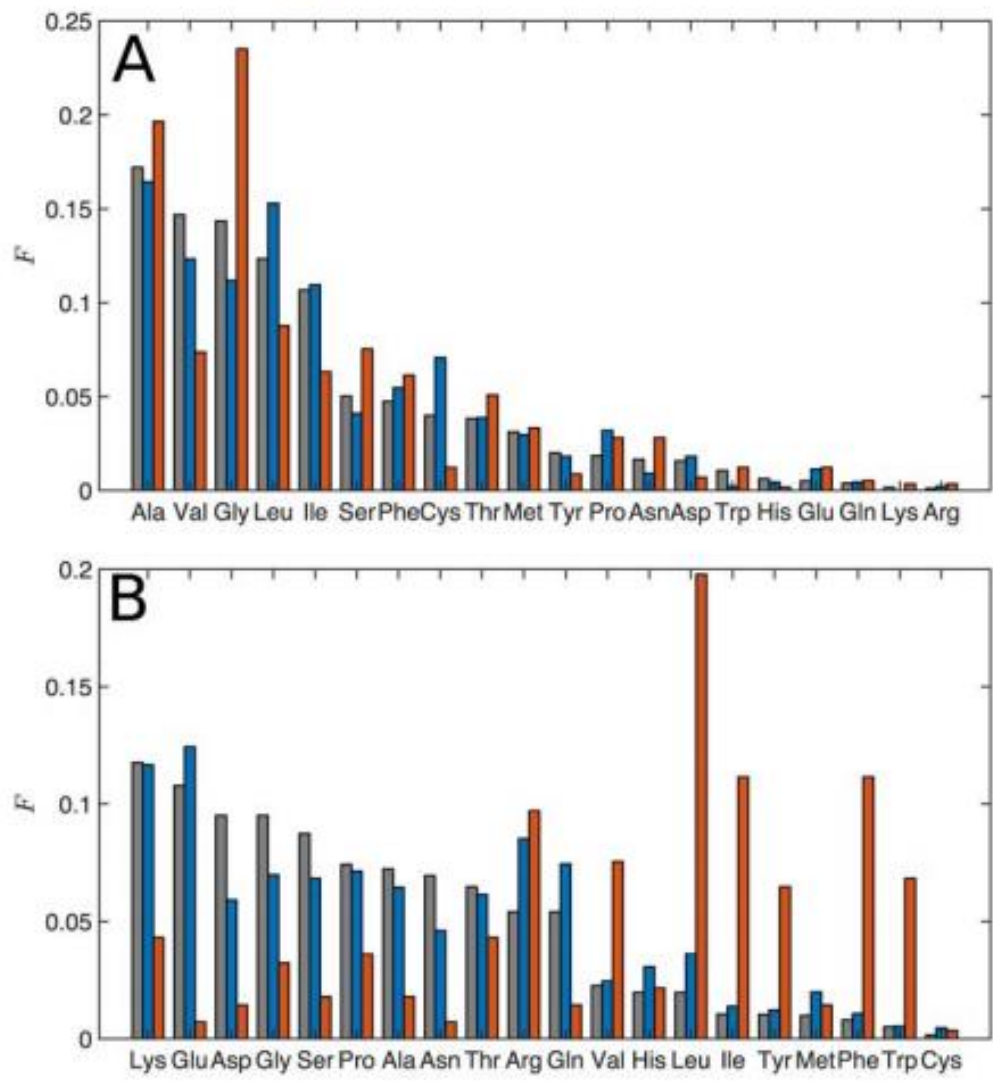
Hydrophobicity index

At pH 2 <sup>a</sup>		At pH 7 <sup>b</sup>	
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

<sup>a</sup> pH 2 values: Normalized from Sereda et al., J. Chrom. 676: 139-153 (1994).  
<sup>b</sup> pH 7 values: Monera et al., J. Pept. Sci. 1: 319-329 (1995).

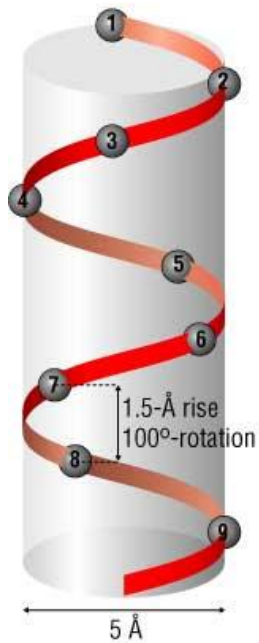
# Solvent Accessible Surface Area and rSASA



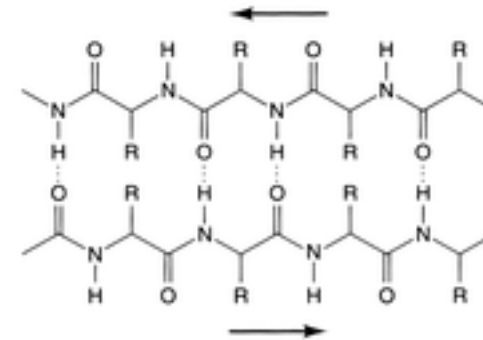
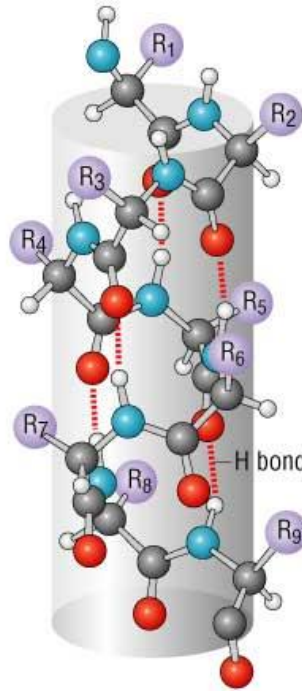
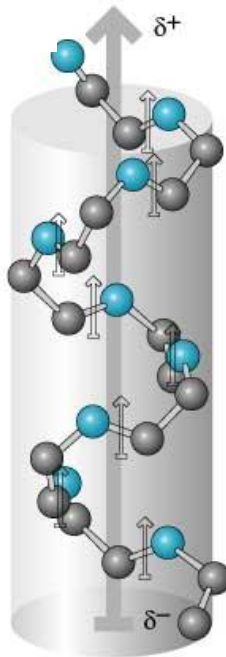


**FIGURE 5** Fractions of amino acids with A,  $rSASA \leq 10^{-3}$  and B,  $rSASA > 0.5$  for residues in the Dun1.0 (grey), PPI (blue), and TM (red) datasets. The fractions are defined relative to the total number of residues in each rSASA category. C, The fractions of core residues (light bars) and non-core residues ( $rSASA > 0.5$ , dark bars) among the 11 non-charged residues (Ala, Gly, Ile, Leu, Met, Phe, Ser, Thr, Trp, Tyr, and Val) [Color figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

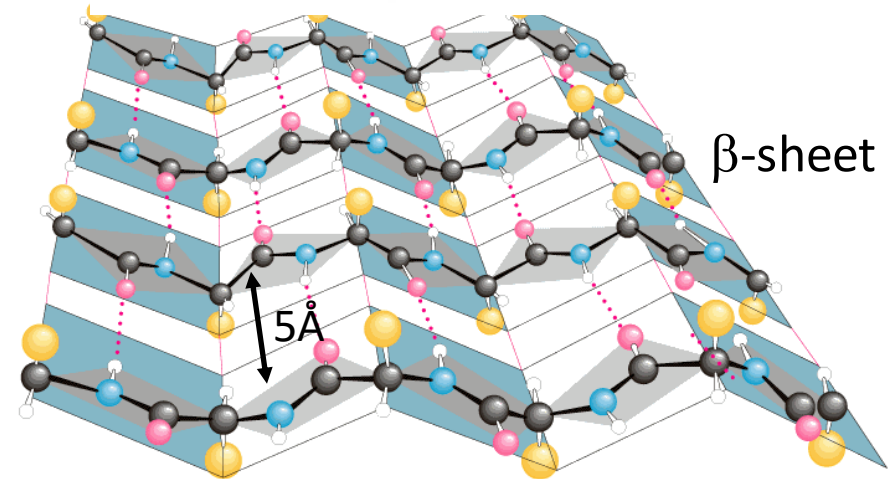
# Secondary Structure: Loops, $\alpha$ -helices, $\beta$ -strands/sheets



$\alpha$ -helix



$\beta$ -strand

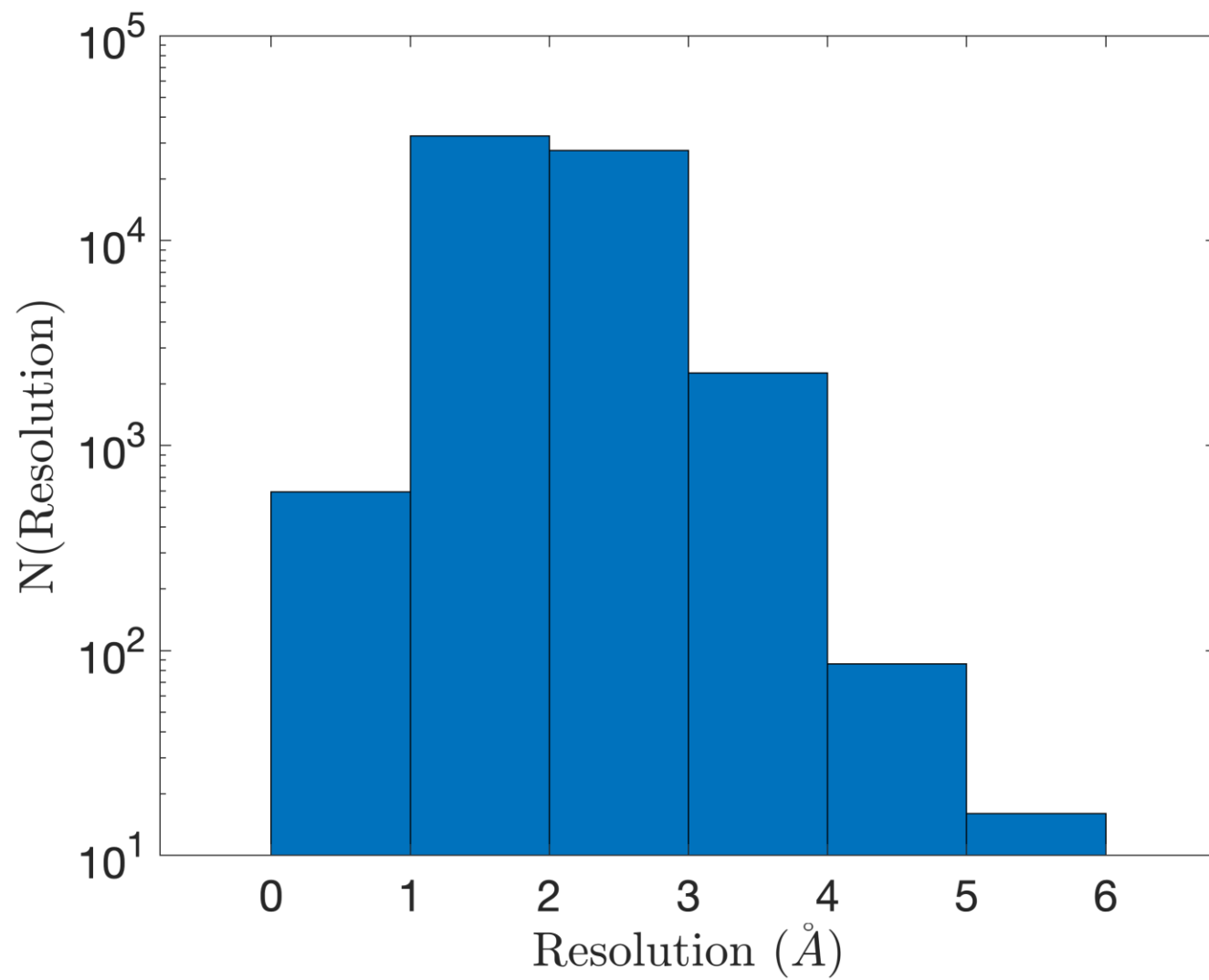


$\beta$ -sheet

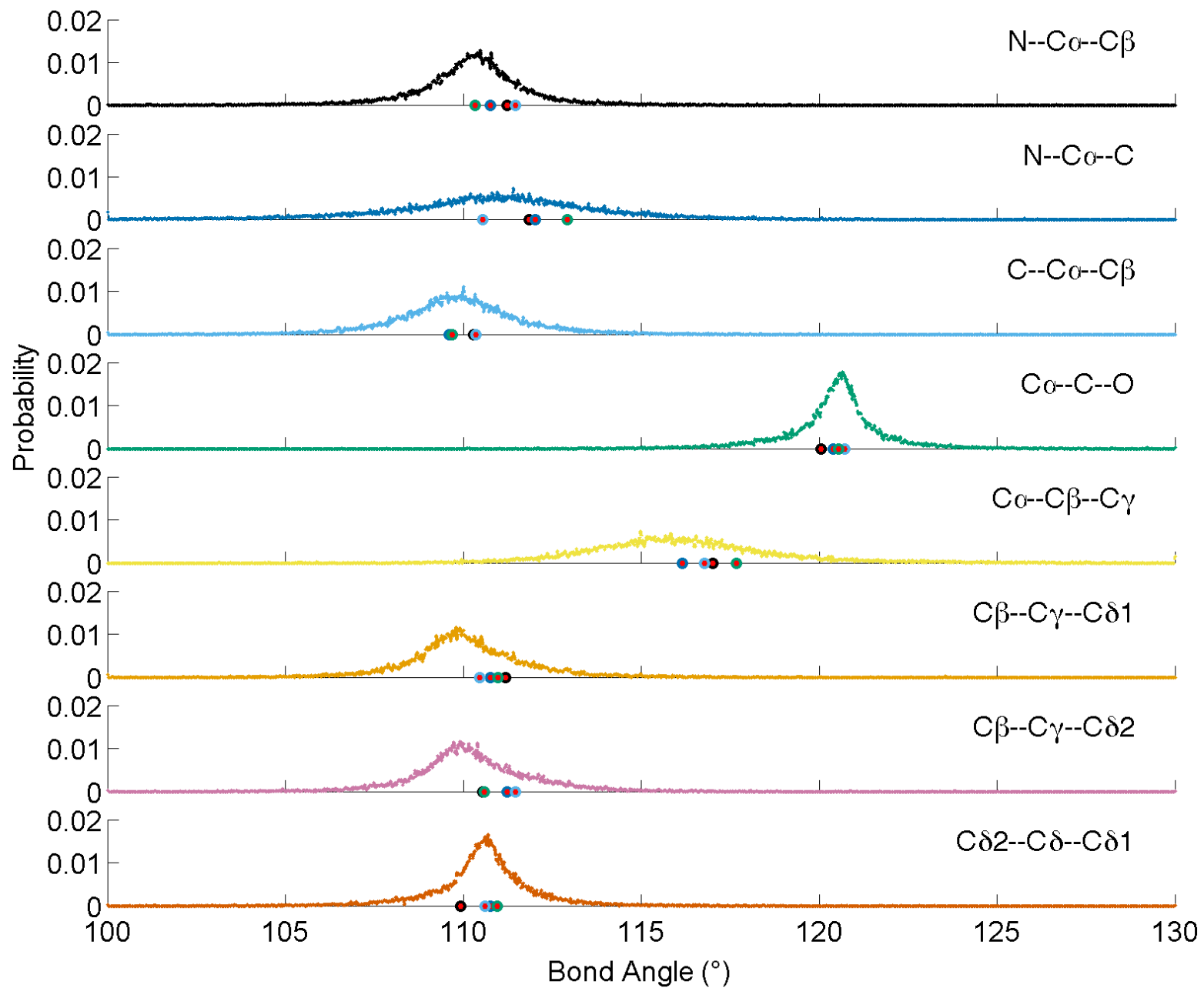
- Right-handed; three turns
- Vertical hydrogen bonds between  $\text{NH}_2$  (teal/white) backbone group and  $\text{C}=\text{O}$  (grey/red) backbone group four residues earlier in sequence
- Side chains (R) on outside; point upwards toward  $\text{NH}_2$
- Each amino acid corresponds to  $100^\circ$ ,  $1.5\text{\AA}$ , 3.6 amino acids per turn
- $(\phi, \psi) = (-60^\circ, -45^\circ)$
- $\alpha$ -helix propensities: Met, Ala, Leu, Glu

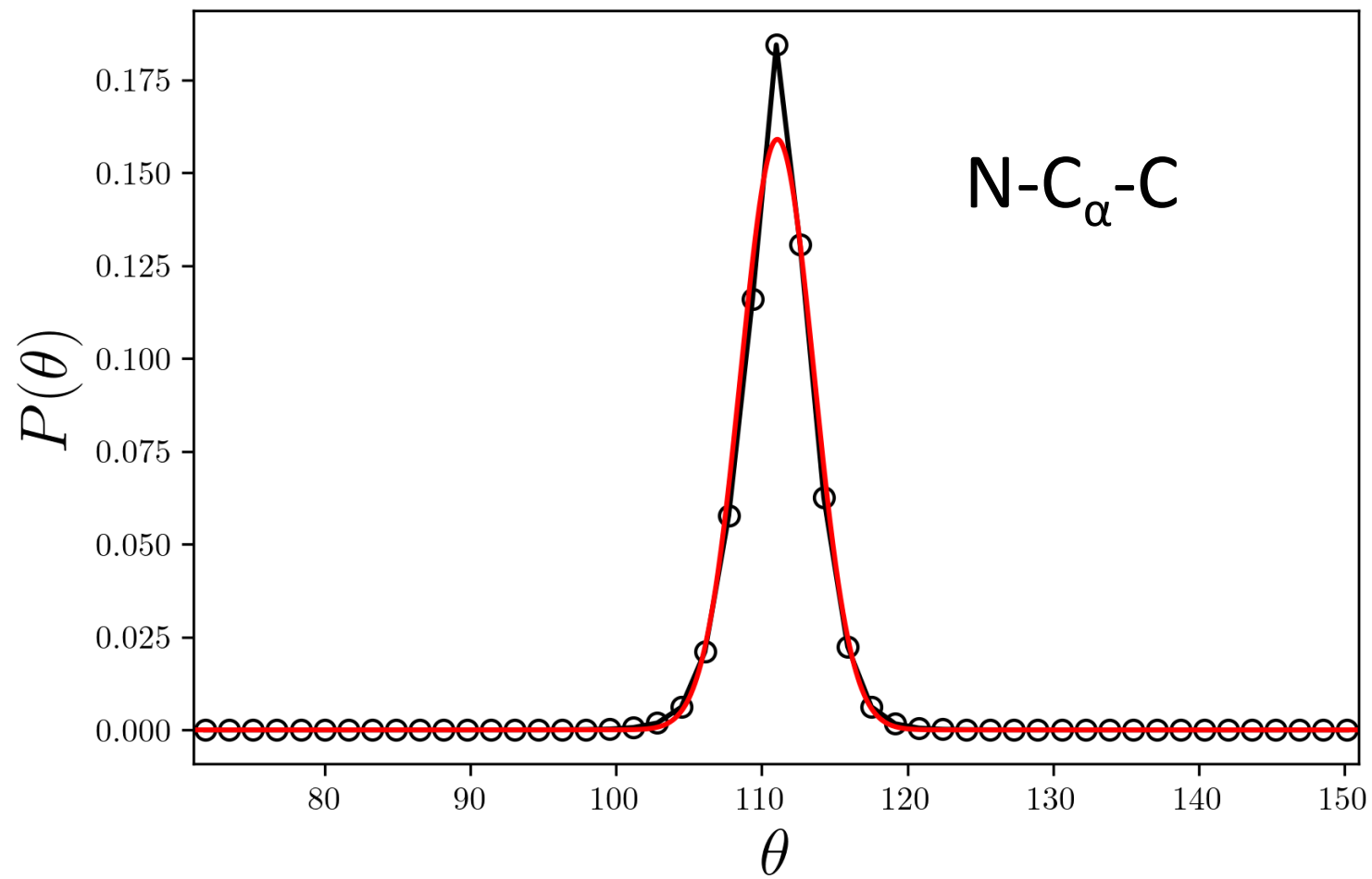
- 5-10 residues; peptide backbones fully extended
- $\text{NH}$  (blue/white) of one strand hydrogen-bonded to  $\text{C}=\text{O}$  (black/red) of another strand
- $\text{C}_\alpha$ , side chains (yellow) on adjacent strands aligned; side chains along single strand alternate up and down
- $(\phi, \psi)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

$N_s=62,938$  monomeric xtal structures



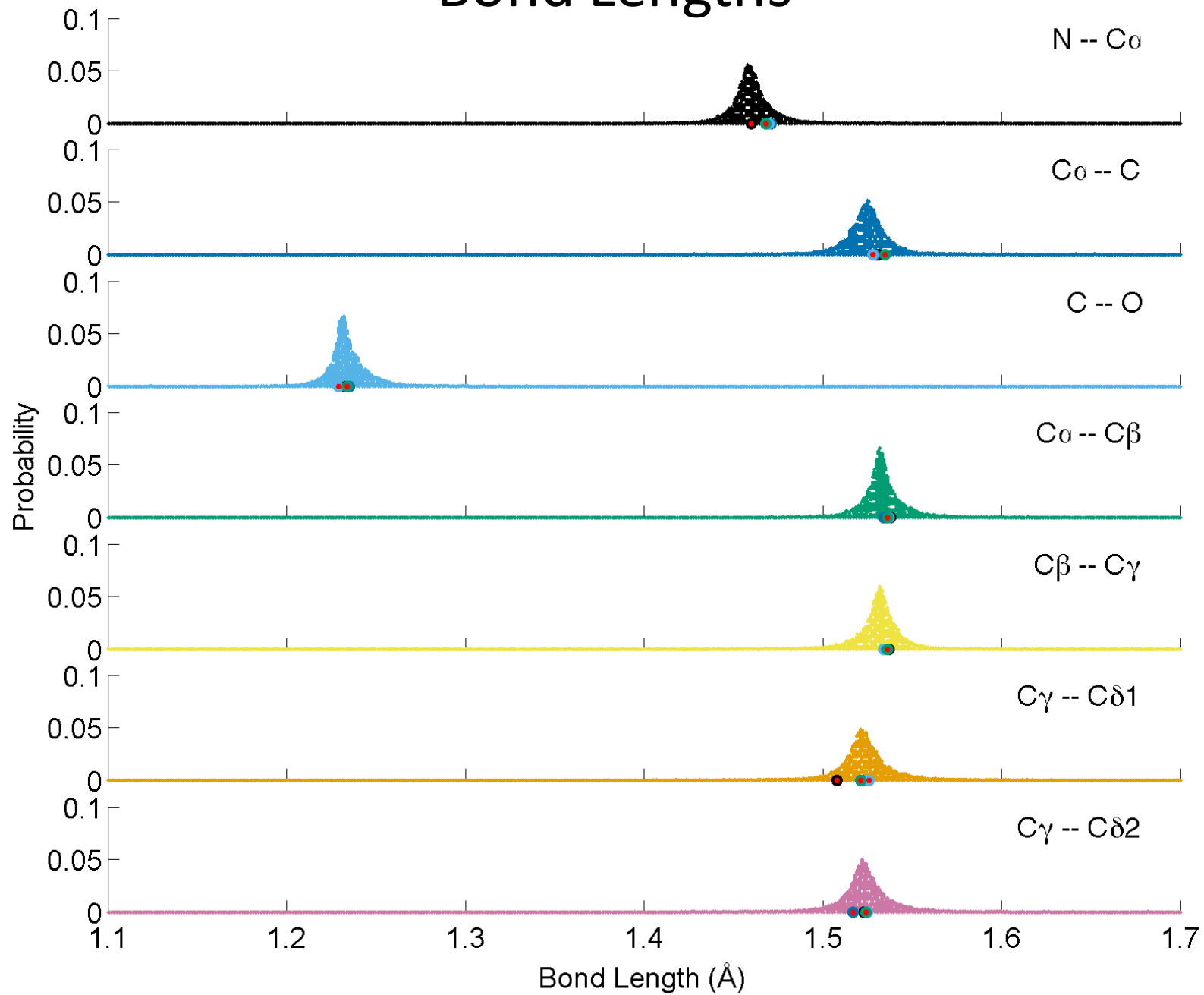
# Bond Angles

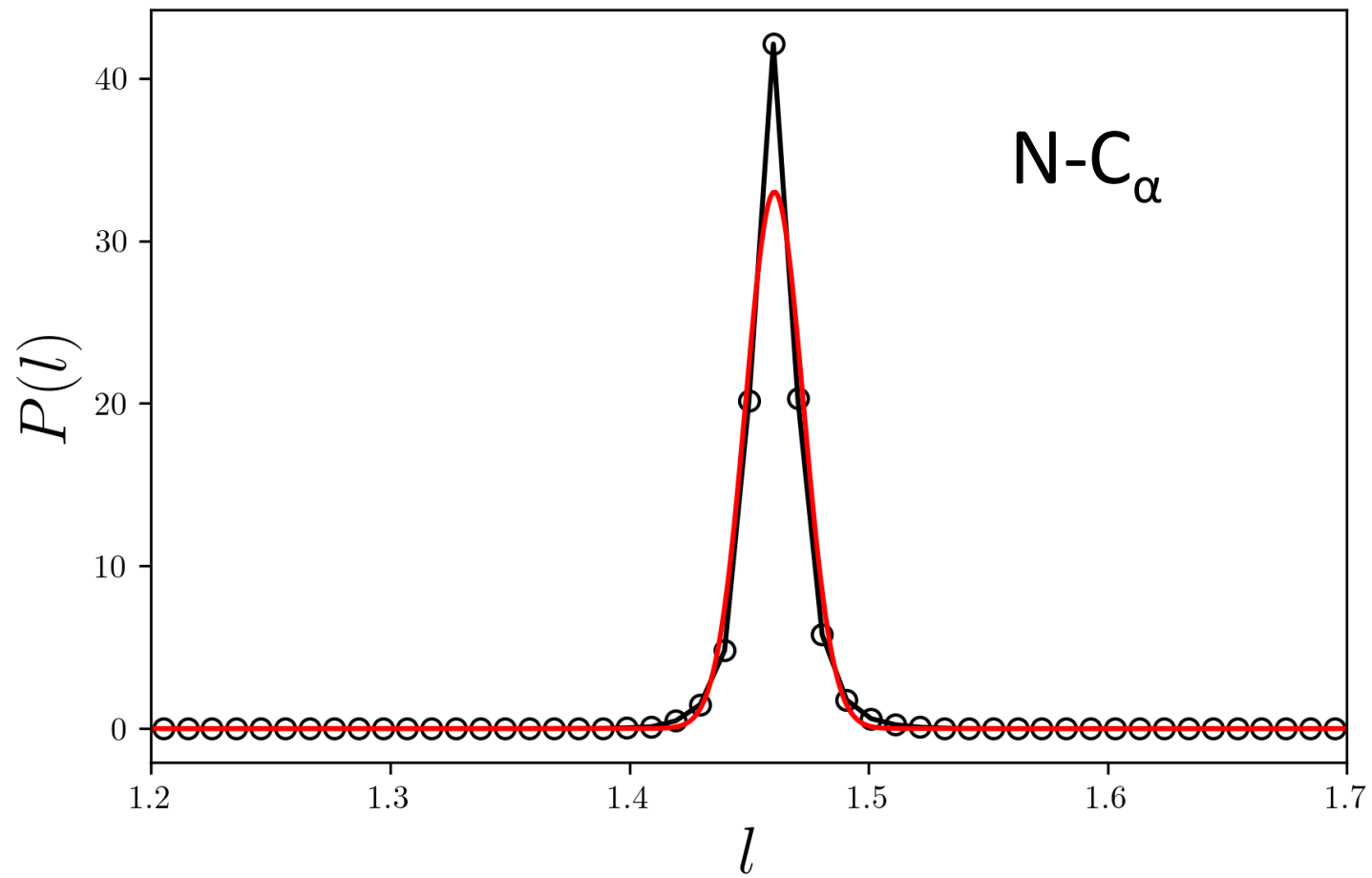




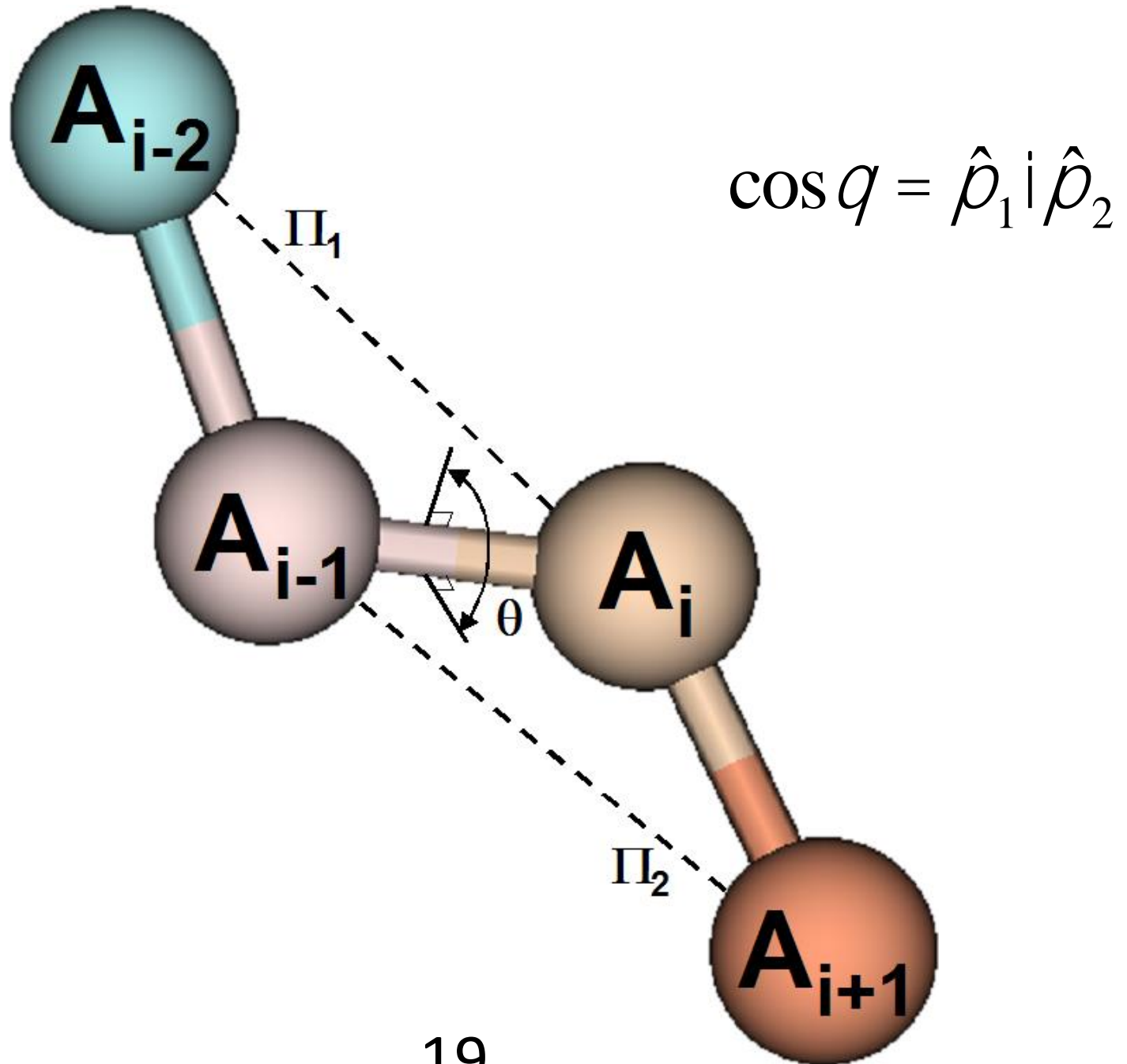


# Bond Lengths





# Backbone Dihedral Angles

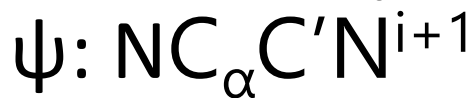
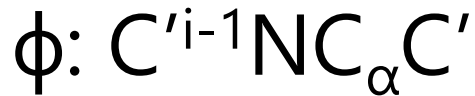
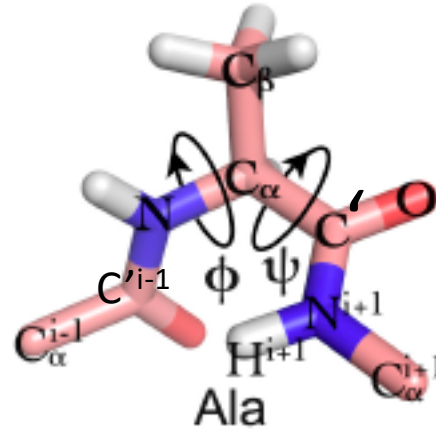


$3N-6$  DoF

$-(N-1)$  Bond lengths

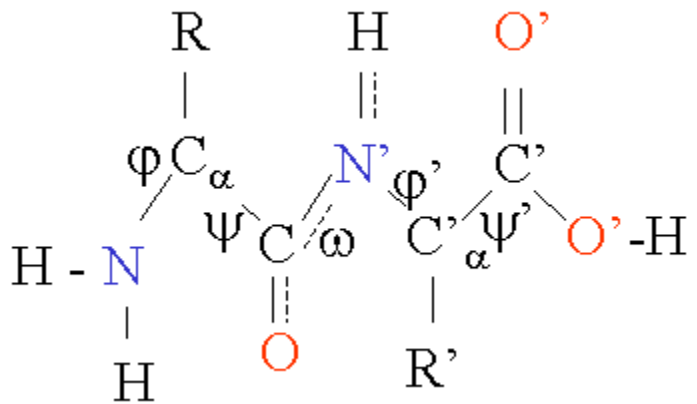
$-(N-2)$  Bond angles

$=N-3$  Dihedral angles

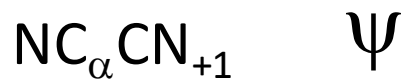
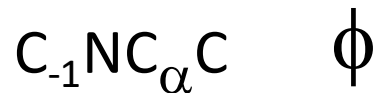


# Ramachandran Plot: Determining Steric Clashes

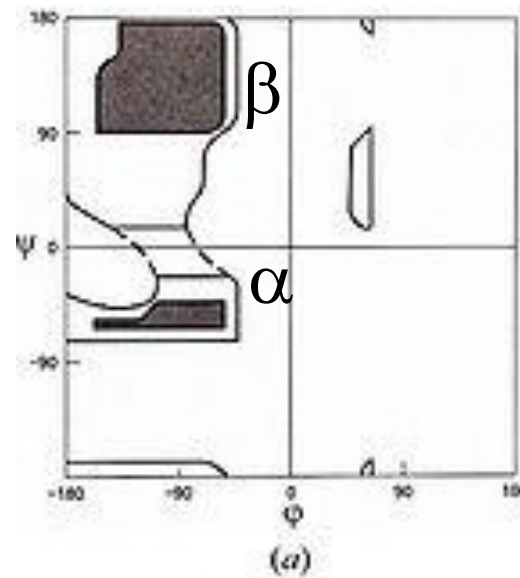
Backbone dihedral angles



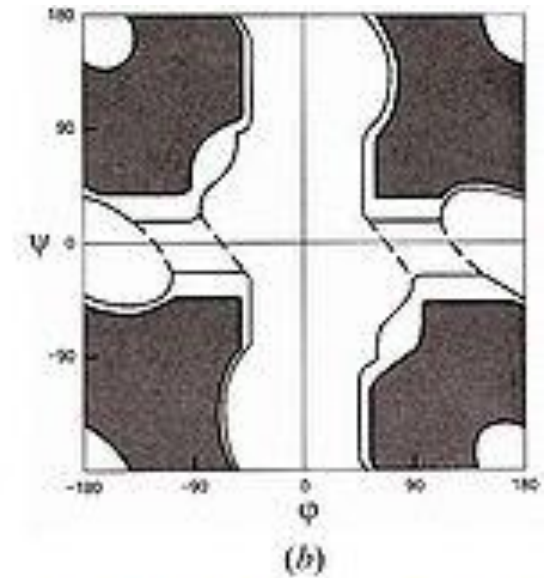
4 atoms define dihedral angle:



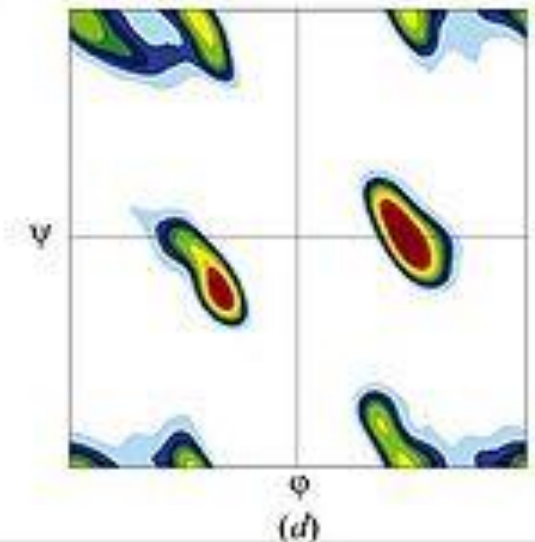
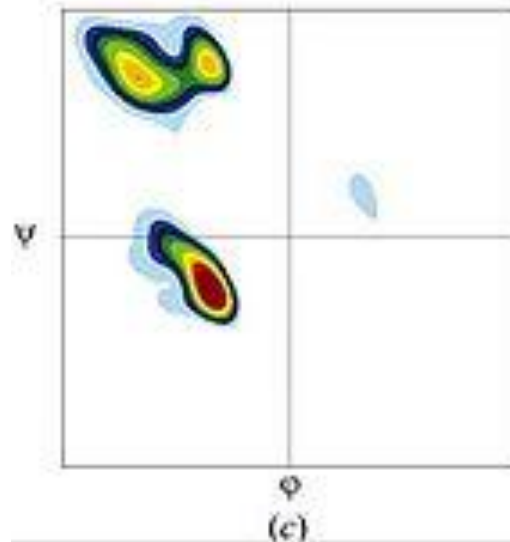
Non-Gly



Gly



theory



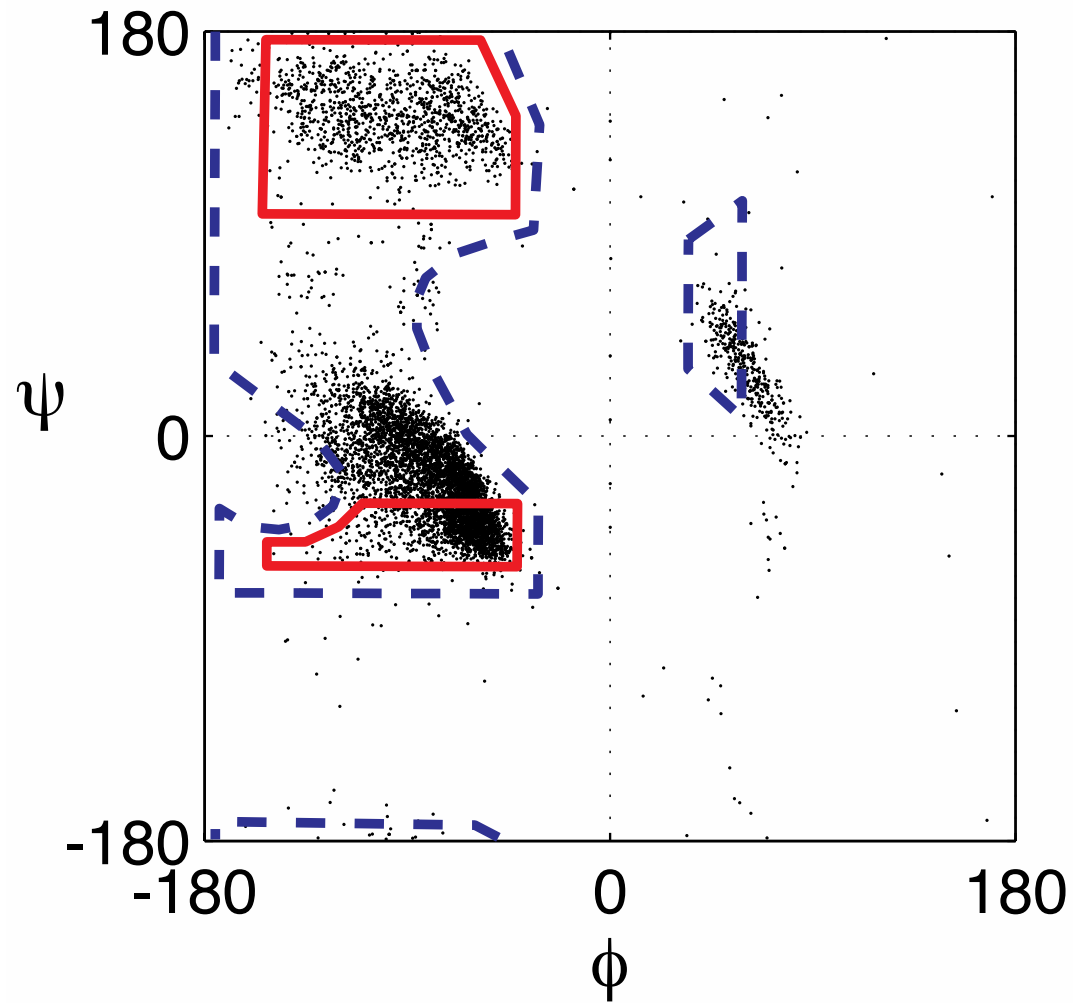
PDB

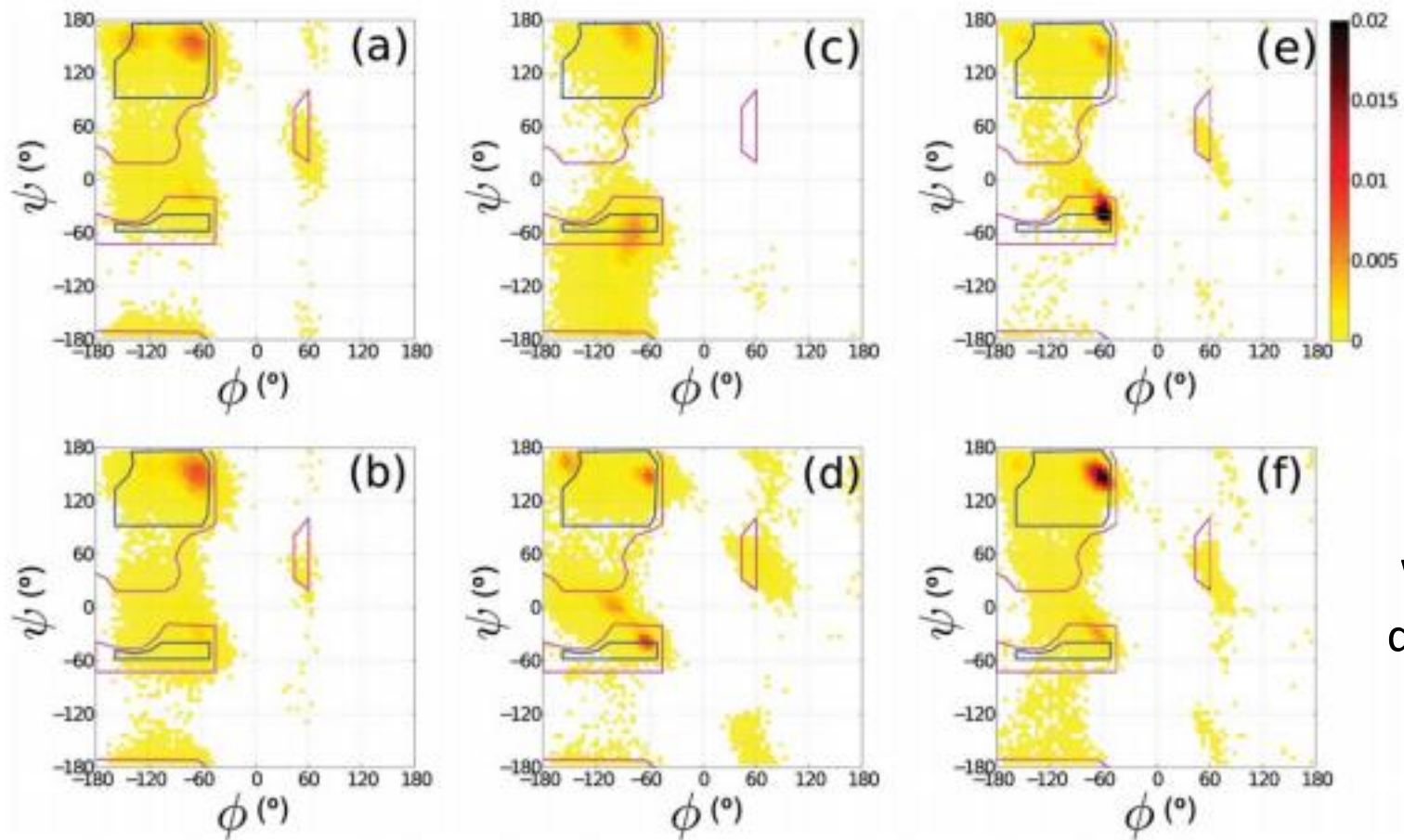
22

■ vdW radii  
— < vdW radii

--- backbone flexibility

# Backbone dihedral angles from PDB



Wu coil  
database

**Figure 5.** Probability distributions  $P(\phi, \psi)$  for the backbone dihedral angles  $\phi$  and  $\psi$  obtained from MD simulations of an Ala dipeptide mimetic using recent versions of the CHARMM and Amber force fields, their associated optimized water models, and with and without the "ILDN-NMR" and "CMAP" dihedral angle potential corrections: (a) Amber99sb + TIP4P-Ew, (b) Amber99sb-ILDN-NMR + TIP4P-Ew, (c) CHARMM27 + TIP3SP, and (d) CHARMM27-CMAP+TIP3SP. Subpanels (e) and (f) correspond to the Ala  $\phi$ - $\psi$  distributions from the Dunbrack Database<sup>3B</sup> and the Wu "Coil-3" library,<sup>10</sup> respectively. The Ramachandran hard-sphere<sup>3</sup> normal and outer limits (pink and blue lines, respectively) for  $\tau = 110^\circ$  are overlaid on each panel. The Amber and CHARMM MD simulations were thermally equilibrated at 303 K and sampled for 500 ns.



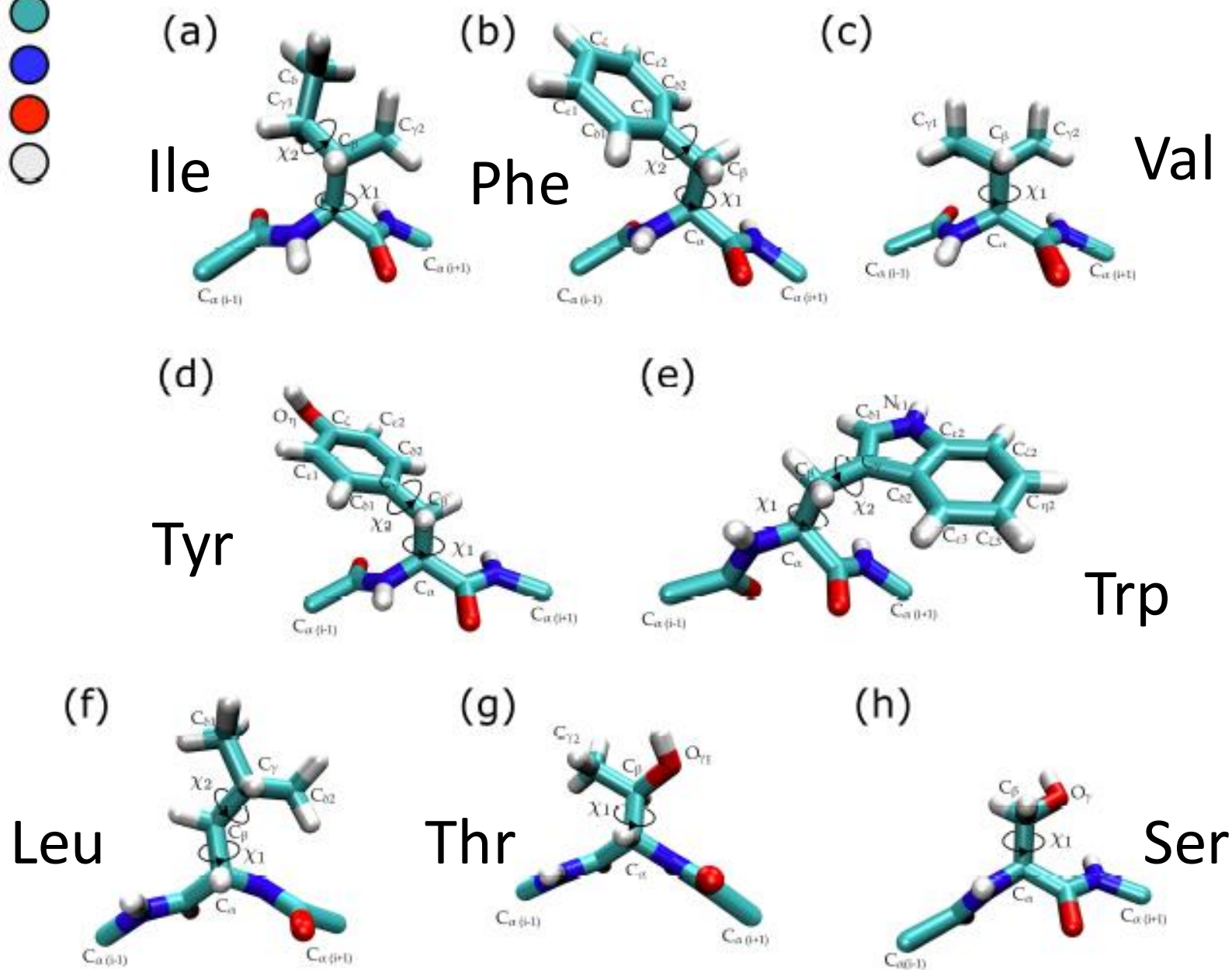
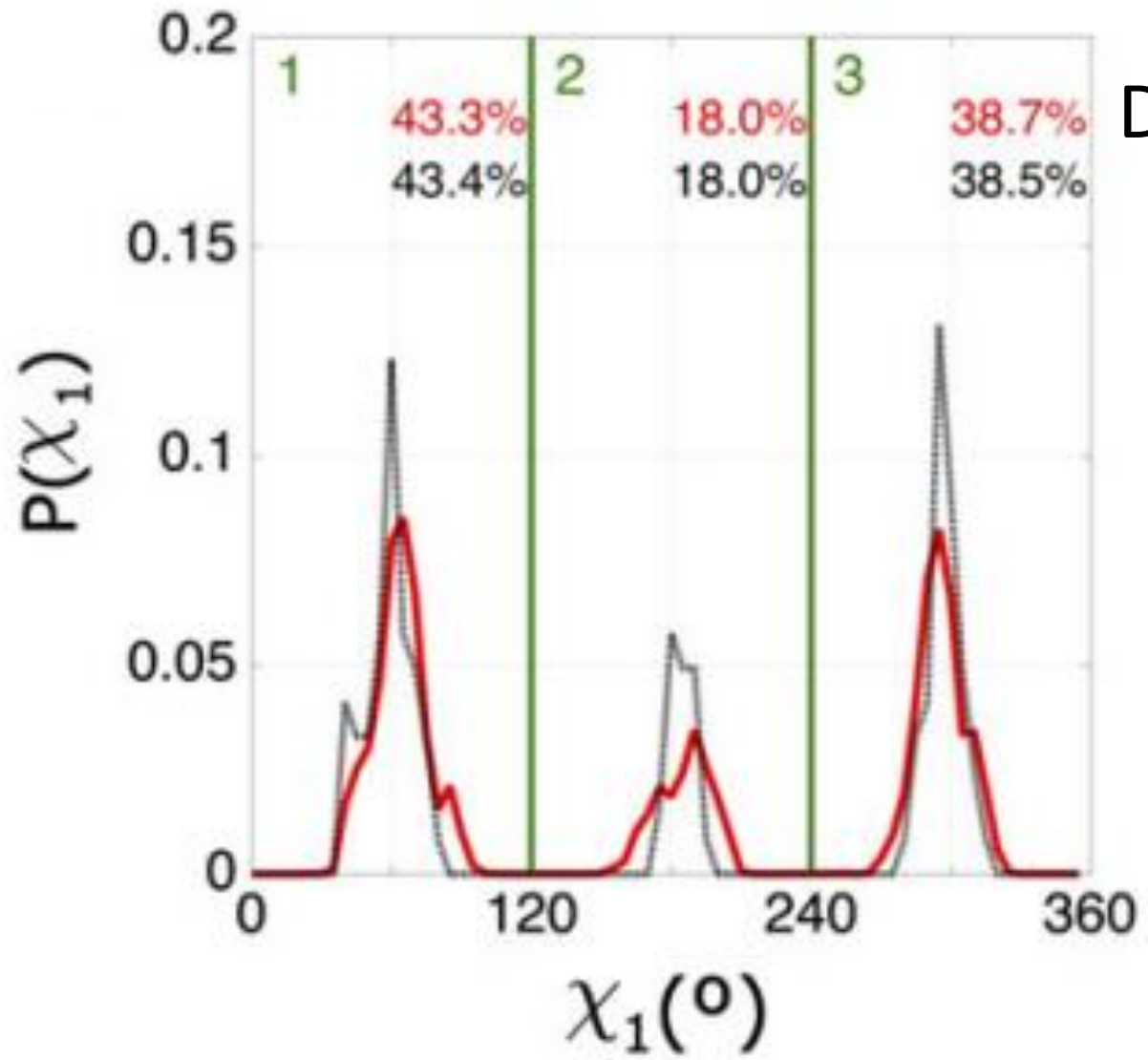


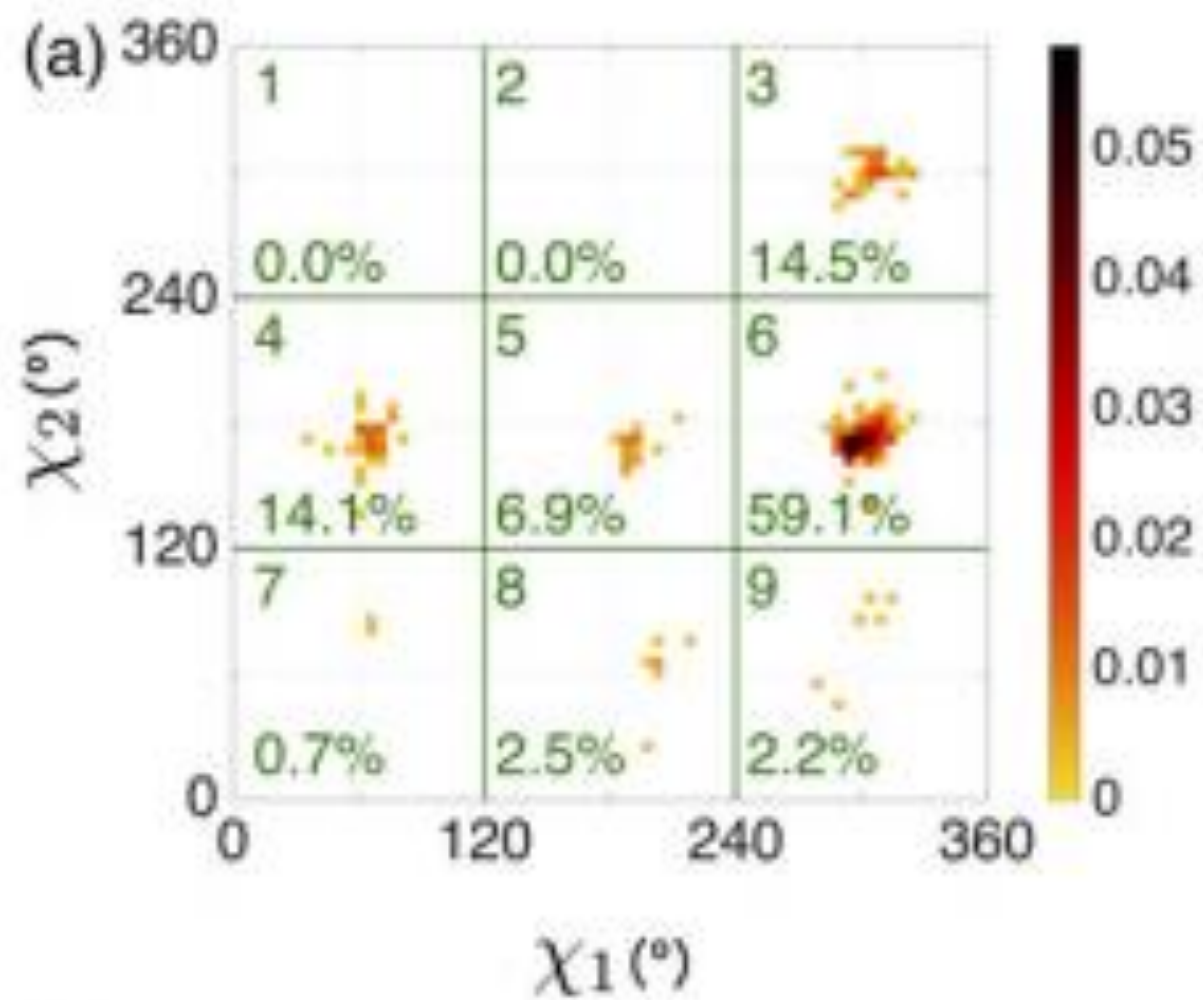
Figure S1: Stick representations of (a) Ile, (b) Phe, (c) Val, (d) Tyr, (e) Trp, (f) Leu, (g) Thr, and (h) Ser dipeptide mimetics. The carbon, nitrogen, oxygen, and hydrogen atoms are shaded green, blue, red, and white, respectively. The side chain dihedral angles  $\chi_1$  and  $\chi_2$  and several key atoms are labeled. The residues before ( $i-1$ ) and after ( $i+1$ ) the  $i$ th central residue are labeled at the  $C_\alpha$  atom.

Thr



Dunbrack 1.0

Ile



1. Can the structural properties of protein cores be quantitatively modeled using hard-spheres?
2. What is the packing fraction in protein cores?
3. Can simple hard-sphere model improve computational design of protein-protein interactions?