# What are proteins?



Linear polymer → Folded state
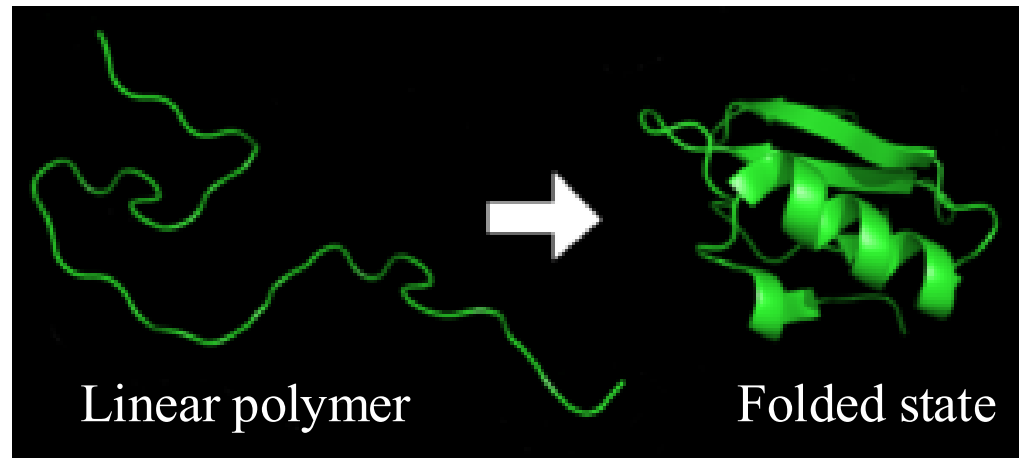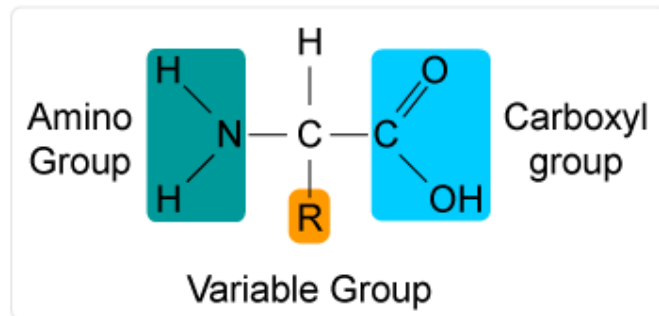
- Proteins are important; e.g. for catalyzing and regulating biochemical reactions, transporting molecules, …
- Linear polymer chain composed of tens (peptides) to thousands (proteins) of monomers
- Monomers are 20 naturally occurring amino acids
- Different proteins have different amino acid sequences
- *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)
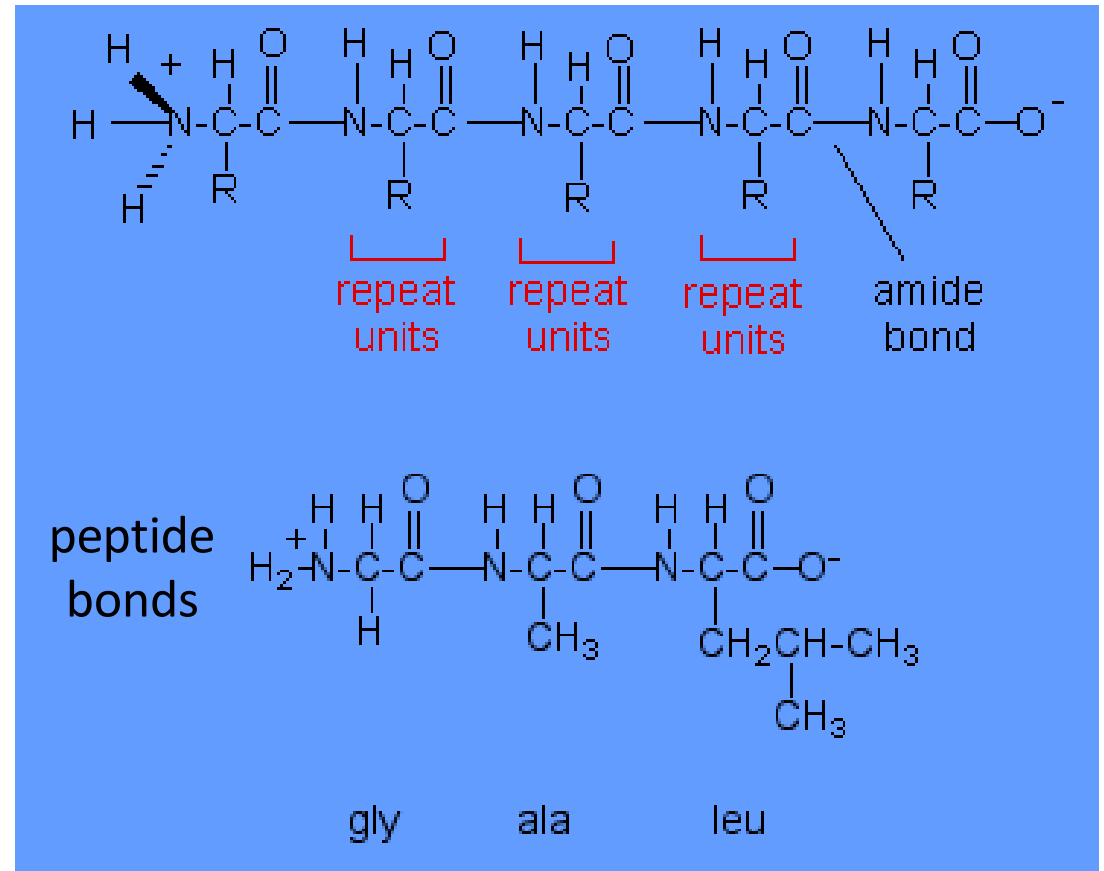- Proteins unfold or denature with increasing temperature or chemical denaturants

1
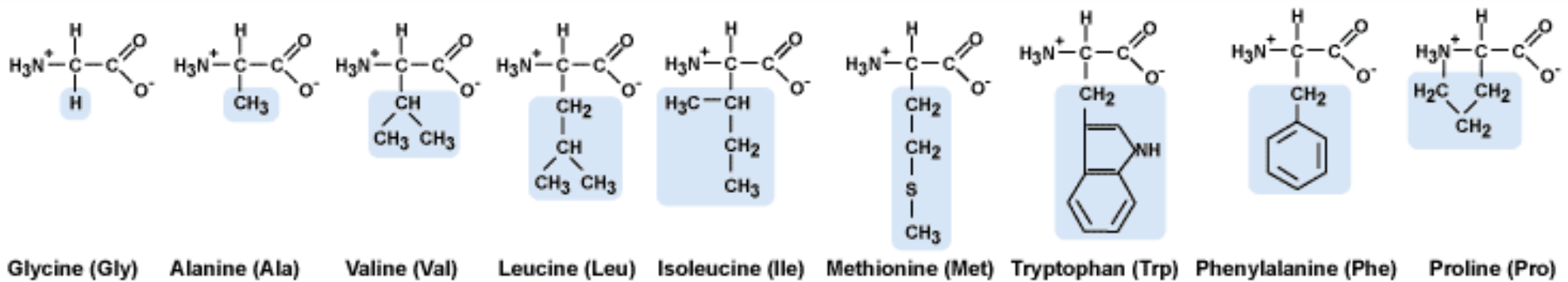
# Amino Acids I

General structure of Amino Acids



Amino Group

Carboxyl group

Variable Group

N-terminal     $C_\alpha$     C-terminal

R
variable
side chain

repeat units     repeat units     repeat units     amide bond

peptide bonds

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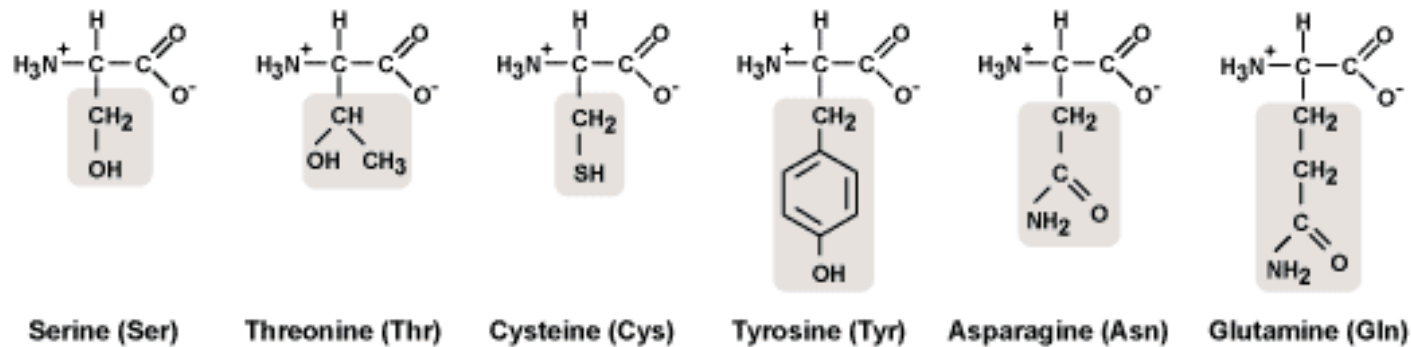
- Side chains differentiate amino acid repeat units
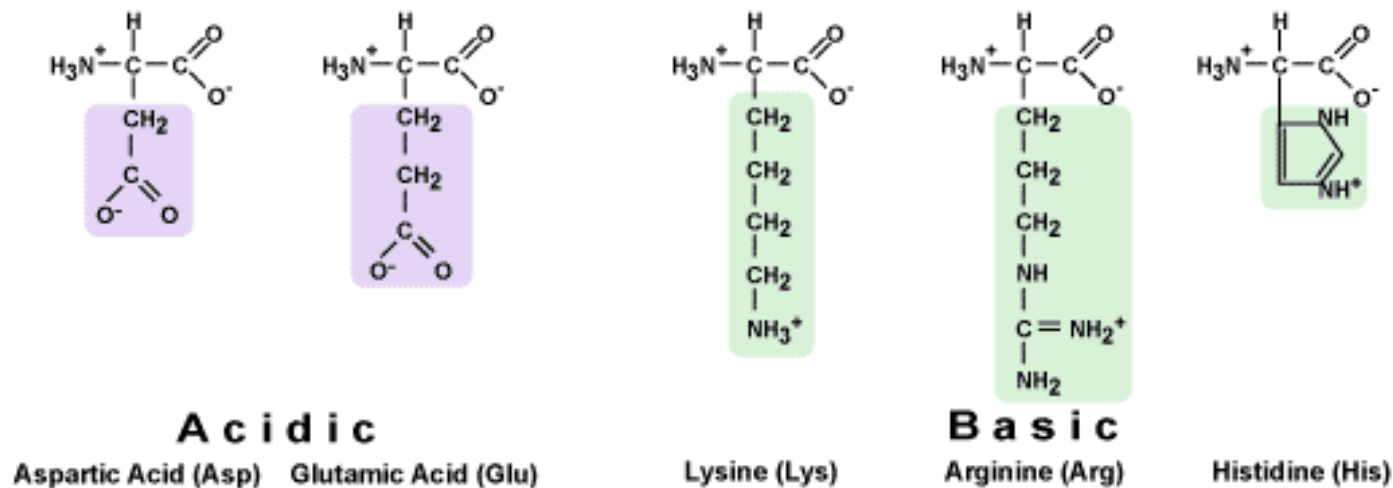- Peptide bonds link residues into polypeptides

2

# 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**

**A c i d i c**

Aspartic Acid (Asp)  Glutamic Acid (Glu)

**B a s i c**

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

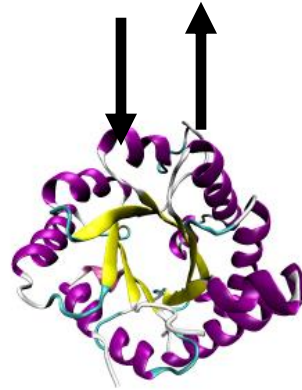(-)          3          (+)

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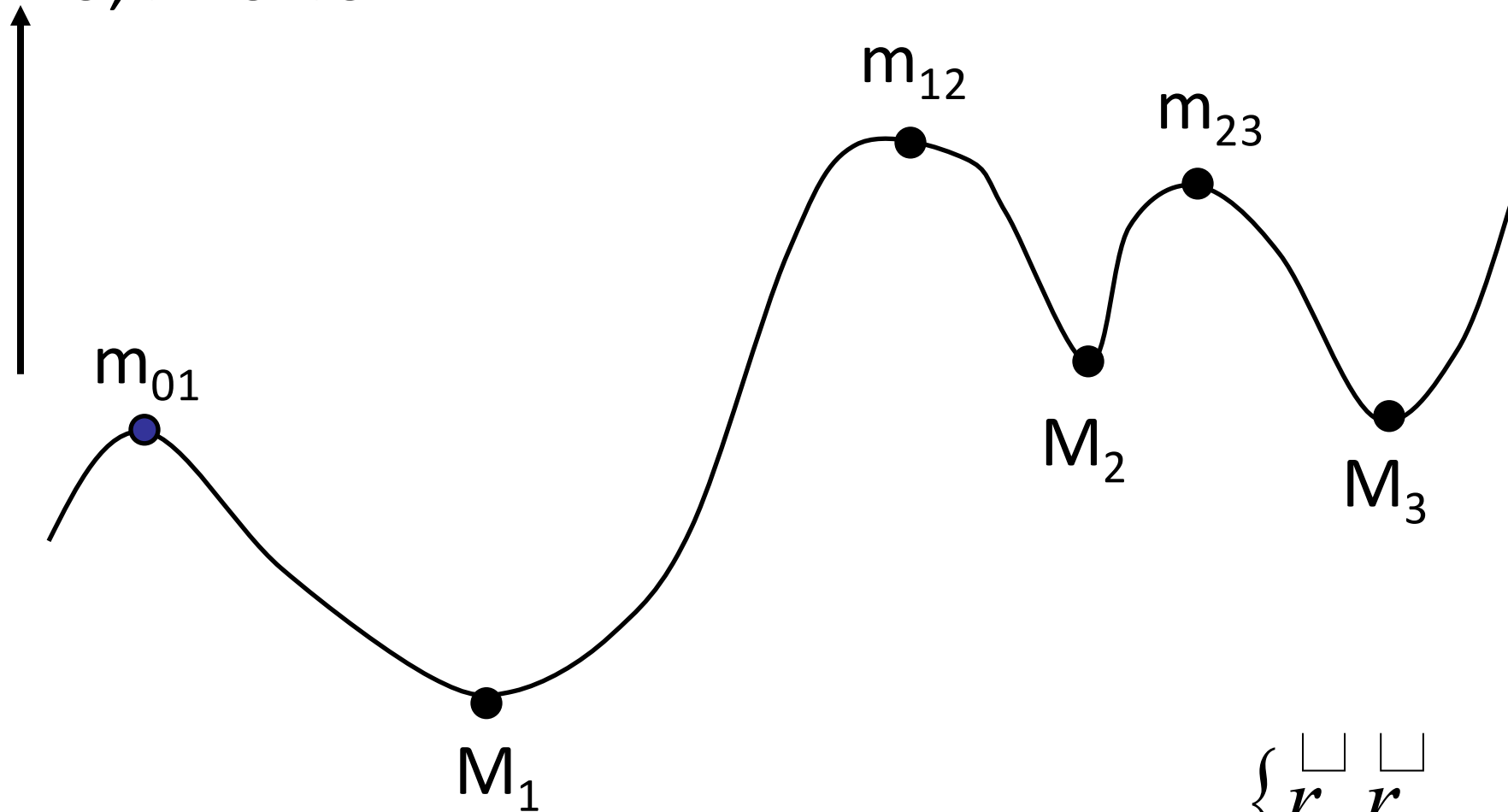
$$N_c \sim \mu^{2N}$$

$\mu$ = # allowed dihedral angles

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

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

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

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

5

# Energy Landscape

$U$, $F = U - TS$

$m_{12}$

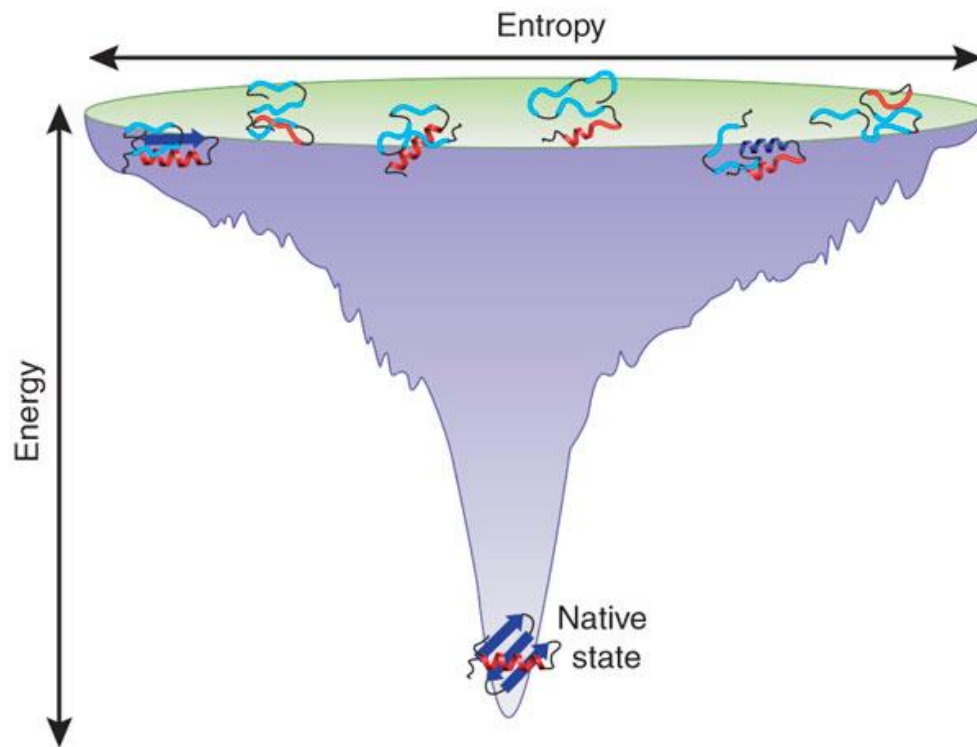$m_{23}$

$m_{01}$

$M_2$

$M_3$

$M_1$

$$\nabla U = 0 \quad \begin{cases} \nabla^2 U > 0 & \text{Minimum (M)} \\ \nabla^2 U = 0 & \text{saddle point} \\ \nabla^2 U < 0 & \text{Maximum (m)} \end{cases}$$

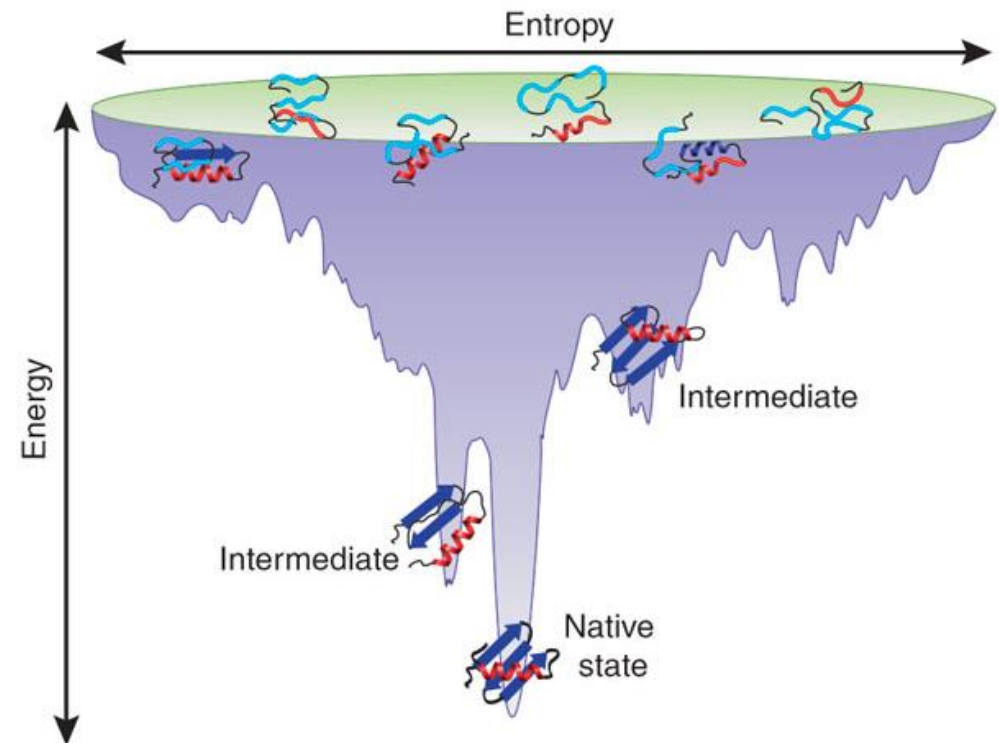$$\{ \vec{r}_1, \vec{r}_2, \ldots, \vec{r}_N \}$$

all atomic
coordinates;
dihedral angles

6

# Roughness of Energy Landscape



smooth, funneled

(Wolynes et. al. 1997)

rough

# Critical Assessment of Structure Prediction (CASP)

Hoval et. al., *Protein Science* (2018)
Moult et. al., *Protein Science* (2018)

# Driving Forces

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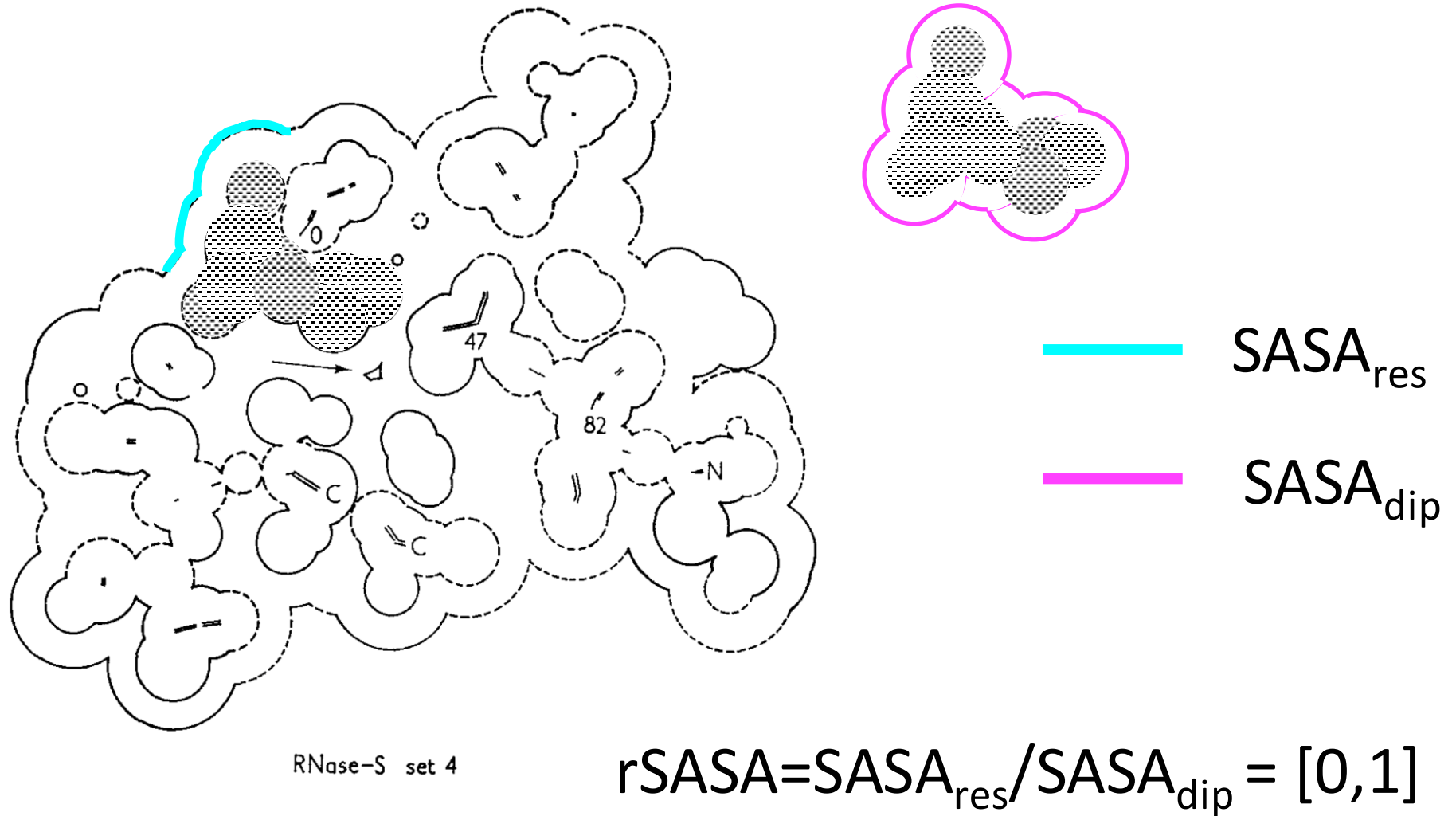
inside     H (hydrophobic)

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).
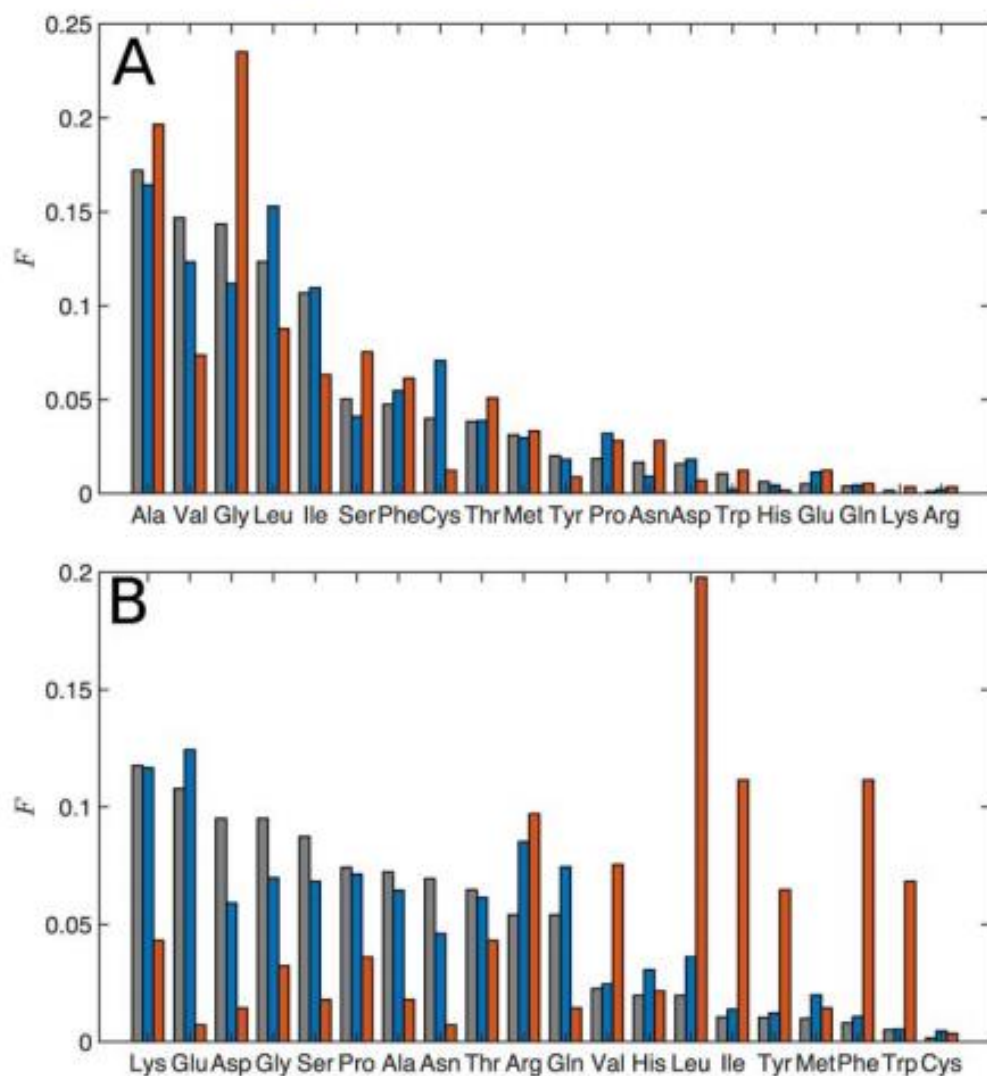
9

# Solvent Accessible Surface Area and rSASA



RNase-S set 4

— $SASA_{res}$

— $SASA_{dip}$

$rSASA = SASA_{res}/SASA_{dip} = [0,1]$

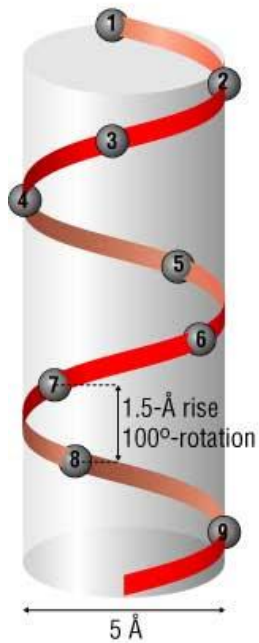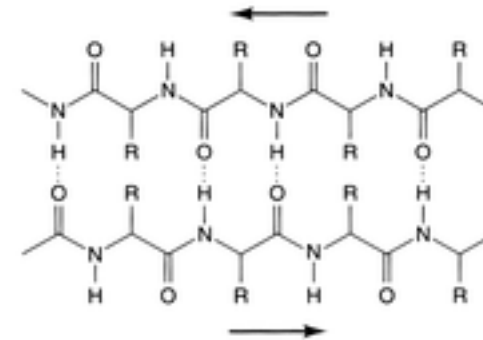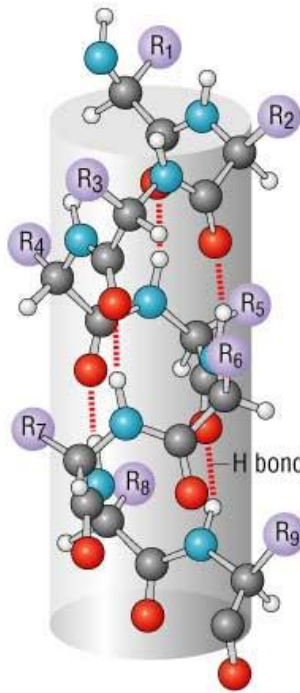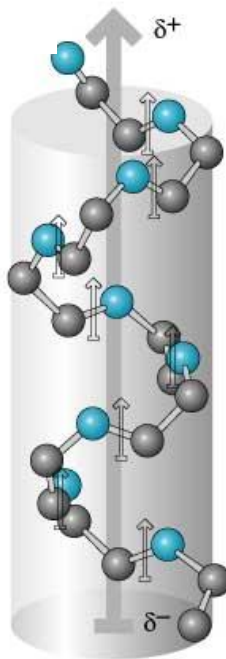**FIGURE 5** Fractions of amino acids with A, rSASA $\leq 10^{-3}$ and B, rSASA>0.5 for residues in the Dun1.0 (grey), PPI (blue), and TM (red) datasets. The fractions are defined relative to the total number of residues in each rSASA category. C, The fractions of core residues (light bars) and non-core 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]

11

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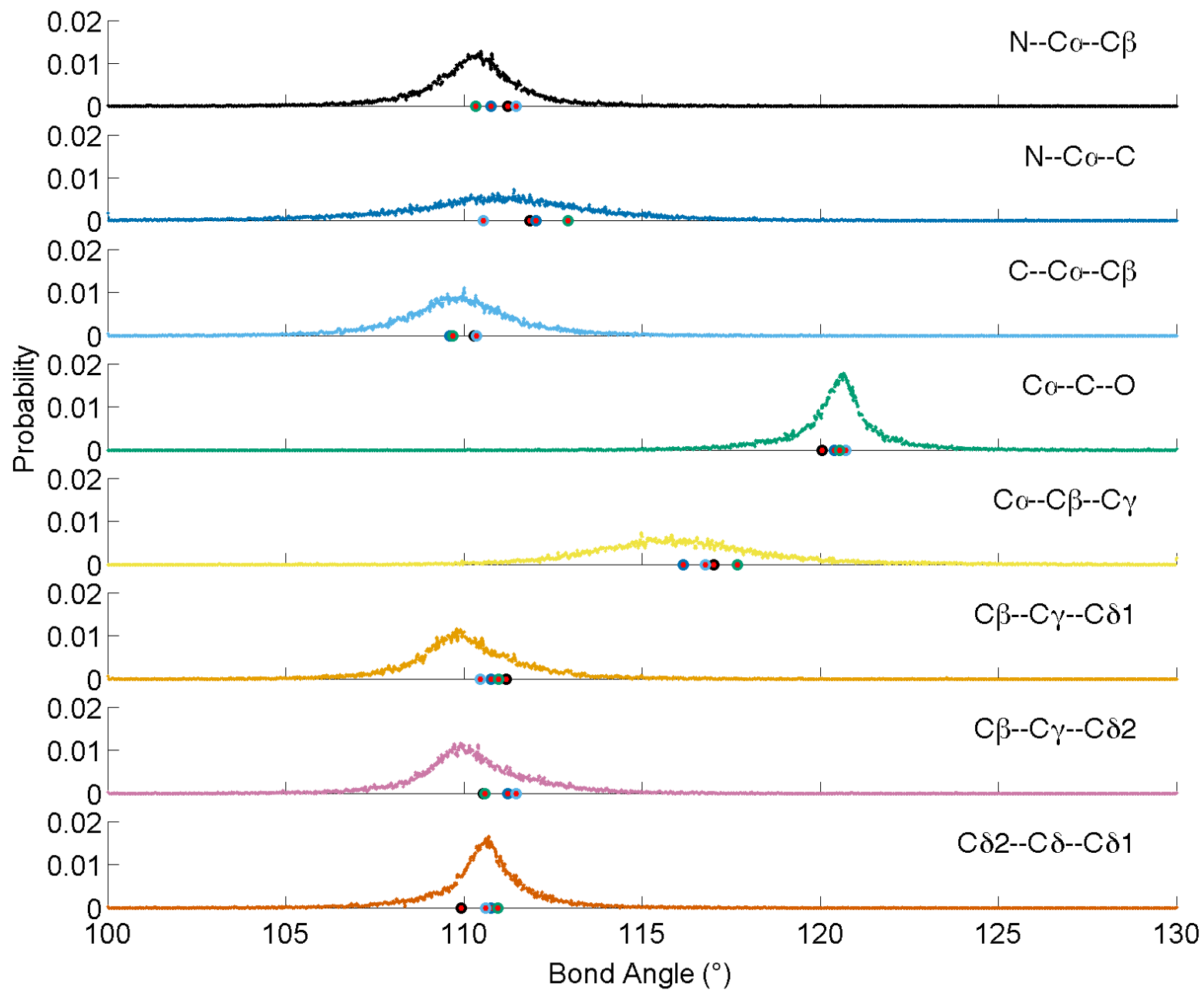
## $\alpha$-helix



$\beta$-strand

$\beta$-sheet

5Å

- 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
- $(\phi,\psi)=(-60°,-45°)$
- $\alpha$-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
- $(\phi,\psi)=(-135°,135°)$
- $\beta$-strand propensities: Val, Thr, Tyr, Trp, Phe, Ile
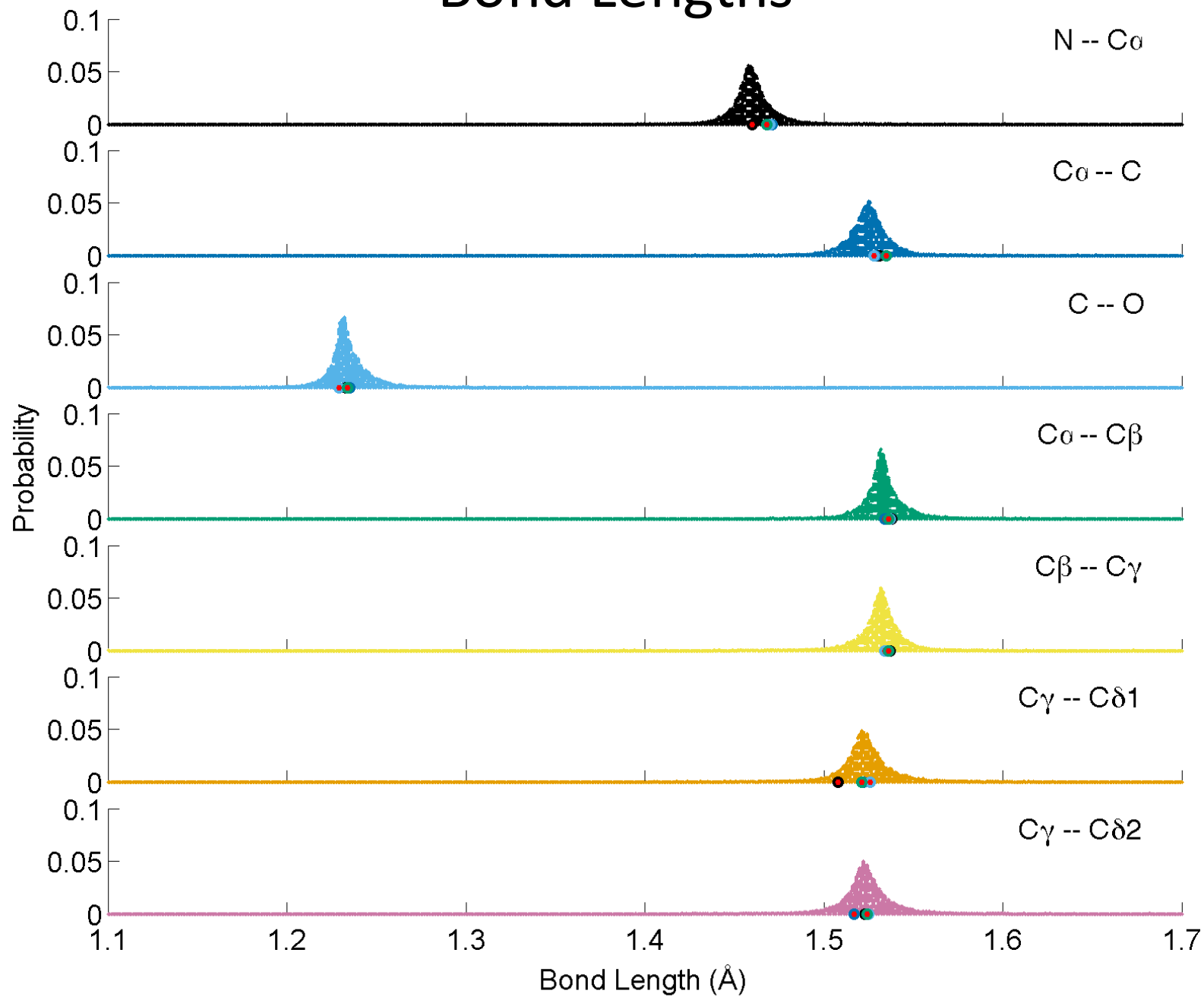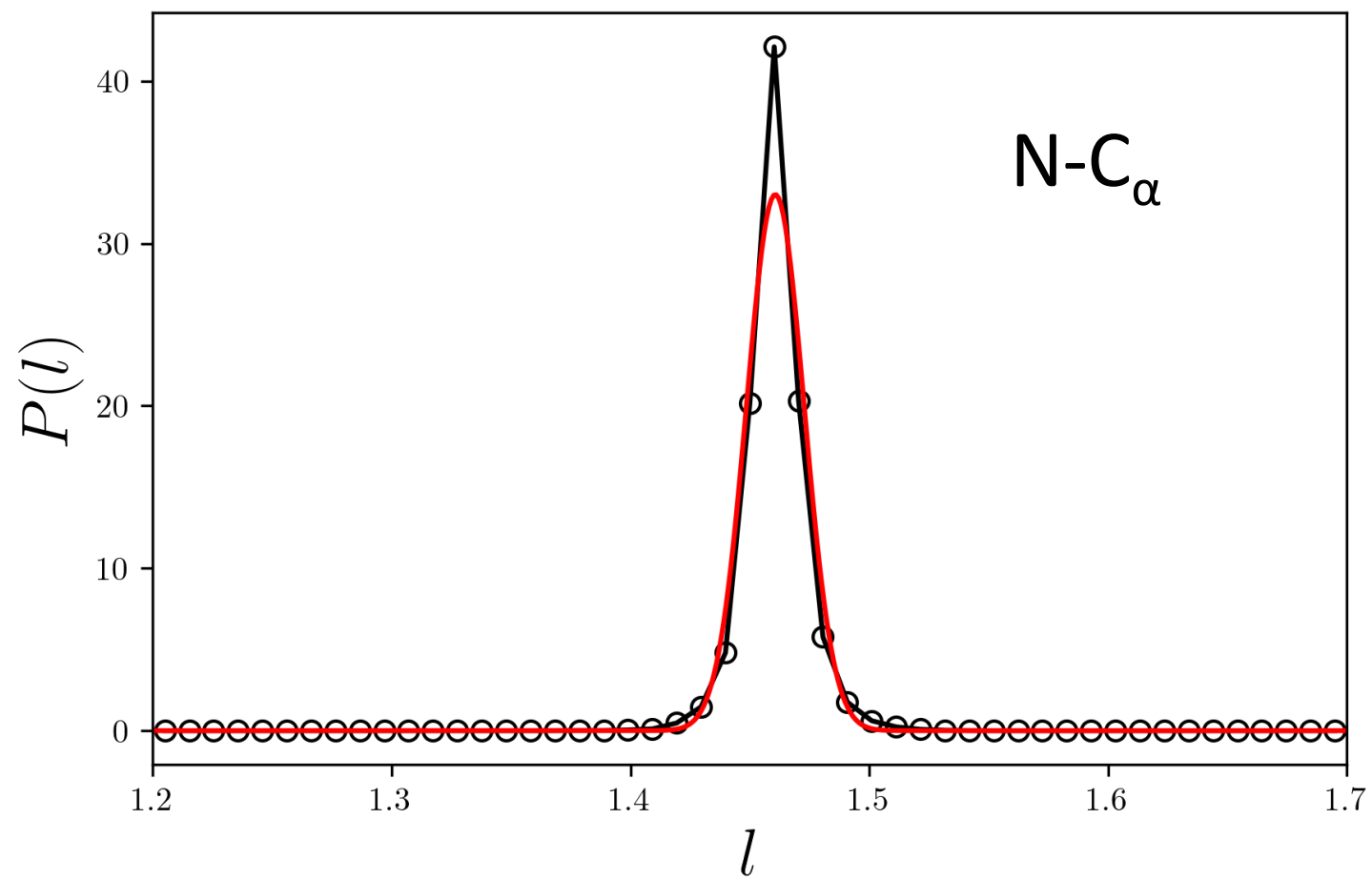
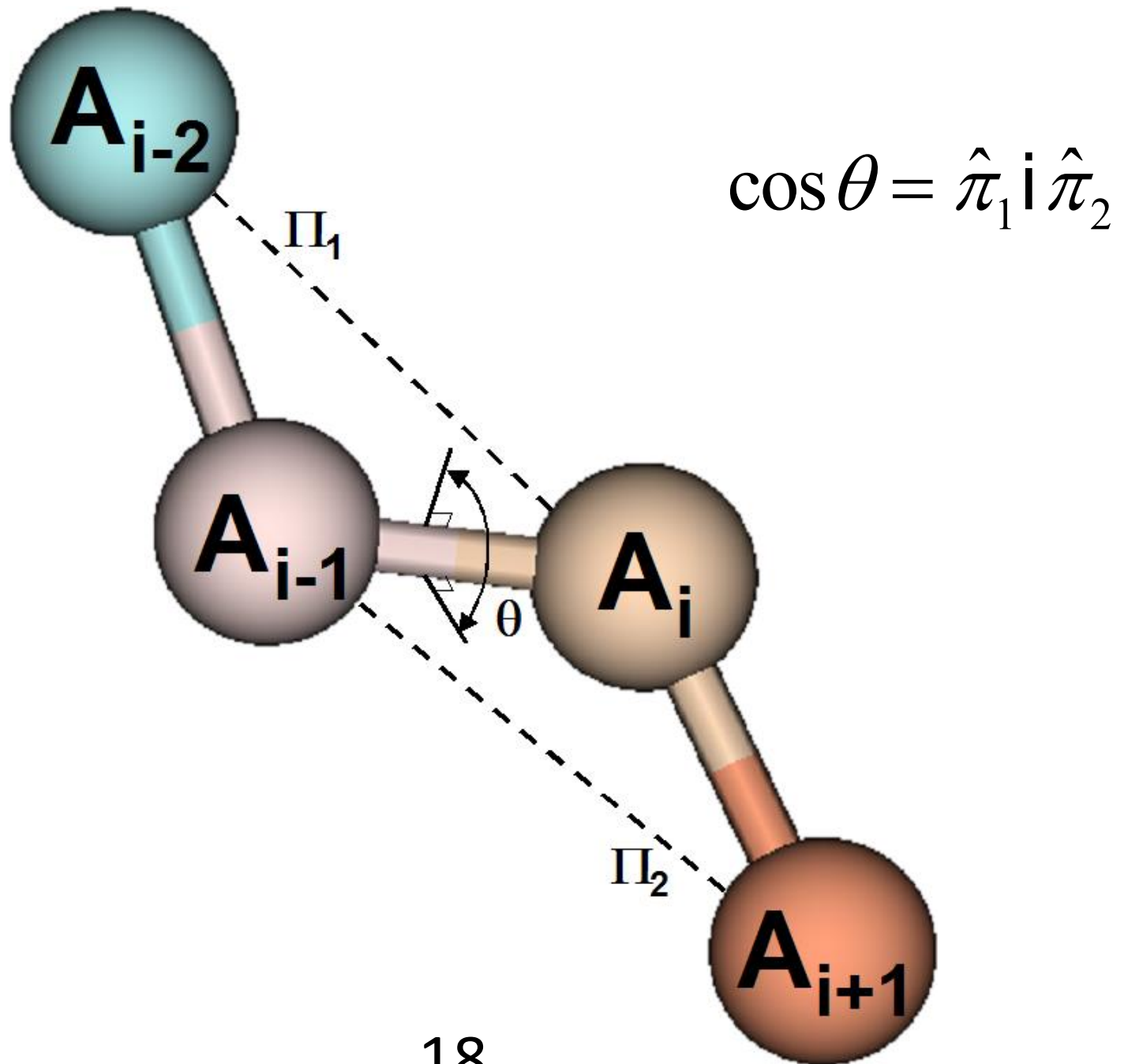12

# $N_s$=62,938 monomeric xtal structures

# Bond Angles

# Bond Lengths
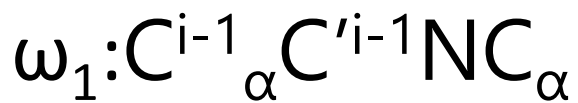
# Backbonde Dihedral Angles
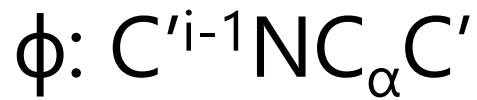


$$\cos\theta = \hat{\pi}_1 \mathbf{i} \hat{\pi}_2$$

3N-6 DoF

-(N-1)　　Bond lengths

-(N-2)　　Bond angles

=N-3　　　Dihedral angles

19

$\phi$: $C'^{i-1}NC_{\alpha}C'$

$\psi$: $NC_{\alpha}C'N^{i+1}$

$\omega_1$: $C^{i-1}_{\alpha}C'^{i-1}NC_{\alpha}$

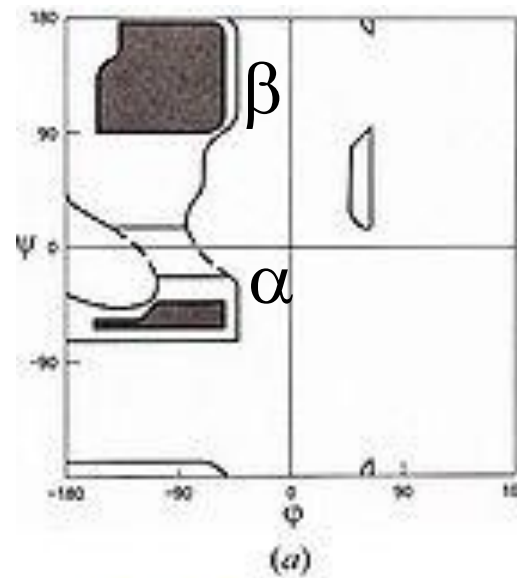$\omega_2$: $C_{\alpha}C'N^{i+1}C^{i+1}_{\alpha}$

20

# Ramachandran Plot: Determining Steric Clashes

Backbone
dihedral angles



4 atoms define dihedral angle:

$$C_{-1}NC_{\alpha}C \qquad \phi$$

$$C_{\alpha,-1}C_{-1}NC_{\alpha} \qquad \omega=0,180°$$

$$NC_{\alpha}CN_{+1} \qquad \psi$$

Non-Gly       Gly



theory

PDB

21

◻ vdW radii

— < vdW radii

- - - backbone flexibility
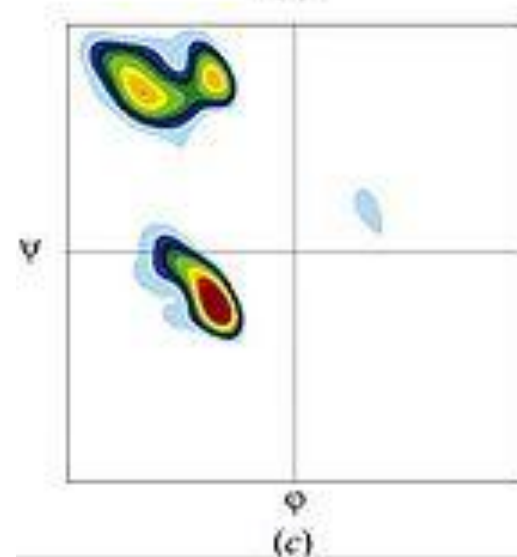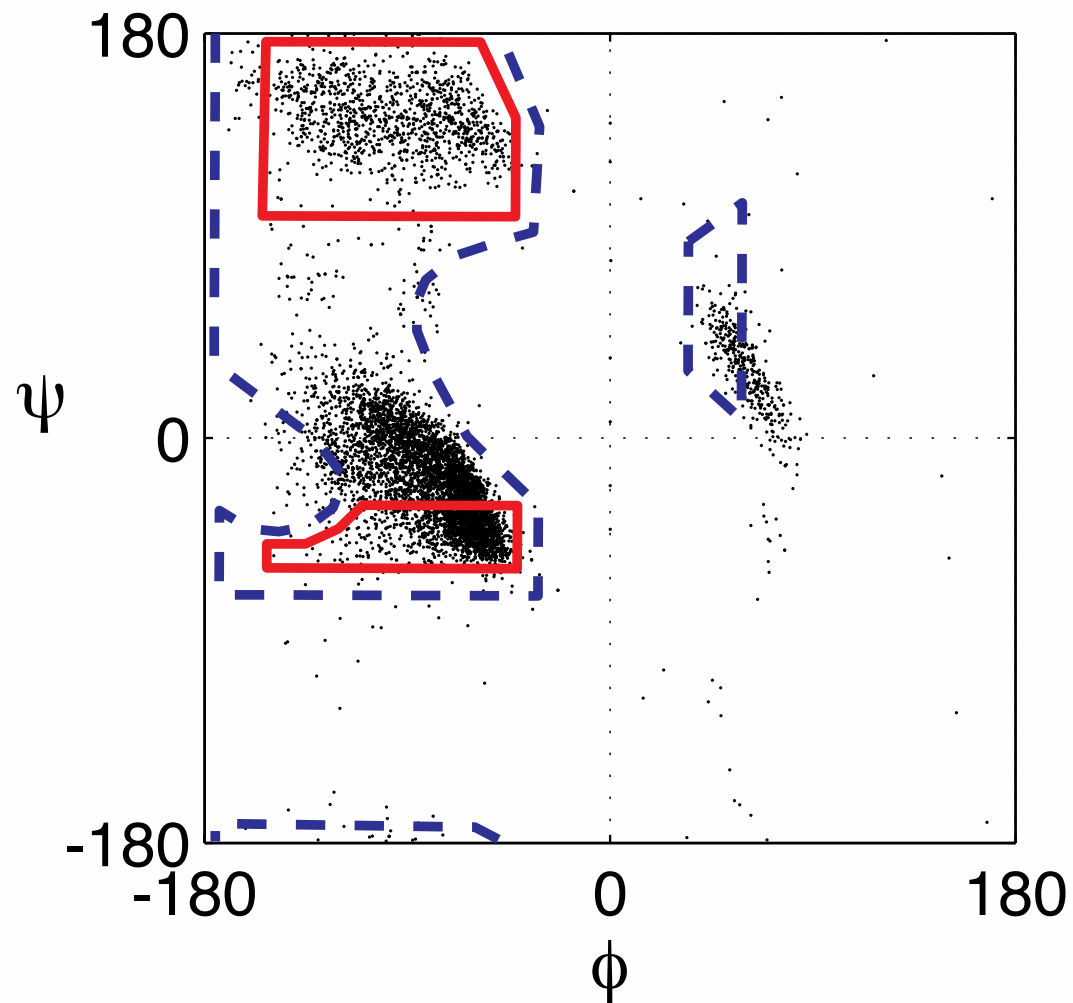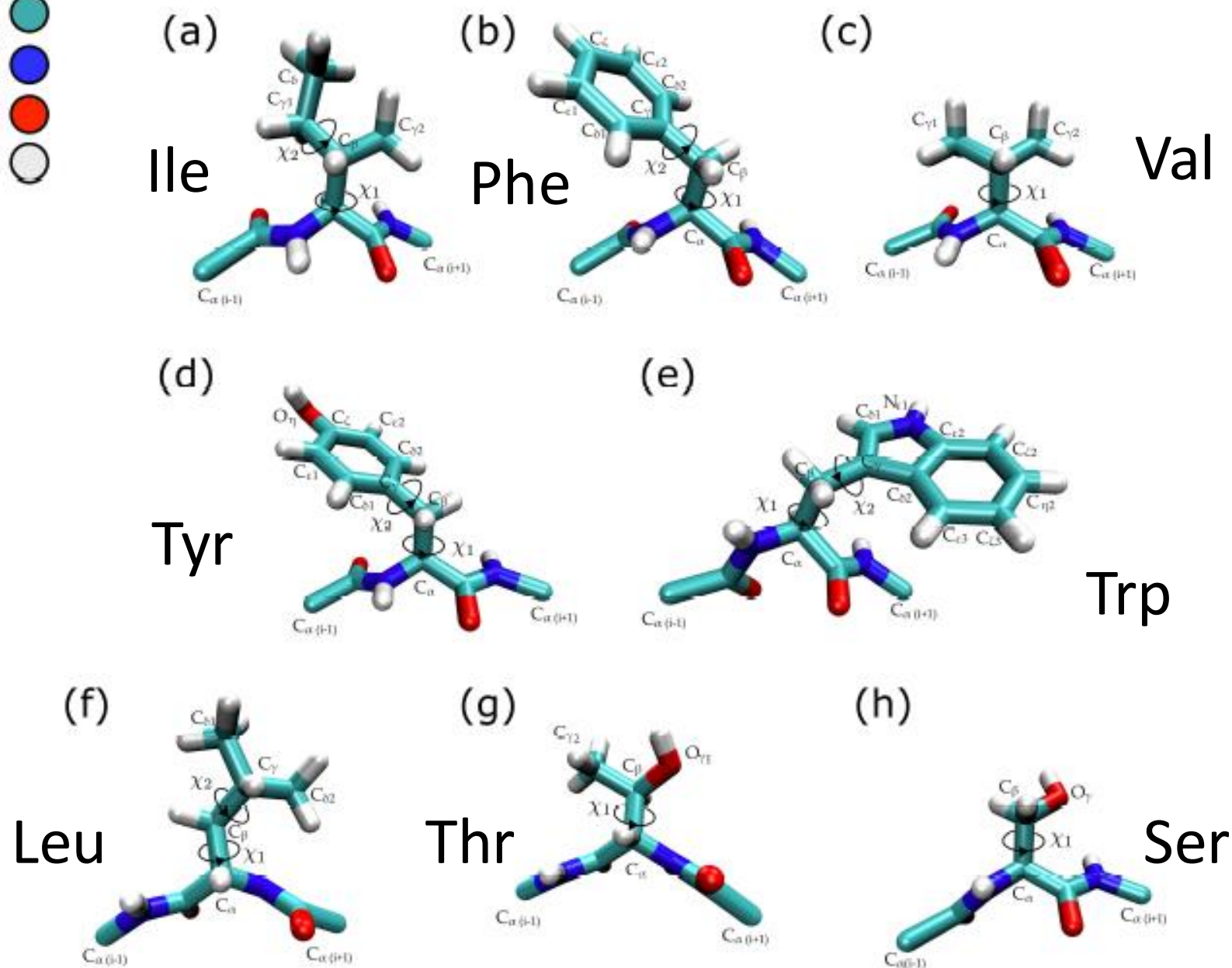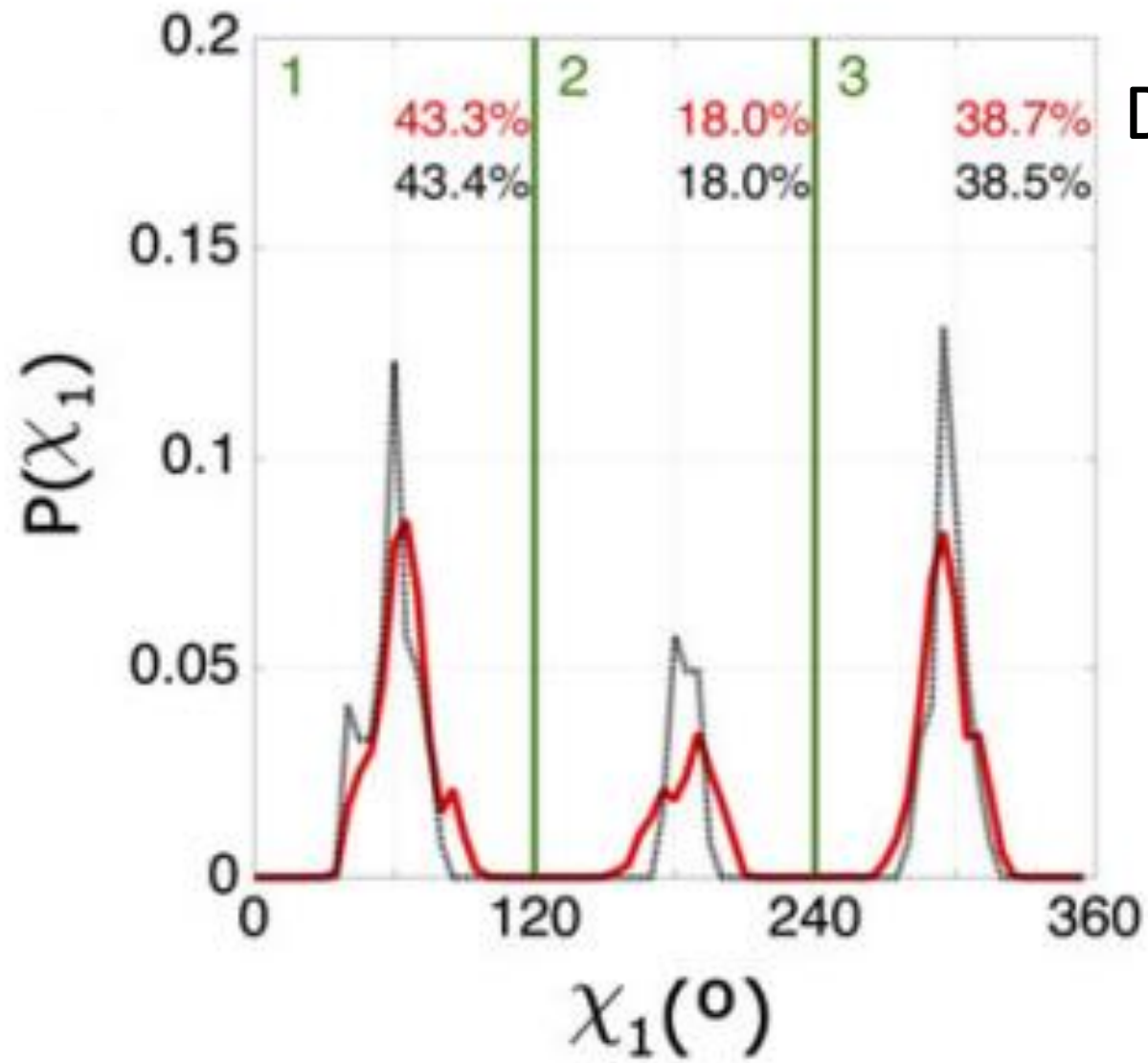
# Backbone dihedral angles from PDB

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



(a)