



INTERNATIONAL ATOMIC ENERGY AGENCY  
UNITED NATIONS EDUCATIONAL, SCIENTIFIC AND CULTURAL ORGANIZATION  
**INTERNATIONAL CENTRE FOR THEORETICAL PHYSICS**  
I.C.T.P., P.O. BOX 586, 34100 TRIESTE, ITALY, CABLE: CENTRATOM TRIESTE



H4.SMR/775-16

**COLLEGE IN BIOPHYSICS:  
EXPERIMENTAL AND THEORETICAL ASPECTS OF  
BIOMOLECULES**

**26 September - 14 October 1994**

*Miramare - Trieste, Italy*

***Structure Prediction Methods in Protein***

**Nicholas D. Soccia  
University of California at San Diego  
La Jolla, - CA, USA**

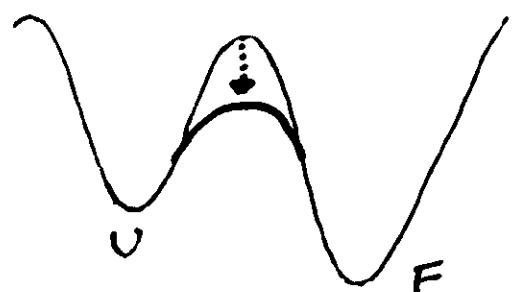
Thermodynamic Hypothesis: seq encodes Struct by determining Free energy Min.

- For some proteins Kinetic factors may also be important.

Although Seq encodes Struct, some proteins need help to fold.

- Disulfide Bonds, normally too strong to break spontaneously. Anfinsen used  $\beta$ -MOH. Cells use Disulfide Isomerase Enzyme with free Cysteine on surface  $\downarrow$   
But faster than
- Prolines  $cis \rightleftharpoons trans$  also slow compared to other steps: prolyl isomerase.

These two work like any other enzyme: lower free energy barrier



In cells (in vivo) additional problems: Aggregation

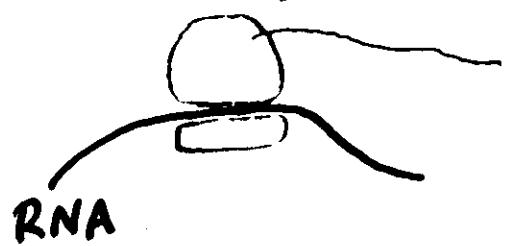
Hydrophobic effect drives not only folding but aggregation of several chains



this reduces surface area even

more so favored over folding. (Also problem in  
vitro, for some  
proteins)

- Cells have high protein concentrations
- Transcription much slower than folding (or aggregate)



For a 100 AA protein  
at the end of transcription 30-40  
AA still in ribosome

(Heat shock, related)

Cells have Chaperones which bind to unfolded proteins.

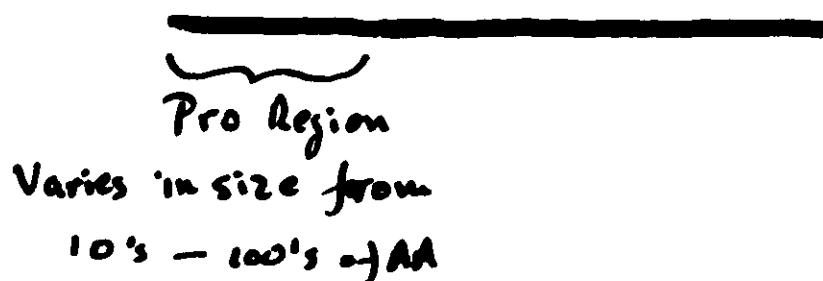
DnaK - Binds to End out of Ribosome

GroEL - Binds to Free proteins

- Many Functions - Heat Shock
    - Translocation
    - Assist in formation of oligomers
- } Precise mechn  
not fully  
understood  
\* SOME USE ATP

## Interesting Kinetics:

① Proteases: Synthesized in an Inactive form (Zymogen), <sup>usually</sup> consisting of an extra piece



To Activate must cleave or cut off pro-region.

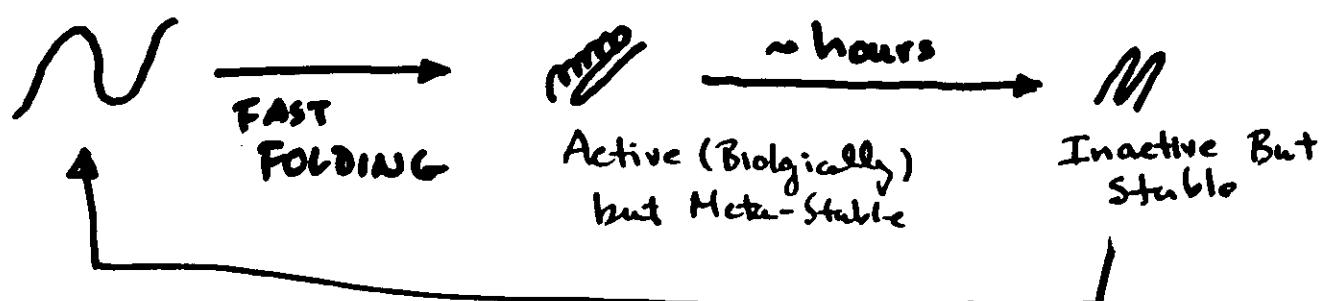
However to fold need pro-region

Examples : • Chymotrypsin (just cleave)  
• α-Lytic (166 AA pro-region)

↓  
Lowers barrier

② Serpins - "Kinetic" Switch

Protease  
Inhibitors



Can Repeat Cycle By Denaturing  
and Renaturing

③ Influenza Virus pH barrier (Not Thermo)

What does this mean for the thermodynamic hypothesis?

① Recall Diversity (a.k.a. Annoying Exception)

Always seem to be special cases  
in biology.

② Are the above special: YES  $\xrightarrow{\text{IMHO}}$

- All seem to have specific FUNCTIONAL Purposes
- Many, Many proteins have the much simpler behavior ( $\sim 25$  yr of experiments)

Some disagree: Baker, Biochemistry 33, 7505 (1994)  
feel it may be a more, general feature

③ Restricted Thermo. Hyp.  $\rightarrow$  Global Free Eny min of accessible states.

But how long are you willing to wait??  
(Finite system, everything accessible)

④ No one has shown, "True" metastability  $\stackrel{\text{My definition}}{=} \text{NO PATHWAYS}$  (ala Leventhal)

Leave issue of folding for now...

Look at the folded structure of proteins

Not only compact but highly ordered

### Structural Hierarchy

- Primary = sequence
- \* SuperSecondary ← [ • Secondary \*
- Domains                          • Tertiary = folded struct
- Quaternary = multi unit systems

Secondary Structure = Local ordering  
(conformation) of chain

- Regular 2° Struct 50-60%
- helix -  $\alpha$  most common ( $\frac{10}{\pi}$ ,  $\frac{3.6}{\text{res}}$ , triple helix, collagen)
  - chiral right handed [no extended left handed]
  - Sheets,  $\beta$  parallel anti-parallel
  - Turns (3 types) → at surface
  - Random Coil

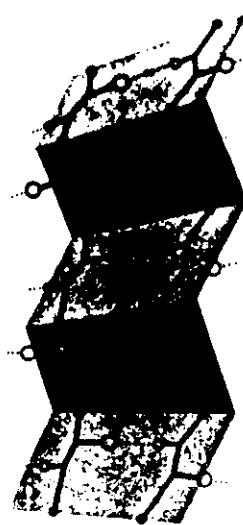
**Figure 1-19**  
Hierarchies of structure.



PRIMARY STRUCTURE  
(Amino acid sequence in the protein chain)



$\alpha$  helix



$\beta$  sheet



Domains (dark color) in  
an antibody molecule

LOCAL FOLDING

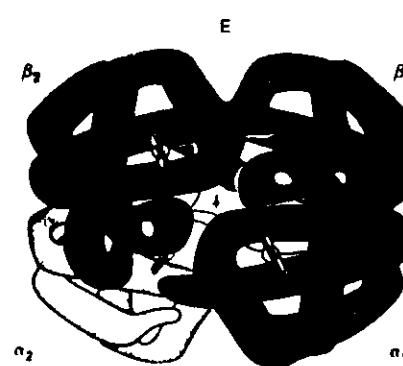
E

SECONDARY STRUCTURE



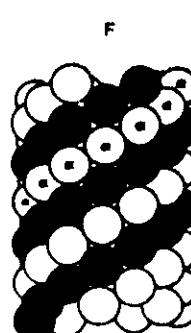
TERTIARY  
STRUCTURE

One complete protein chain  
( $\beta$  chain of hemoglobin)



QUATERNARY STRUCTURE

The four separate chains  
of hemoglobin assembled  
into an oligomeric protein



MACROMOLECULAR  
ASSEMBLY

$\alpha$  (white) and  $\beta$  (color)  
tubulin molecules in a  
microtubule.

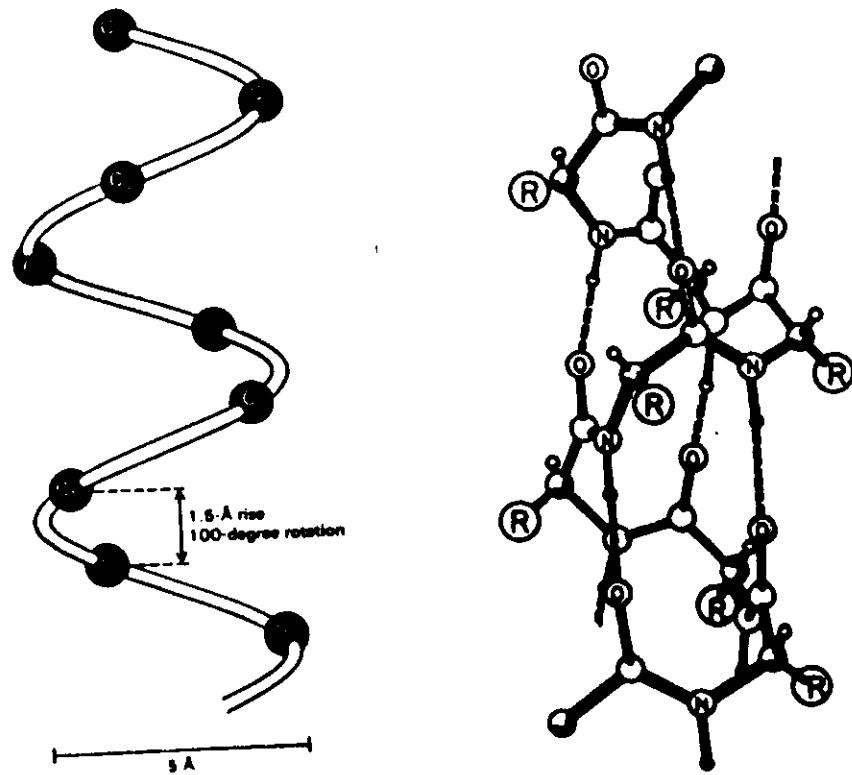


Figure 4.4: Two representations of the  $\alpha$ -helix. The left picture is a simple structure showing only the  $\alpha$ -carbons and the trace of the main chain. The picture on the right is an (almost) all atom picture showing all but the side chains in detail. Note the hydrogen bonding between every fourth residue. Figures from [77, 69].

**PROOFS**

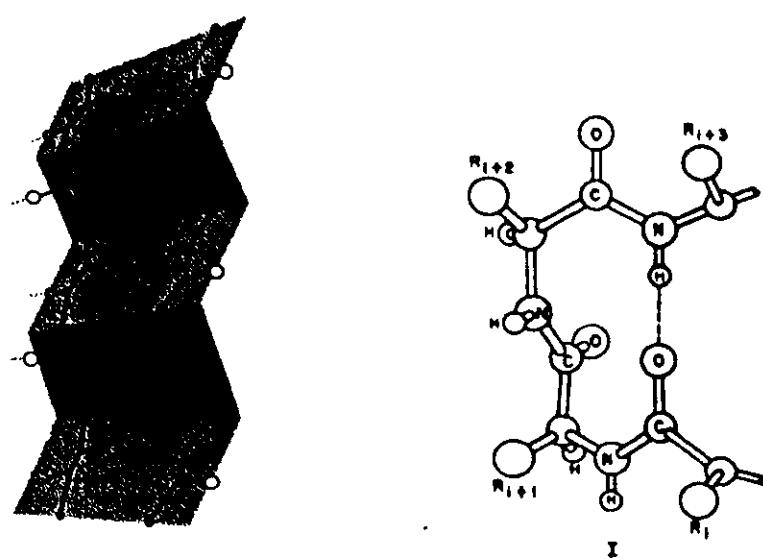
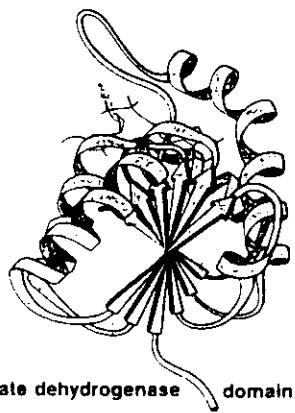


Figure \*5: The structure on the left is an anti-parallel  $\beta$ -sheet. The amide plans have been shaded in. The figure on the right is a reverse turn.

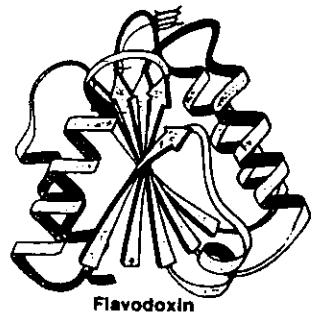
**PROOFS**



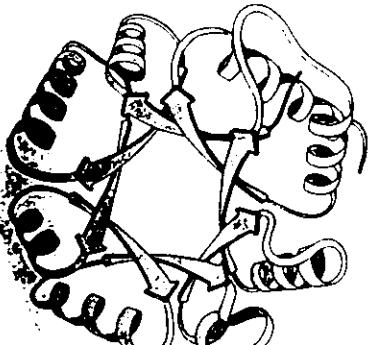
Triose phosphate isomerase



Lactate dehydrogenase domain 1

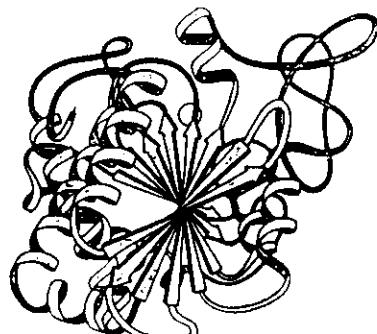


Flavodoxin

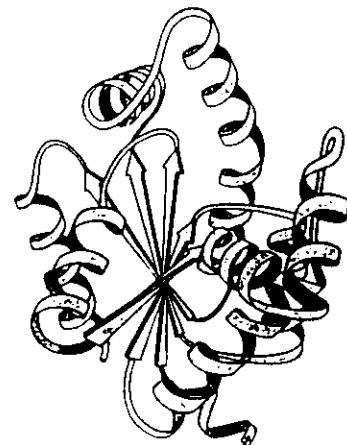


Pyruvate kinase domain 1

(a)

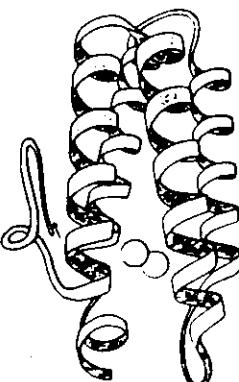


Carboxypeptidase

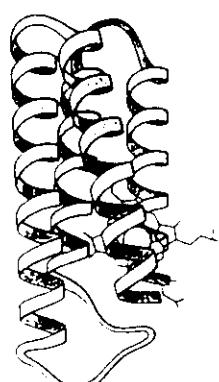


Adenylate kinase

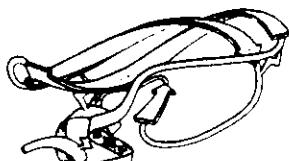
(b)



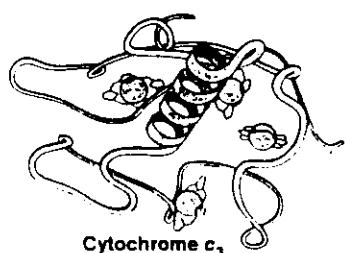
Myohemerythrin



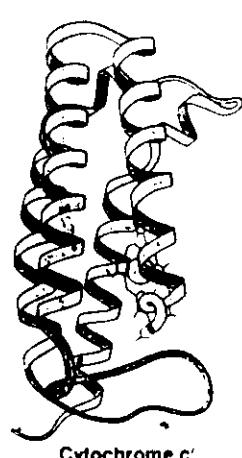
Cytochrome b<sub>562</sub>



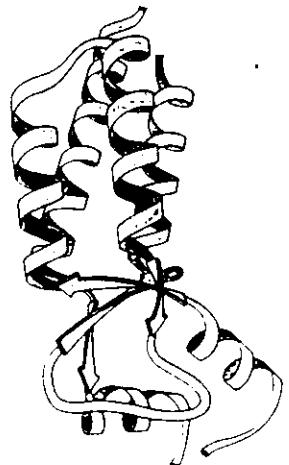
Pancreatic trypsin inhibitor



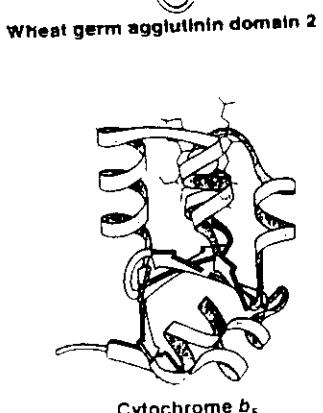
Cytochrome c<sub>3</sub>



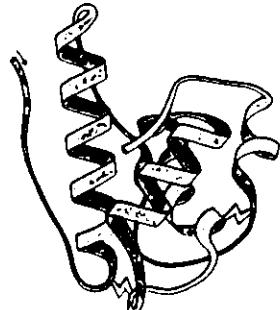
Cytochrome c



Tobacco mosaic virus protein



Cytochrome b<sub>5</sub>



Papain domain 1

## Examine proteins from a physical viewpoint

Much of what we know about proteins comes from biochemists. This information tends to be quite detailed and specific.

Let's use ideas from physics and examine various aspects of proteins.

- First use a simplified model for proteins.
- Map real proteins onto this model
- Measuring various physically interesting quantities

Proteins are complex systems: to understand them it is very useful to examine simple protein-like models.

- Remove all but the essential "ingredients"
- What is essential usually depends on what you are interested in studying.
  - Structure → side chain packing
  - Dynamics / Folding
    - Long time case
    - Short time case
  - Function { GT  
Catalysis
- Probably there is NO simplest Universal (good for all problems) protein-like model

For Folding Really do need to simplify

Eg: All atom + solvent simulations

cf: Wiethrich 2 nanoseconds  $\cong$  3 months

supercomputer  
CPU Time

## Perspective

1 hr  $\sim \$1,000$  (1,500,000 L) Very generous  
 $\sim 2,000$  hrs or  $\$2,000,000$  (3. Billion L)

Not a calculation for mortals

But even if you had money still no good  
Proteins take at least ~1 millisecond to fold

$\Rightarrow$  12.5 years at current computer power

Assume that computers double in speed every two years, it will be 20 years until we can fold all atom simulations on Month time scales.

# Protein Folding: A General View

Nicola Soccia

penultimate lecture → { Tomorrow Will  
NOT Be about  
Structure prediction }

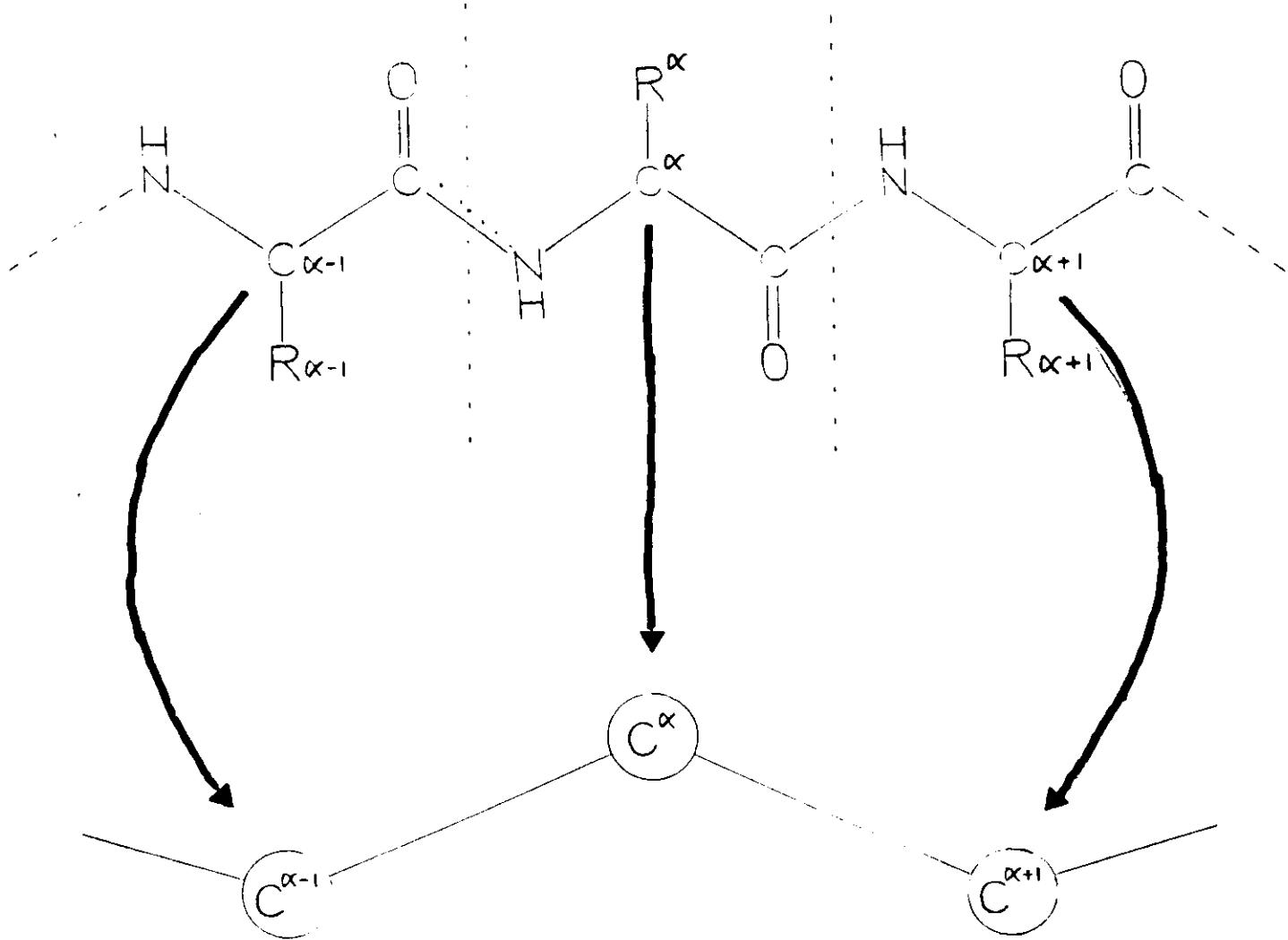
---

Look at proteins again, but with a different perspective (new eyes)

- Put our glasses on (better metaphor, take them off if you are myopic (nearsighted))
- Coarse Grain: qualitatively (generalize)
- Fancy way of saying we are going to ignore most of the atomic level detail.
- However; remember not to throw away too much
  - STUMP THE BAND

## Simple models for proteins.

- Reduced description of proteins.



First Step: Use the simple model to understand real proteins.

Use the model to simplify actual protein structures and then study various statistical properties.

Objective: Look for universal features shared by all\* proteins, rather than specific features of individual proteins

- i.e. What are average (typical) protein properties
- Specific Aim: Define 2° structure in a simple generic way

\* OLD CAVAT: "all" means globular, water soluble, crystallizable proteins since we will be using the PDB

eliminate  
redundant  
examples

Take the Protein Data Bank of roughly 110 proteins and measure the following quantities average over this ensemble:

1. Radius of gyration. Scaling behavior.
2. Pair correlation function. Potentials
3. Angle distribution functions
4. Angle correlation functions. (Secondary Structure)

Goal:

- Understand the universal features of proteins

- Future Work - Need to define Hamiltonian<sup>9</sup> (potential): Only have half of the model. By looking at real protein determine the potential "parameters"

## Radius of Gyration & End to End Length

$$R_{\text{Gyr}} = \sqrt{\frac{\sum_{i=1}^N |\vec{r}_i - \vec{r}_{\text{com}}|^2}{N}}, \quad \vec{r}_{\text{com}} = \frac{\sum_{i=1}^N \vec{r}_i}{N}$$

For sphere   $R_G = R$ ,  Uniform  $R_G = \sqrt{\frac{3}{5}} R$

End to End plot number monomers along chain vs separation in space



Do it for all subsegments of the chain

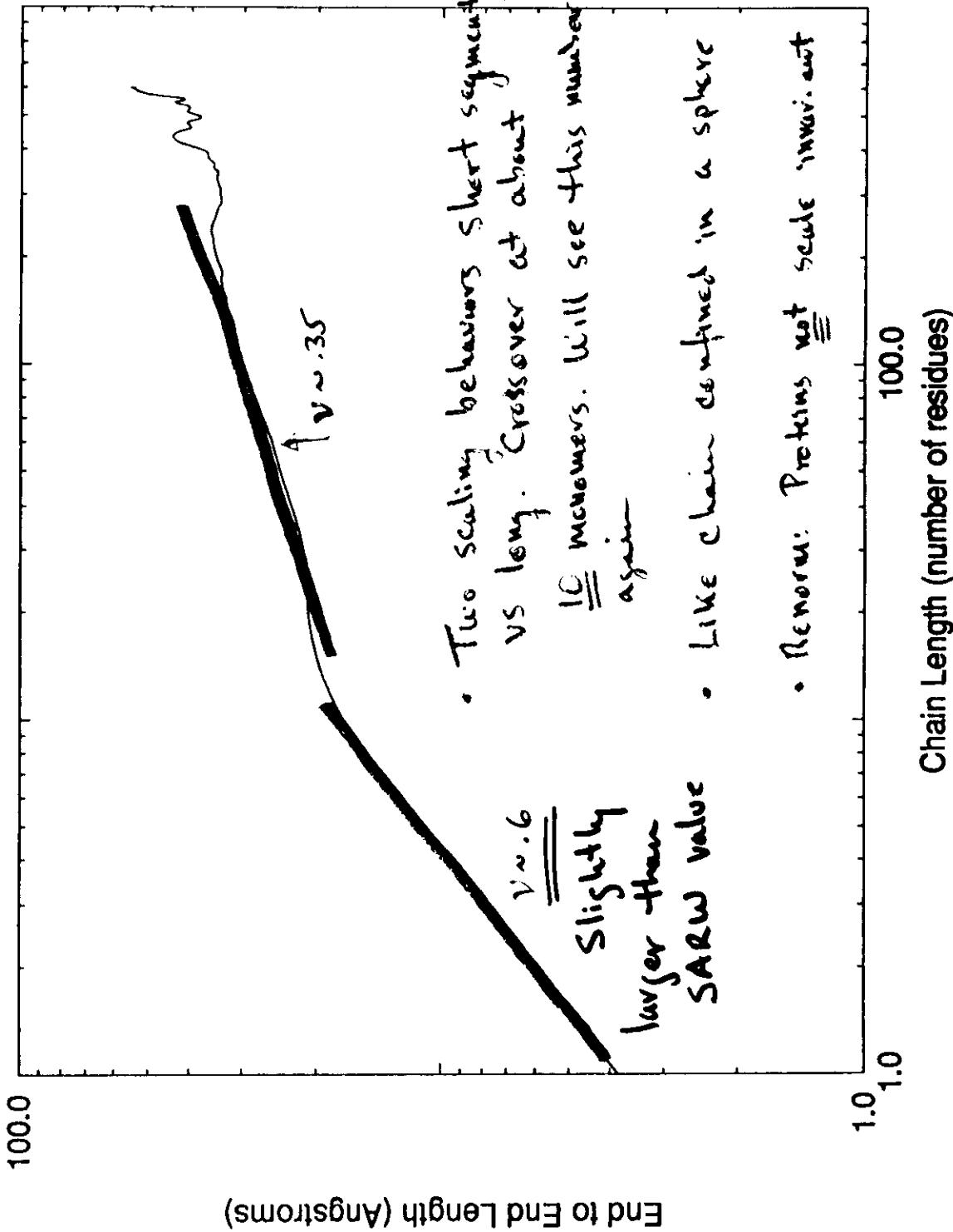
ie distance between all 2 monomer



3  
4  
⋮  
 $N$

## ChainLength vs EndtoEnd Dist

Real Proteins



# Pair Correlation Function (radial): (Radial Distribution Function)

Gives Probability two monomers separated by a distance  $R$

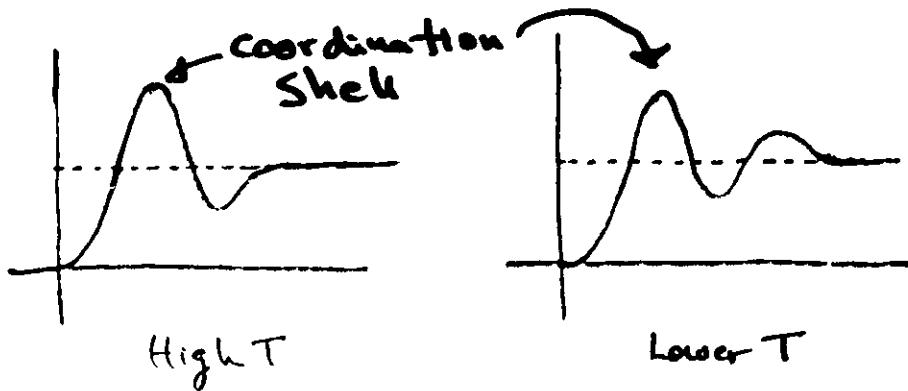
$$\rho(R) \dots R$$

Note Normalized by spherical shell  $\frac{1}{4\pi R^2}$  Dr. Bohr showed unnormalized ones

For Crystal sharp peaks



For a Liquid

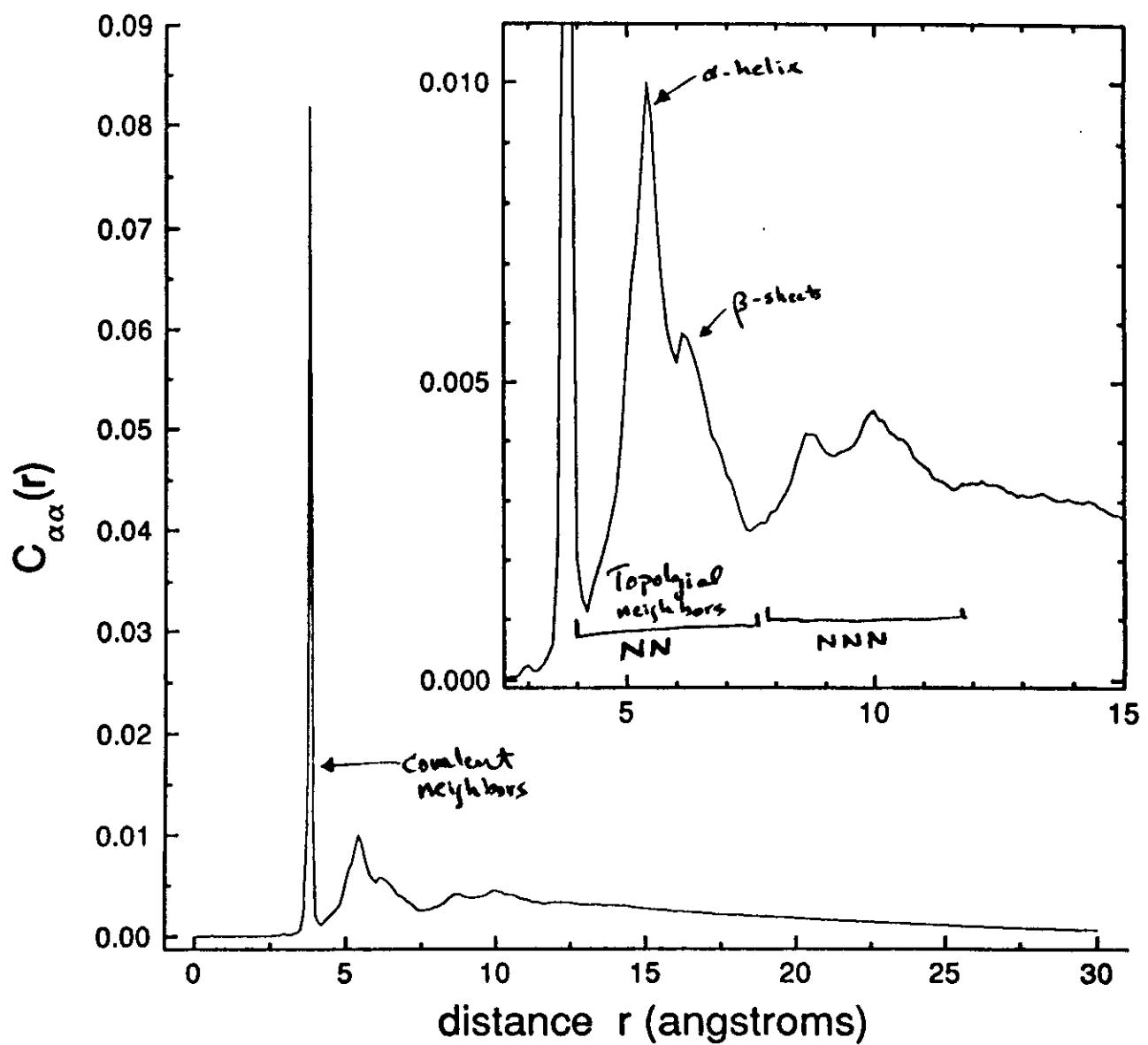


Can Also measure mixed correlation function

i.e  $\rho_{HP}(R)$

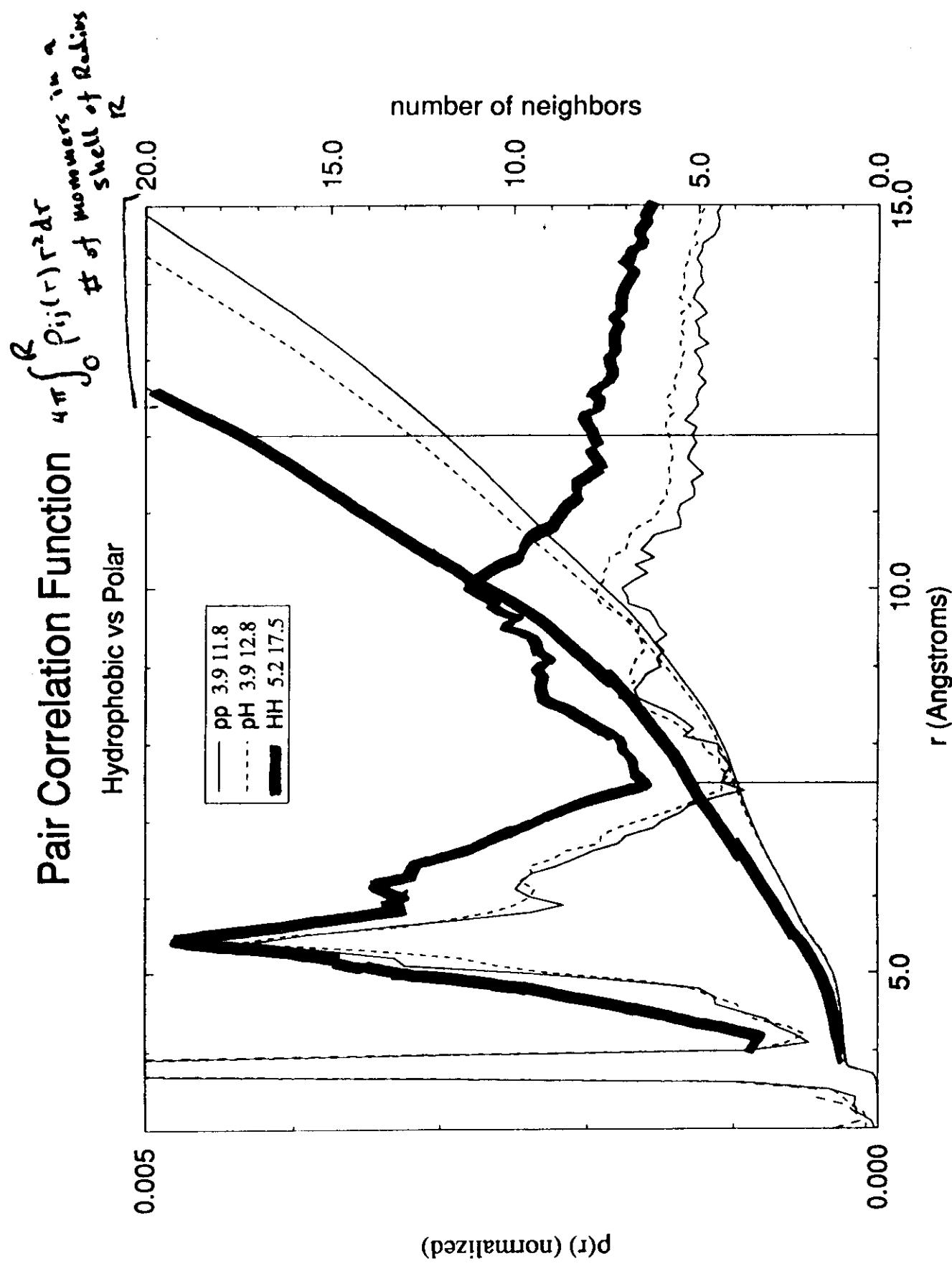
- Polar
- Hydrophobic
- PP
- HH

What is the prob that a Polar monomer is  $R$  from a hydrophobic one



- ① Two nearest neighbor length scales (Lattice problem (simple))
- ② Well defined coordination shell
- ③ Bifurcation of peaks  $\Leftrightarrow \text{Z}^0$  structure
- ④ Potentials BBGKY on cheap  $\log$

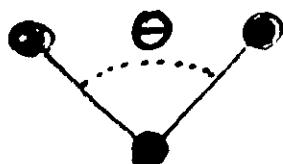
Real Proteins



# Angle Distribution Function.

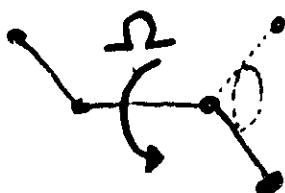
- Simple model also has ~2 angle degrees of freedom per residue just like proteins. But slightly different.
- Examine the distribution of angles to study the local structure of protein chains

① Bond angles: in proteins bond angles fixed here it becomes free



- Excluded Volume
- Bifurcation: 2 peaks ( $\alpha, \beta$  + very) but very different in size and shape  
 $\alpha$ -helix - better defined than  $\beta$ -sheets

② Torsion (Dihedral Angle): Not new used back in 1972!  
to define 2° struct.

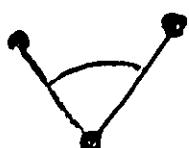
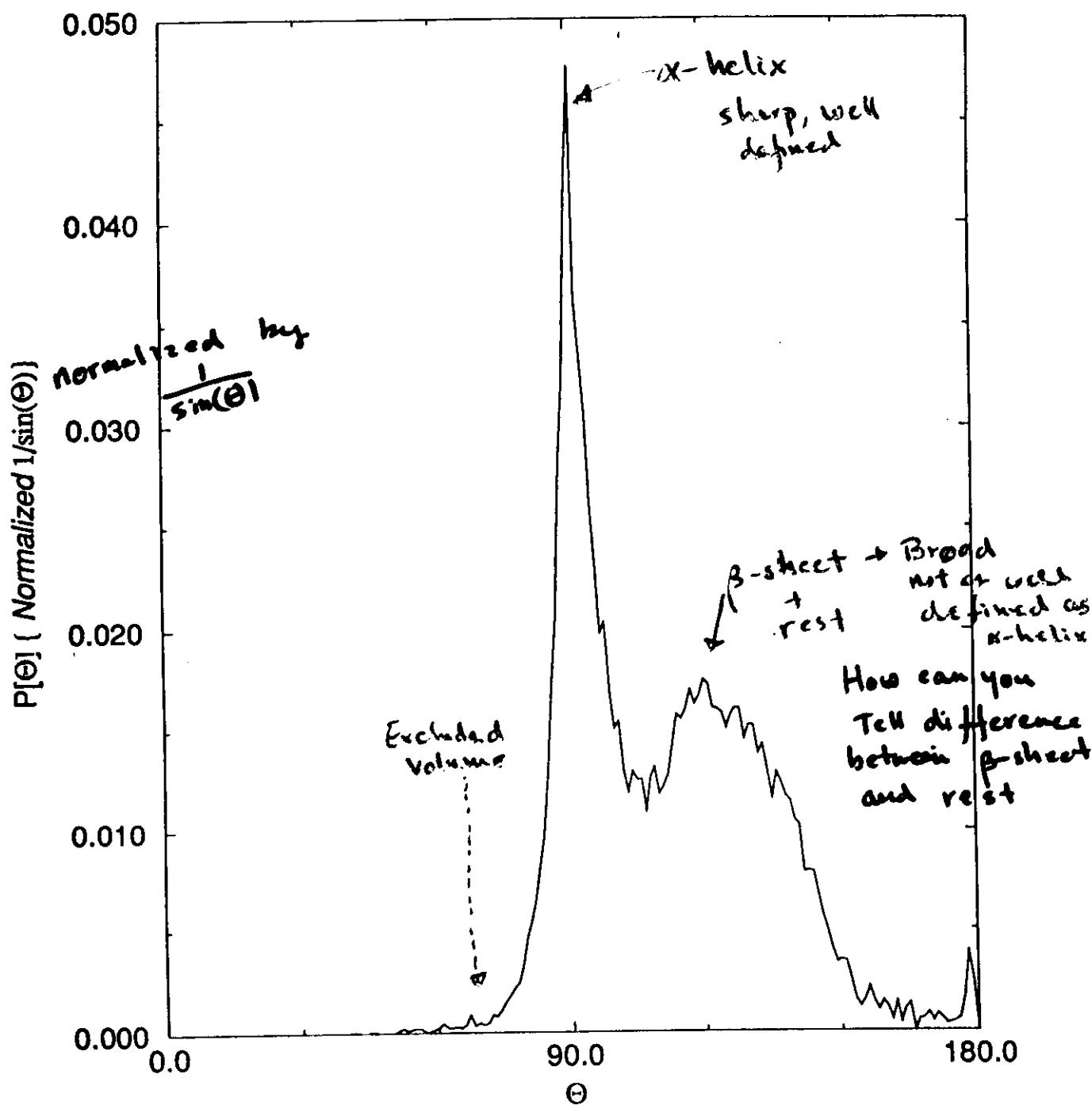


- Bifurcation again. Peaks different
- Chiral symmetry broken ( $\frac{1}{2} \rightarrow \frac{1}{2} \pm \pi$ )

③ Quasi Ramachandran Plot

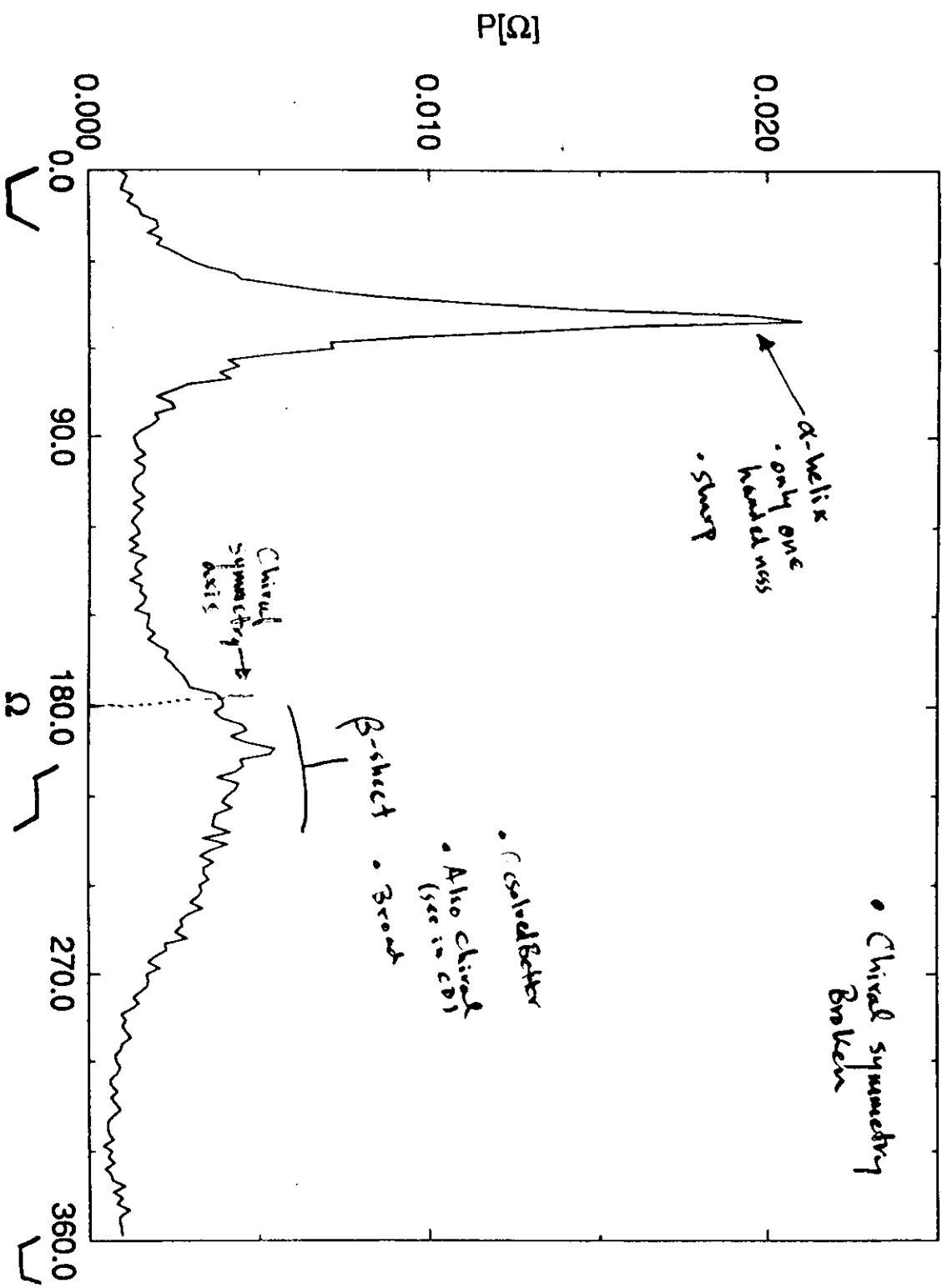
# Planar Angle Distribution Plot

Real Proteins



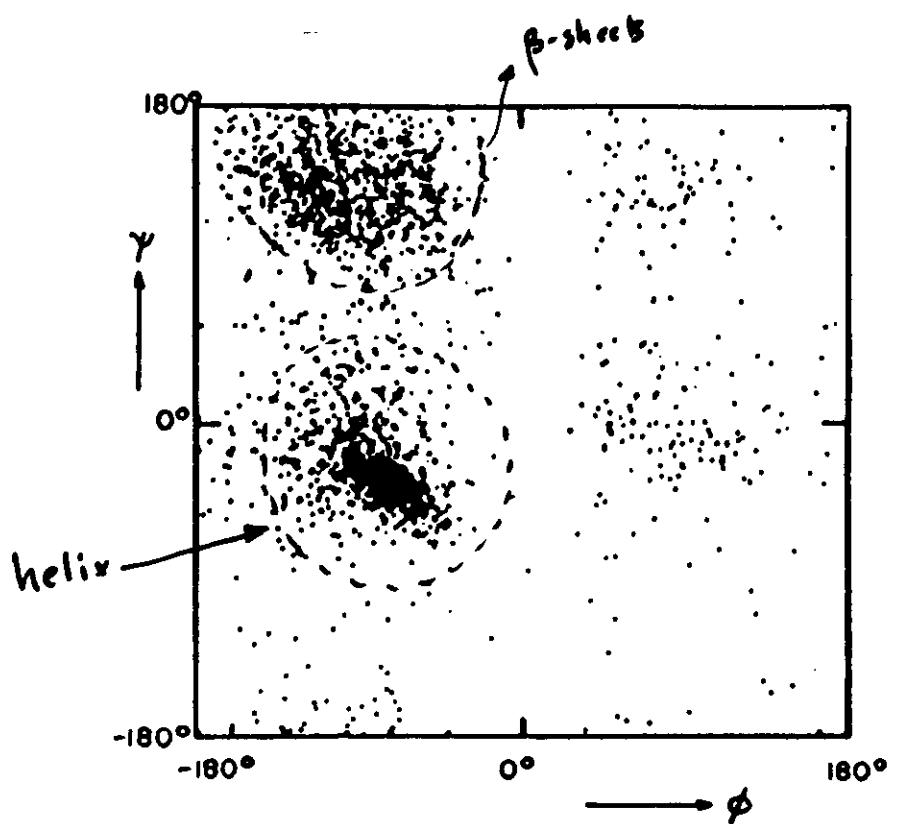
# Torsion Angle Distribution

Real Proteins



25

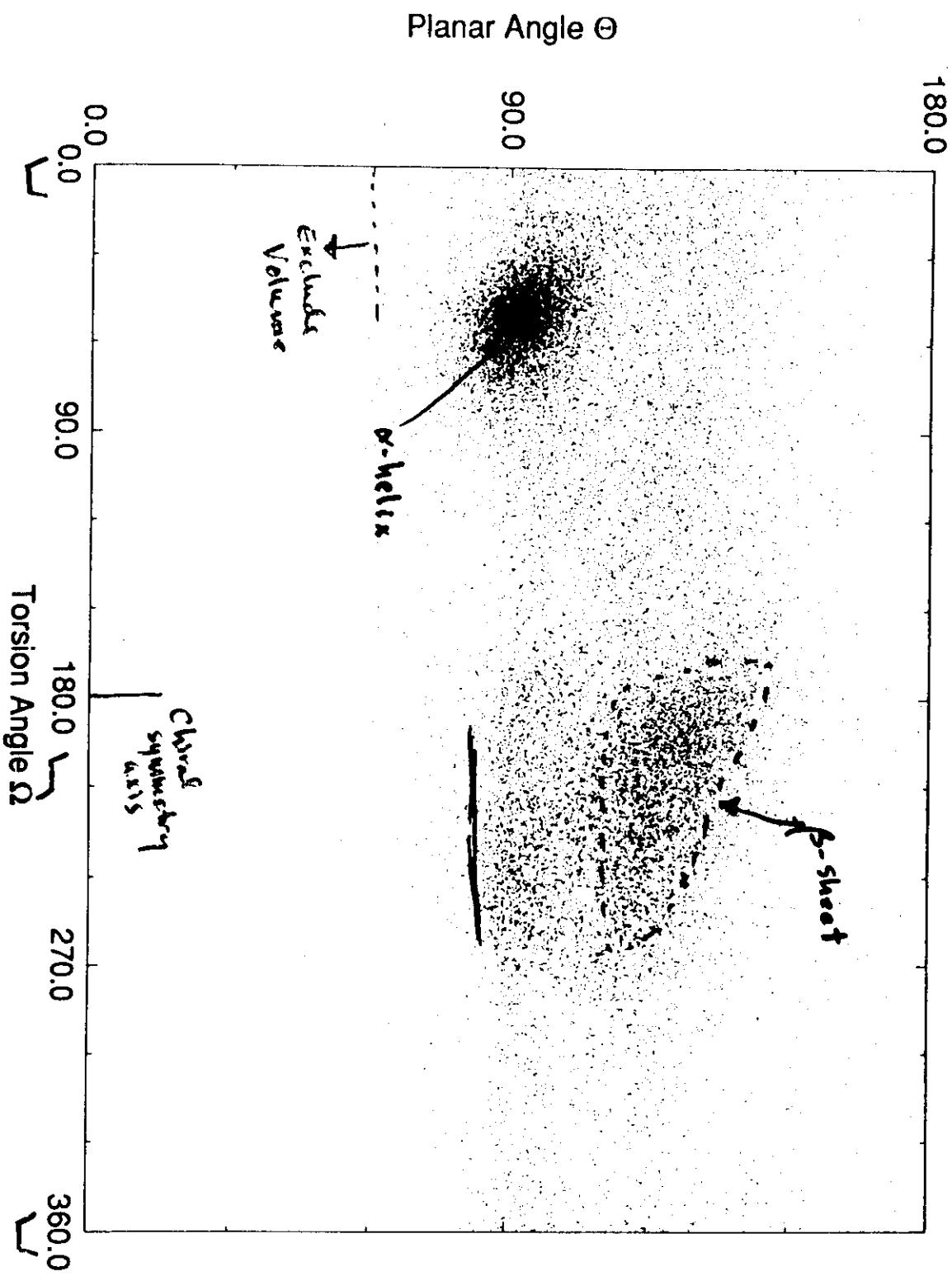
## Ramachandran Map



- Note usually map is made by looking at  
steric interactions in di, tri peptides
- This density was taken from the PDB

# Torsion/Planar Rham Plot

Real Proteins



# Chain Correlation Functions

correlation =  
• Green Function  
• Kernels  
• Propagator

previous corr. functions were through space

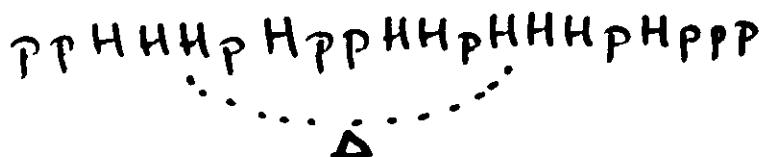


now let us look at correlations along the chain



$C(\Delta)$  where  $\Delta$  is the separation  
(in number of residues) along  
the chain

- First Look at Sequence Correlations: just polar vs Hydrophobic

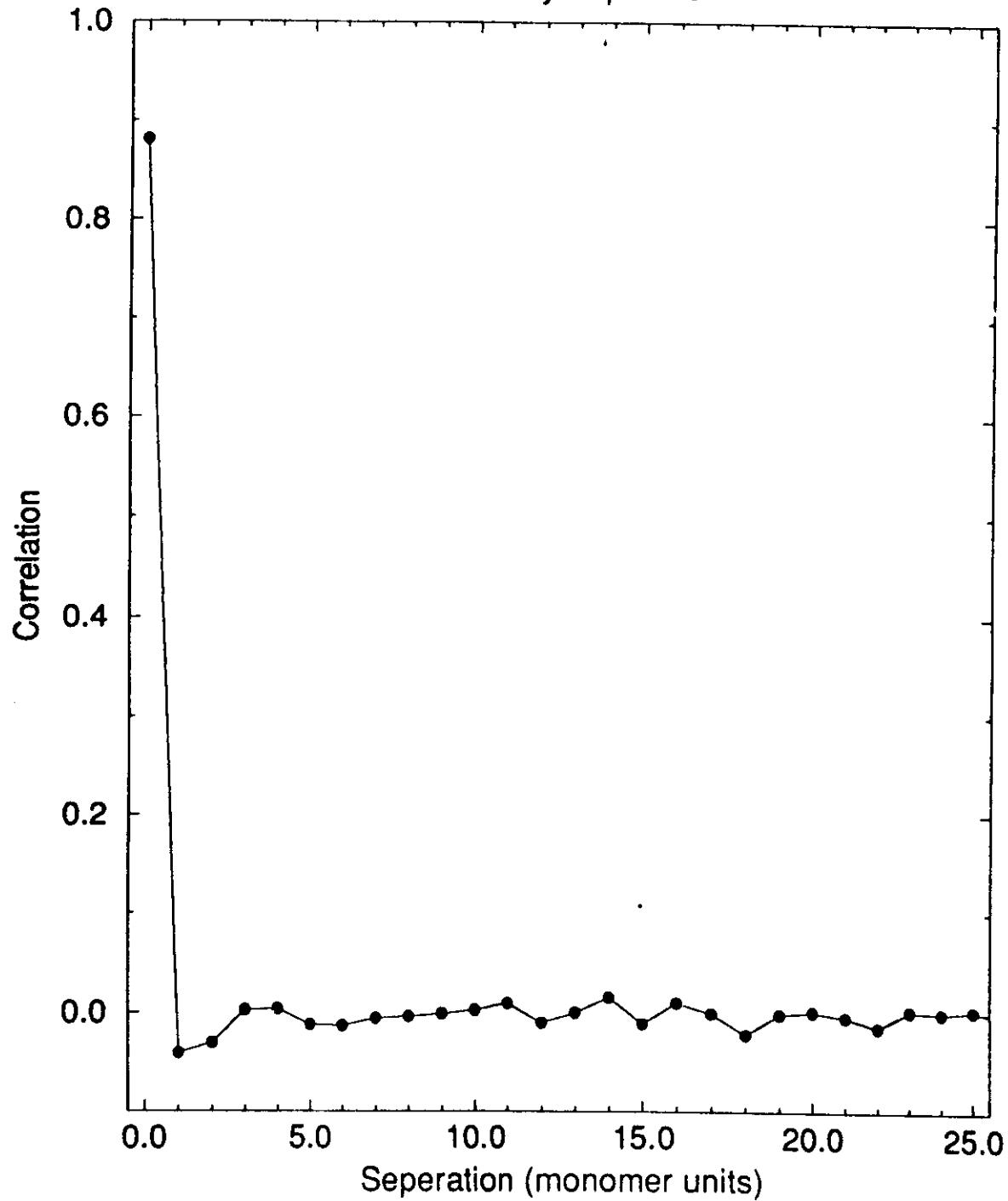


- No correlations!
- Same is true if one measures corr of the 20 letters AVTGIN....
- This is one difficulty in predicting 2° structure from local sequence (Identical pentapeptides)

N.B. Random Sequences of AA not proteins, yet  
Protein Sequences "appear" random!??

# Chain Correlation

Polar-Hydrophobic

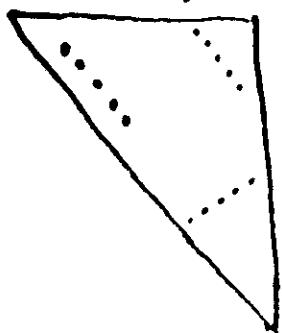


## Chain Correlations Cont.

What is secondary structure?

Biology Definition: Keys in on H-Bonds

- Define 2° structure (recognize) by look for H-bonding patterns. (Kabsch & Sander)



$i \dots i+4, i+1 \dots i+4+1$   $\alpha$ -helix  
 $i \dots i+j, i+1 \dots i+j+1$  parallel  $\beta$ -sheet  
 $i \dots i+j, i+1 \dots i+j-1$  anti-parallel  $\beta$   
 $i \dots i+3 (i+4)$  turn

- Not quite so simple: what is H-bond to what not also easy to determine
- Nice definition but can be too specific, nor is it easy to calculate,
- Bigger Problem: our simple model does not have Hydrogens (no H-bonds)

Can we define secondary structure in our simple models??

Yes, we can

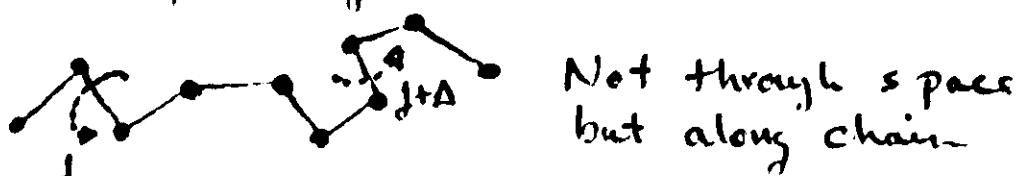
ReRead Bio Def: This time key in on  
regular.

Regular Arrangement of chain, i.e Correlated along the chain.

But what is correlated?

- Not Sequence
- Angles

Measure Angle - Angle correlation Function



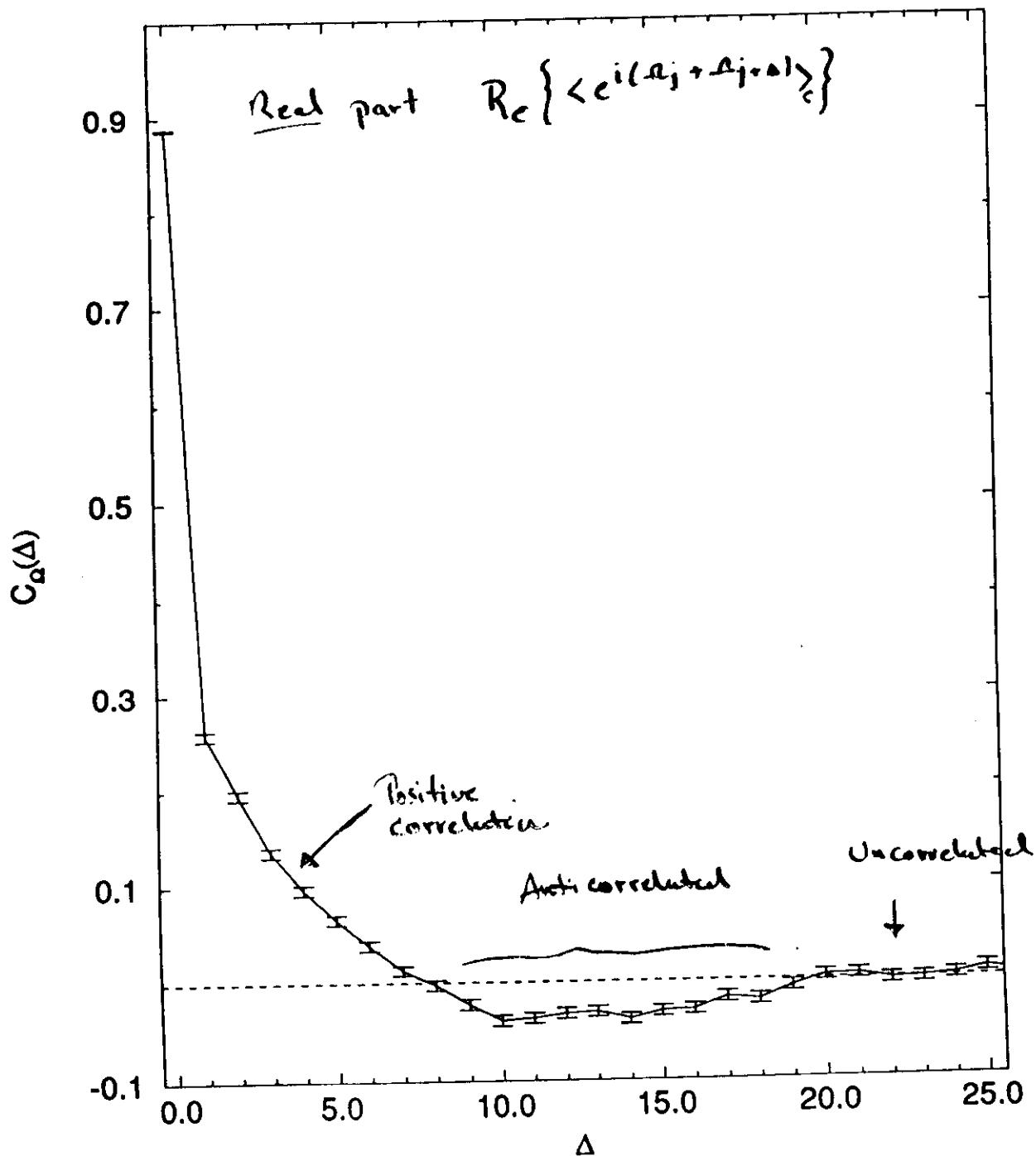
$$\langle e^{i(\Omega_j - \Omega_{j+\Delta})} \rangle$$

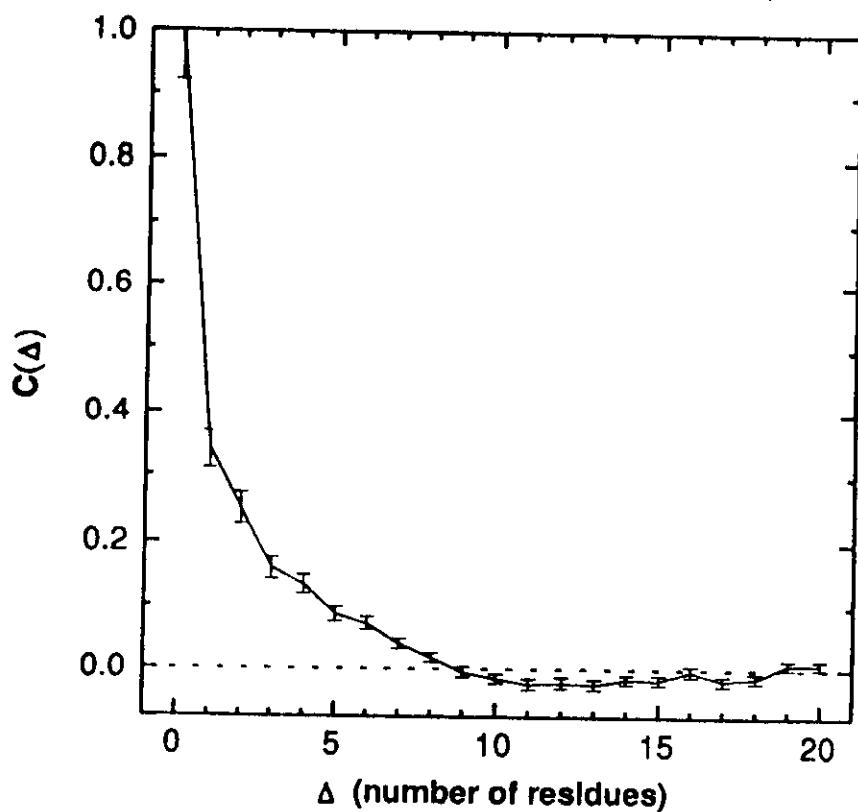
C ← connected correlation  
(i.e subtract mean)  
 $|\langle e^{i\Omega_j} \rangle|^2$

- Long correlations: up to 10 AA  
10 is the "average" size of  $\text{Z}^\circ$  in proteins
- Anti correlated piece: super secondary struct

## Torsion Angle Correlation

Real Proteins





Real ( $\theta$ )

## Simple models continued:

- Simplified Potentials

- Dominant force in protein folding: compactification of structure due to hydrophobic effect.

$$\mathcal{H} = \sum_i \frac{1}{2} k_{\text{cov}} (|\vec{r}_i - \vec{r}_{i+1}| - l_{\text{cov}})^2 + \frac{1}{2} \sum_{ij} V(r_{ij})$$

*Covalent Term*                                   *Every thing else*  
 $l_{\text{cov}} = 3.78 \text{ \AA}$                                     $r_{ij} = |\vec{r}_i - \vec{r}_j|$   
 $k_{\text{cov}} = 1$

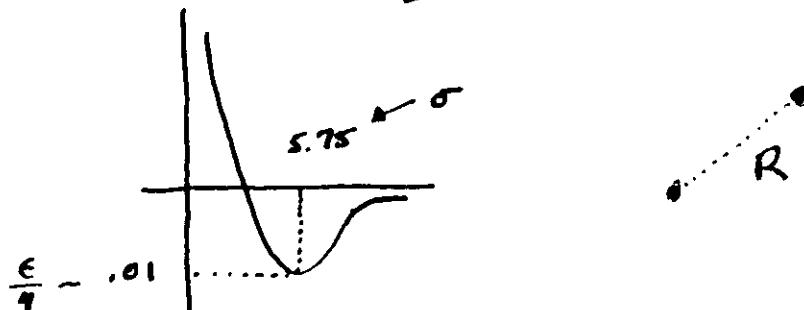
- Pick  $V(r_{ij})$  to compact protein.
- Start with simple two body interactions.

Goal: Find a protein-like potential for  
later study in folding and dynamics

### 3 potentials

#### ① LJ

$$V(R) = \epsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right]$$



#### ② Rad Gyr

$$V(\{\vec{r}_i\}) = \epsilon \left\{ \sum_{i < j} \left( \frac{\sigma}{r_{ij}} \right)^{12} + \sum_{i=1}^N |\vec{r}_i - \vec{r}_{com}|^2 \right\}$$

$\epsilon = .01$     $\sigma$  pick by look at  $p_{rad}(R)$

#### ③ Rad Gyr + Chiral Local (chain) Int



$$V_C = \hat{s}_{i-1} \cdot (\hat{s}_i \times \hat{s}_{i+1})$$

"Fold" (compact) the Chain.

Minimize the Energy

Use a simple function minimization routine  
to minimize the potential energy.

Possible Problem: Model is a homopolymer  
will have many ( $c^N$ ) low energy structure  
that look very different.

Make the follow conjecture

Each one of the many different low energy structs of  
the homopolymer will correspond to the Native  
structure of a different seq in a related  
heteropolymer model

So by generating a set of compact structures we  
can then average over sequence, and repeat the analysis  
we did for real proteins

## Score Card.

How protein-like are the potentials

- No potential 100% protein-like  
Each seems to capture different structural features of real proteins
- No Secondary Structure !

Need Better Potentials: Possibilities

- ① Angle Potentials (Explicit)
- ② Combine the two (LJ - Radgyr)
- ③ Heteropolymer. Maybe Conjecture is incorrect.
- ④ Model too simple :
  - Side chains
  - ....

## More Questions (Physics)

- Hetero vs Homopolymer
  - ↓
  - # of gs small
  - O(1)
- ↓
  - Many gs  $\sim N^{\gamma}$

How many monomers do you need to get unique  
gs? RNA-4    Experiments  $\approx$  (HP)  
Lattice Models 2 -

- Real Seq vs Random -  $2^{0^N}$  seq what % folds?
  - Qualitative Ans ??
- Minimal Information to Solve folding Problem
  - Complexity NP?? (Lattice Proof)  
Problems? ↗
- Inverse Folding Problem
  - may be easier
  - possible more useful commercially

## Conclusions (Philosophical)

- Proteins are an extremely complex system, with a wide diversity in structure and function.
- On close examination, they appear to share common (universal) features. Theory(ies) of proteins seems possible.
- Simple systems can sometimes have very complex and interesting properties  
Maybe much of the apparent complexity of proteins can be understood from looking at simpler model systems
- We study proteins to study how living organisms work (to study life):  
But understanding life by looking at protein function may be as difficult as trying to understand proteins from e=photon interactions

