

Thousand Faces of Proteins

2004 1 13

1. (protein) 가?
2. (protein synthesis) Central Dogma
3. Protein folding problem vs. Traveling salesman problem
4. Protein structure (prediction)

,

•

,

,

,

•

,

,

•

,

•

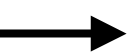
,

•

•

.....

Protein Folding Problem



1.

(protein)

가?

- DNA, RNA, (carbohydrate), (lipid)

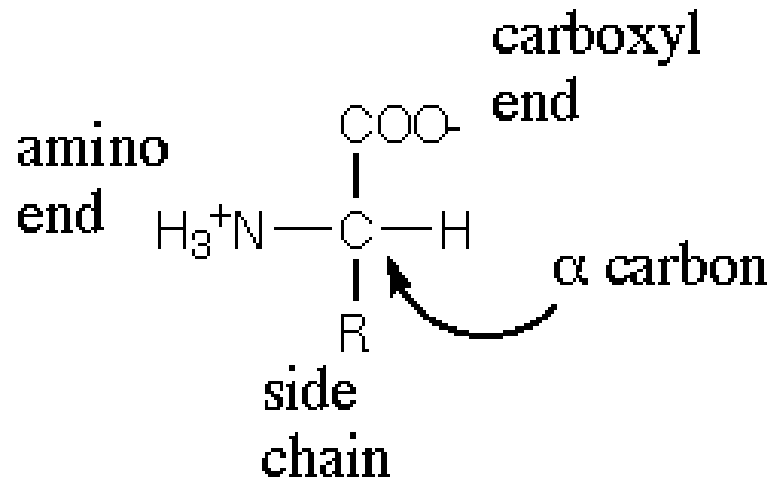
Introduction

- ?

cys-gly-val-ala-ala-leu-met

- : 20가
- 가

- (Amino acid):



side chain = H , CH_3 , ...

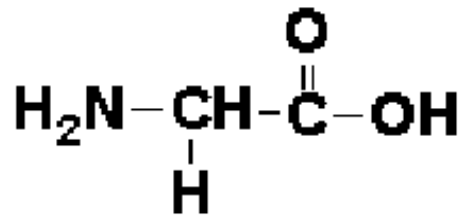
- Side chain

.

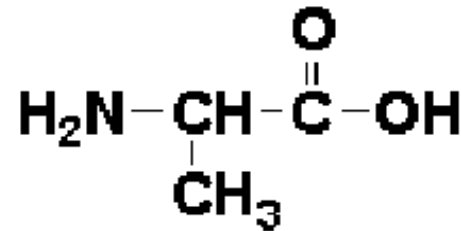
- 가 .

20가

가



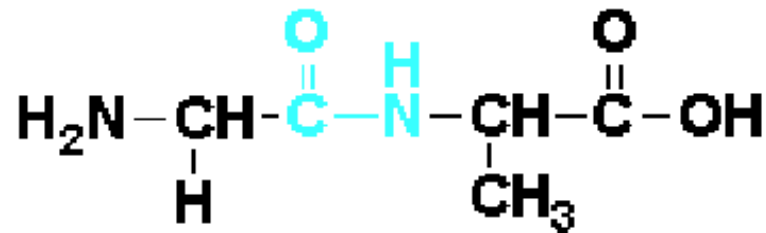
Glycine



Alanine

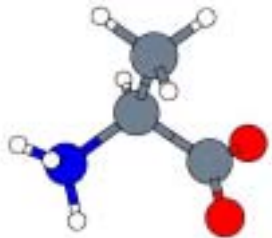


(peptide bond)

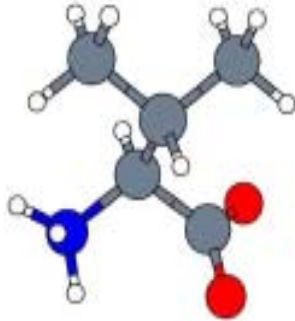


Glycylalanine

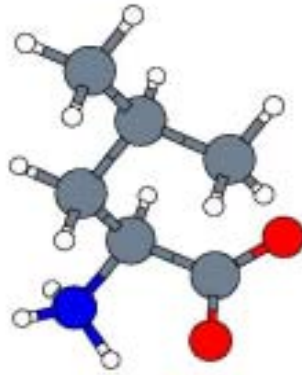
(hydrophobic)



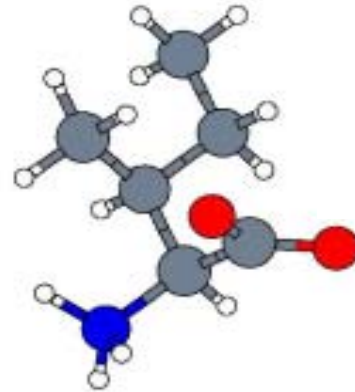
alanine



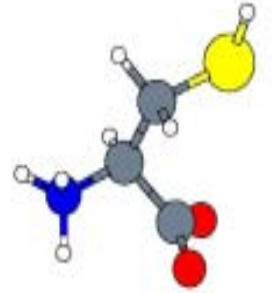
valine



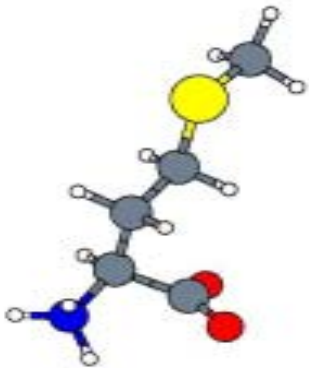
leucine



isoleucine



cysteine



methionine



phenylalanine

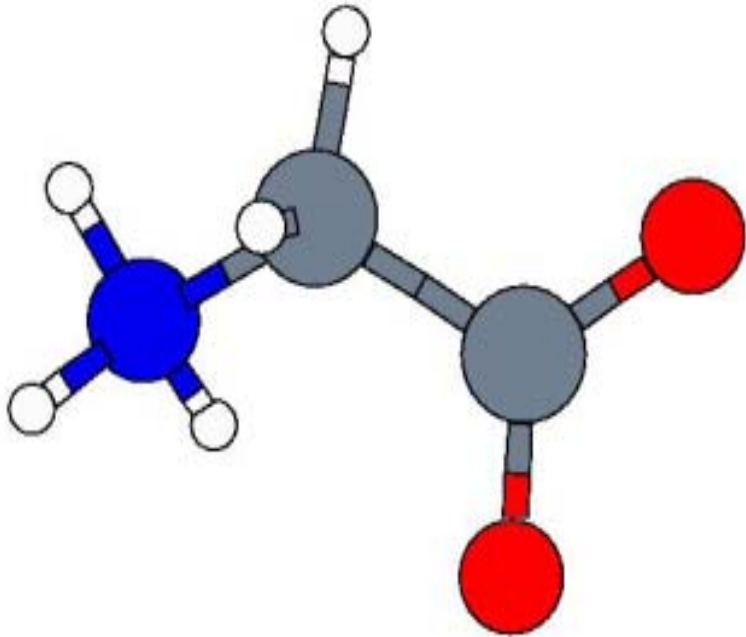


tyrosine

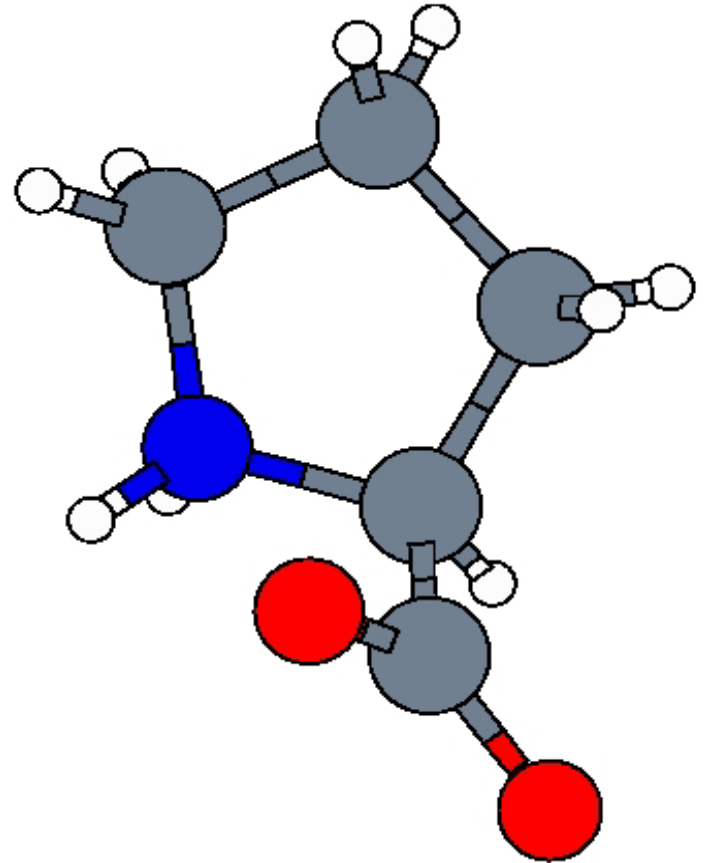


tryptophan

**neither hydrophilic nor hydrophobic
but rather H₂O-soluble**

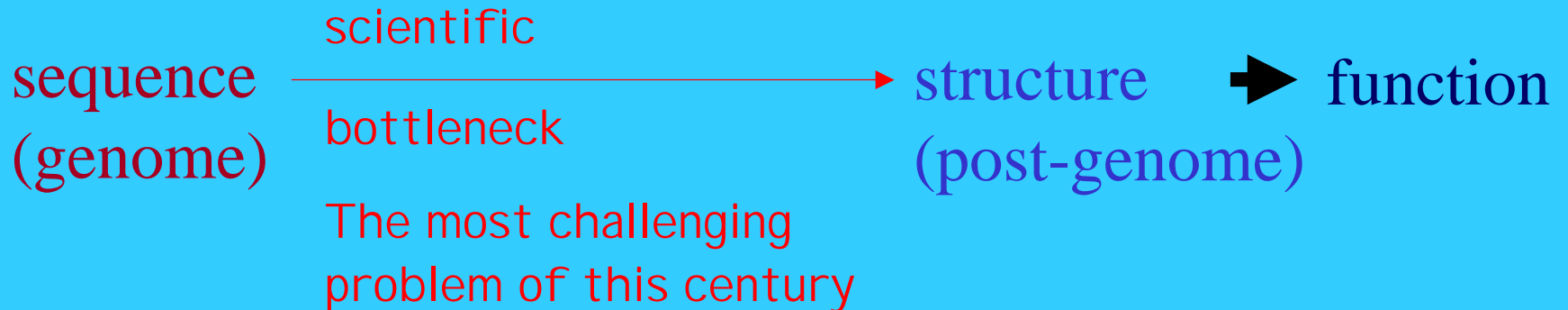
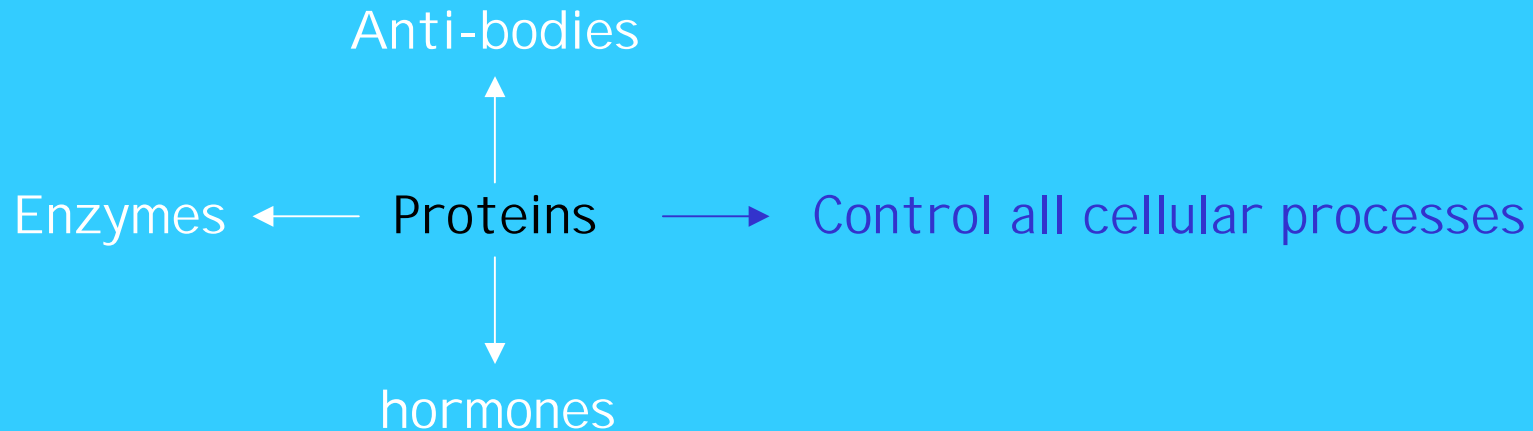


glycine



proline

Why are proteins important ?



(Protein Folding Problem)

Protein Folding Problem

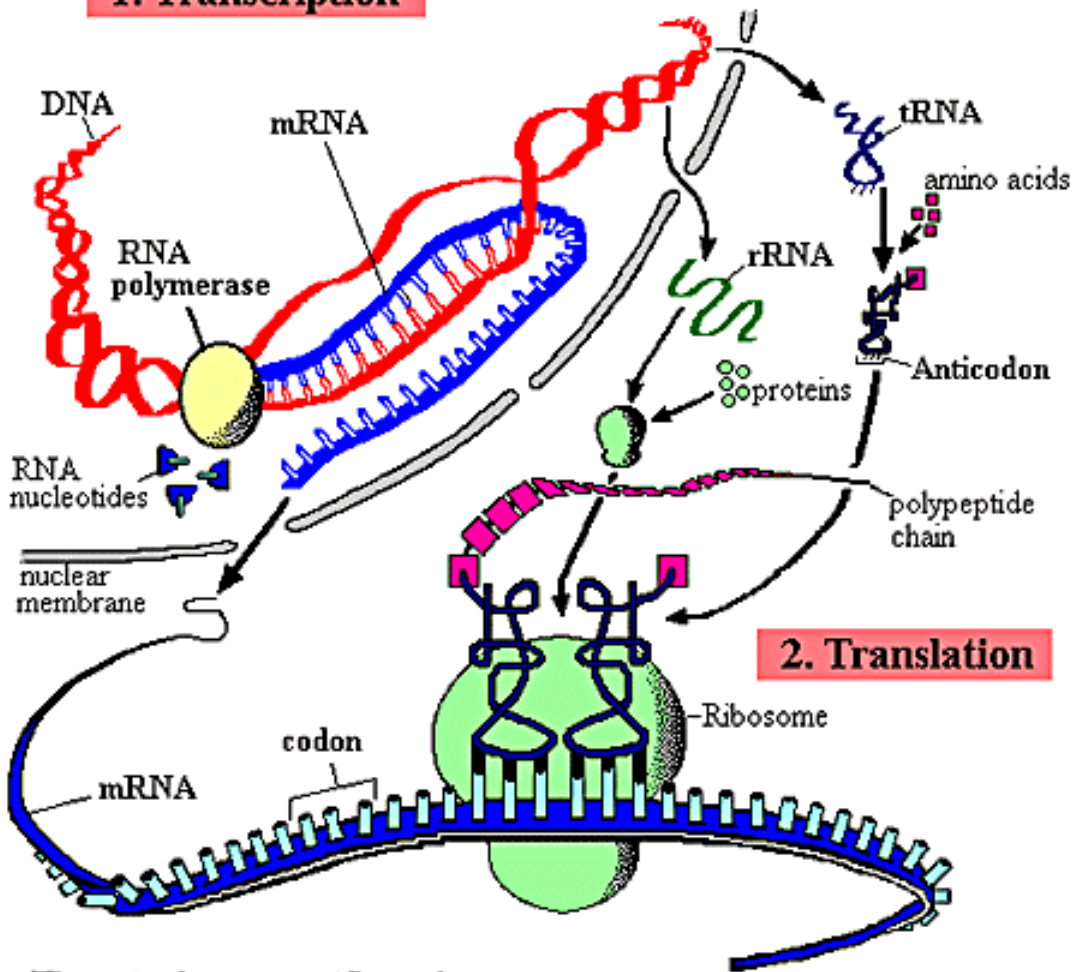
1. (protein) 가?

→ 2. (protein synthesis) Central Dogma

Central Dogma



1. Transcription



Human hemoglobin

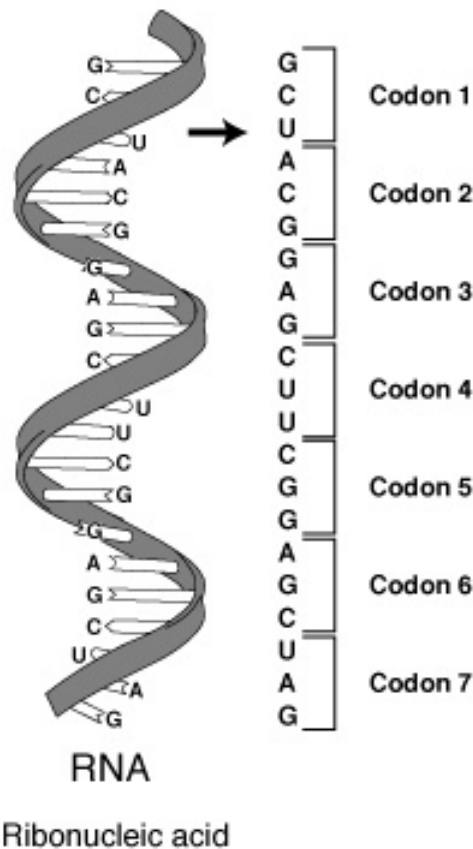


Protein folding problem

- For a given amino acid sequence (of size n), find the native structure of the protein.
- Total # of protein structures: 10^n
- mathematically NOT well defined problem

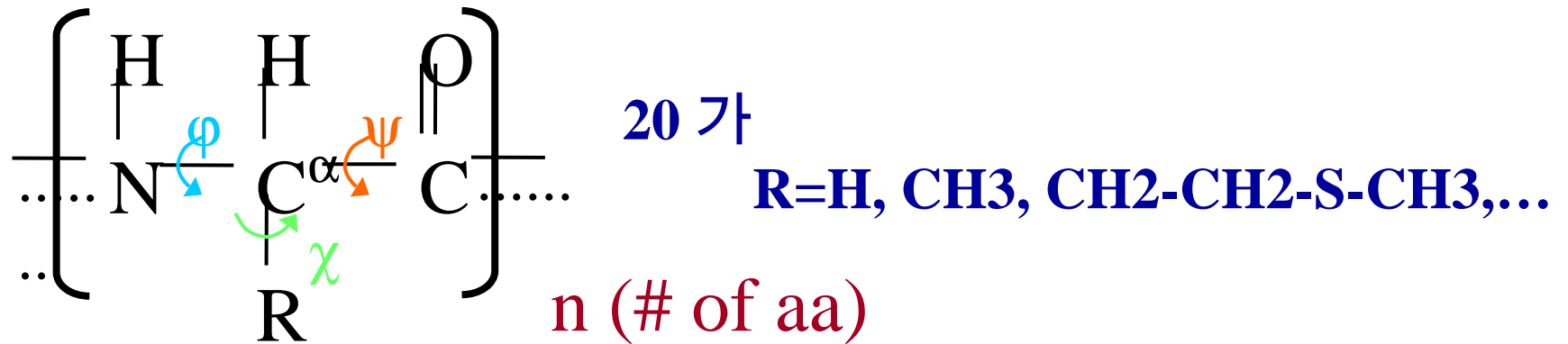
sequence → structure → function

DNA (RNA), Codon and Amino acids



1st	2nd			
	T	C	A	G
T	TTT 0.43 TTC 0.57 Phe TTA 0.06 TTG 0.12 Leu	TCT 0.18 TCC 0.23 Ser TCA 0.15 TCG 0.06	TAT 0.42 TAC 0.58 Tyr TAA 0.22 TERM TAG 0.17	TGT 0.42 TGC 0.58 Cys TGA 0.61 TERM TGG 1.00 Trp
	CTT 0.12 CTC 0.20 Leu CTA 0.07 CTG 0.43	CCT 0.29 CCC 0.33 Pro CCA 0.27 CCG 0.11	CAT 0.41 CAC 0.59 His CAA 0.27 Gln CAG 0.73	CGT 0.09 CGC 0.19 Arg CGA 0.10 CGG 0.19
	ATT 0.35 ATC 0.52 Ile ATA 0.14 ATG 1.00 Met	ACT 0.23 ACC 0.38 Thr ACA 0.27 ACG 0.12	AAT 0.44 AAC 0.56 Asn AAA 0.40 Lys AAG 0.60	AGT 0.14 AGC 0.25 Ser AGA 0.21 Arg AGG 0.22
	GTT 0.17 GTC 0.25 Val GTA 0.10 GTG 0.48	GCT 0.28 GCC 0.40 Ala GCA 0.22 GCG 0.10	GAT 0.44 GAC 0.56 Asp GAA 0.41 Glu GAG 0.59	GGT 0.18 GGC 0.33 Gly GGA 0.26 GGG 0.23
C				
A				
G				

(Amino Acid Residue)



- (polypeptide):

• :

(10⁵)

• 3 •

The Protein Folding Problem

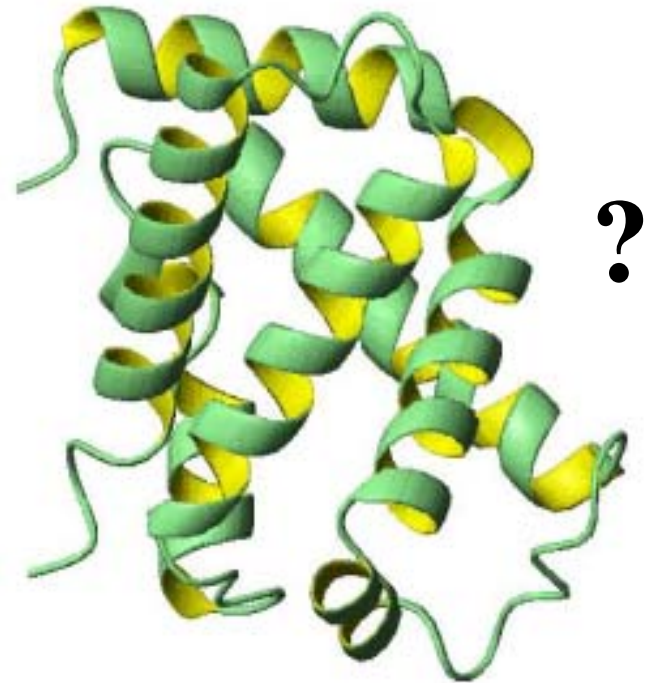
1.

(
가?

: Native structure)

3

141 amino acids
2152 atoms
deoxy human
hemoglobin
(oxygen transport)



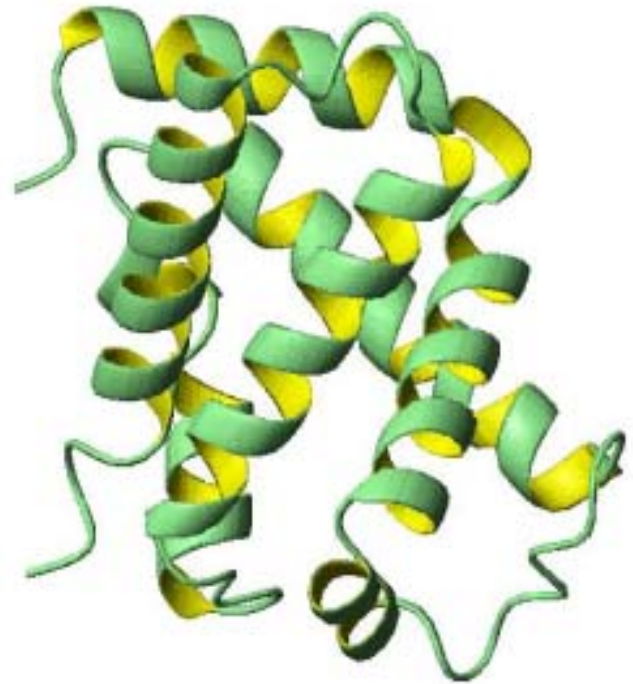
The Protein Folding Problem

2.

(Protein-Folding Mechanisms)

가?

141 amino acids
2152 atoms
deoxy human
hemoglobin
(oxygen transport)



Protein-Folding Mechanisms

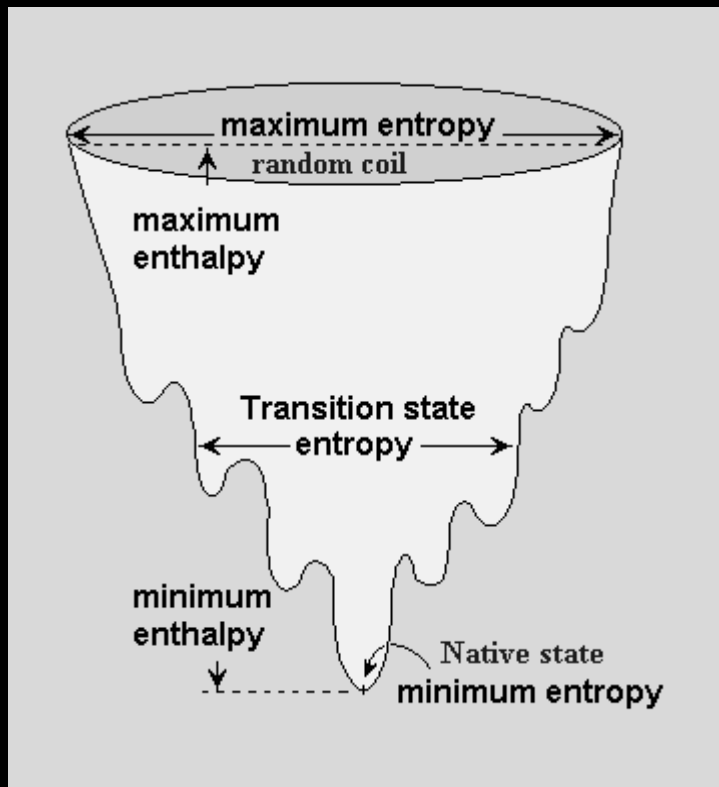
Random search of all conformational space requires an immense amount of time (longer than the age of universe).

In vitro refolding normally takes seconds or minutes.

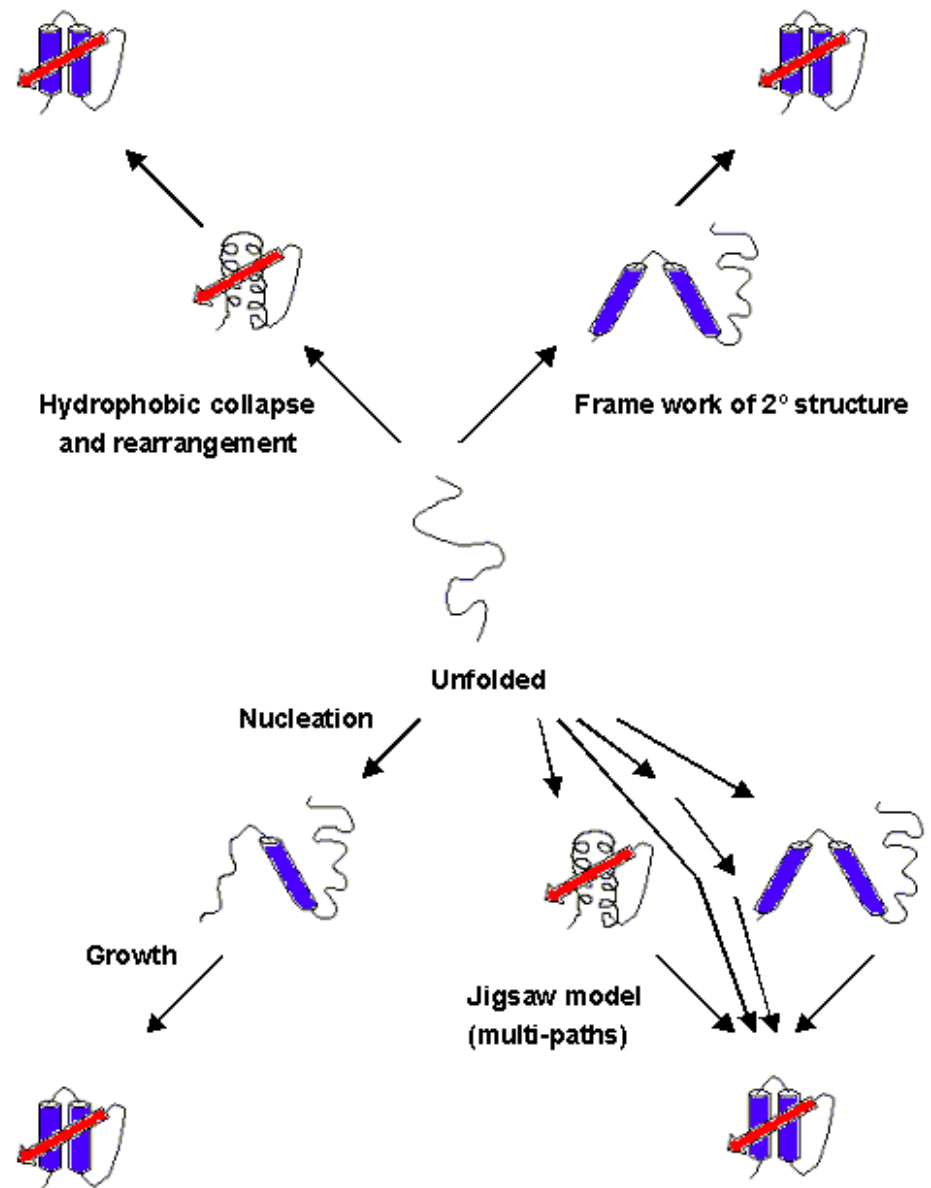
-Levinthal paradox

Ex) for a protein with 100 amino acids, random search, helix, sheet, random search, native !
 3^{100} (or about 10^{48}).
 10^{14} s^{-1} (bond rotation)
 10^{34} , random search 10^{26} → native !

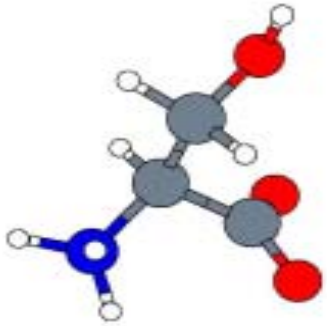
→ **Folding pathway** problem: identifying intermediates and constructing folding mechanism



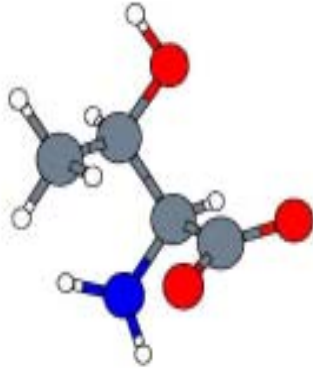
Various models for protein folding



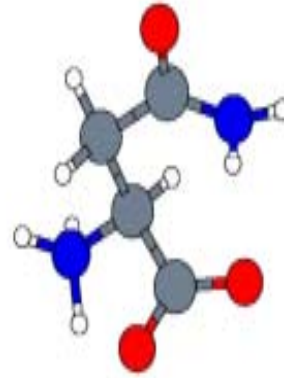
(hydrophilic)



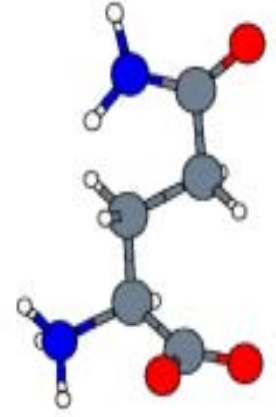
serine



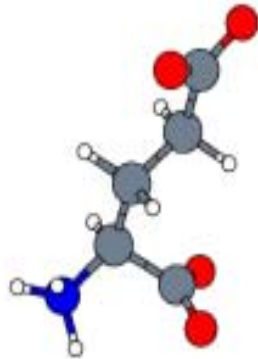
threonine



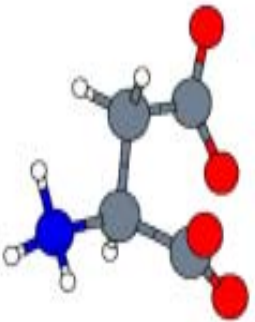
asparagine



glutamine



glutamic acid



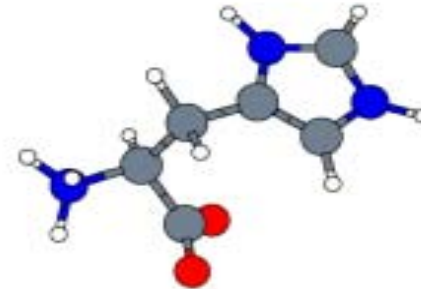
aspartic acid



lysine



arginine +



histidine +

- DNA (sequence)

e.g.) ATT - ACG - CAG - CCA - CGG - CGG - ATT

- RNA

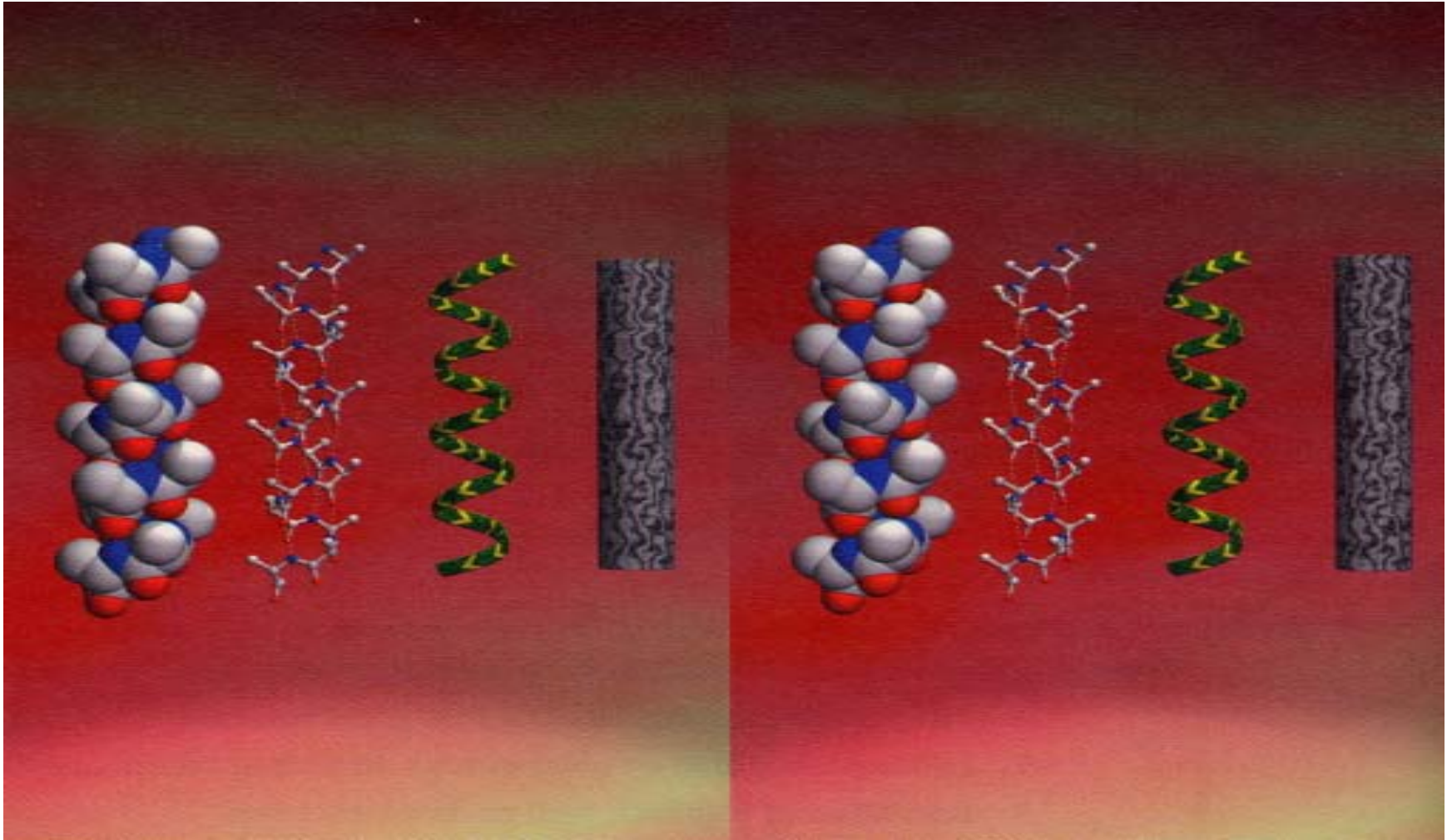
UAA - UGC - GUC - GGU - GCC - GCC - UAA

- (Amino acid sequence)

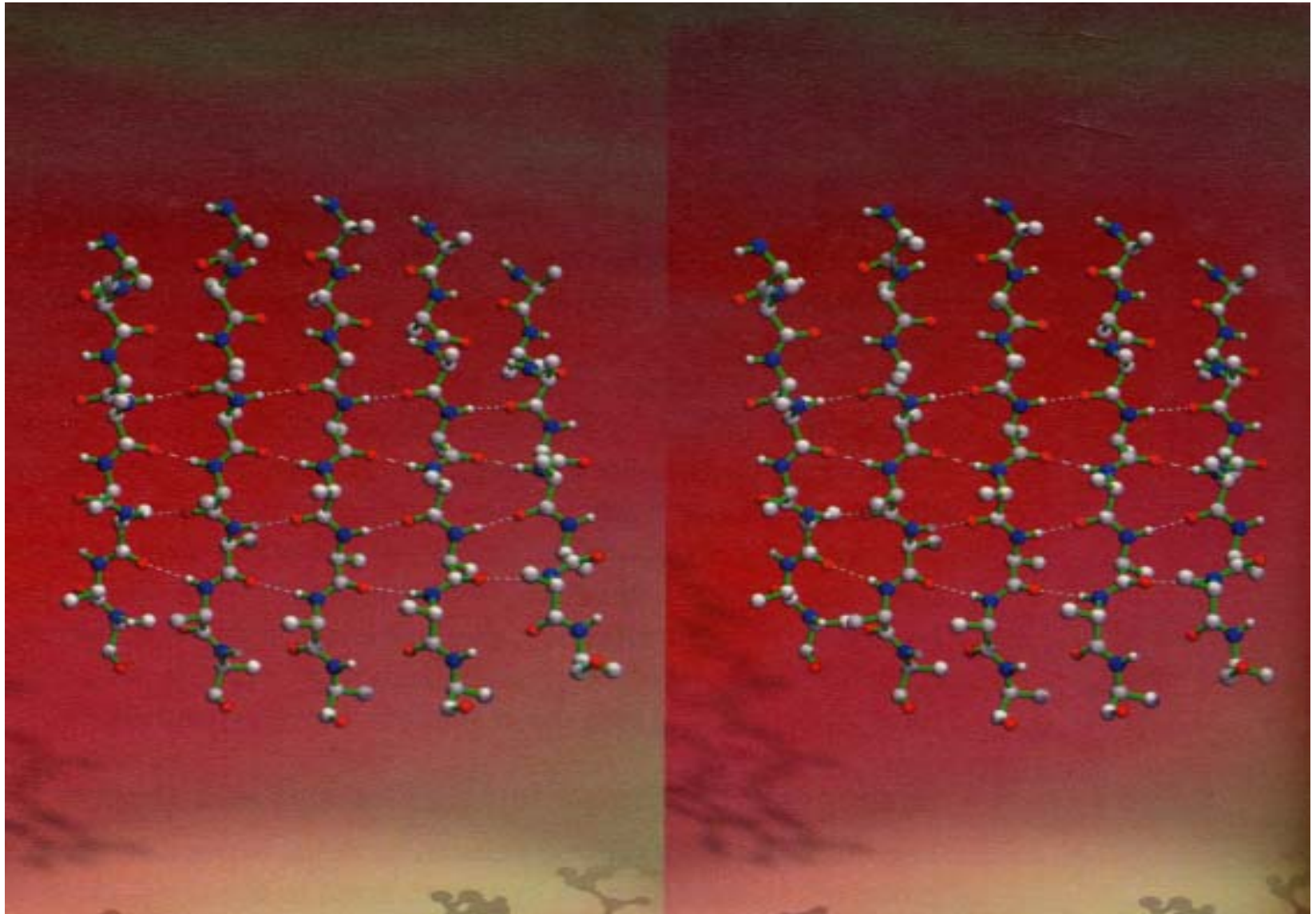
cys - gly - val - ala - ala

→ 1 (Primary Structure)

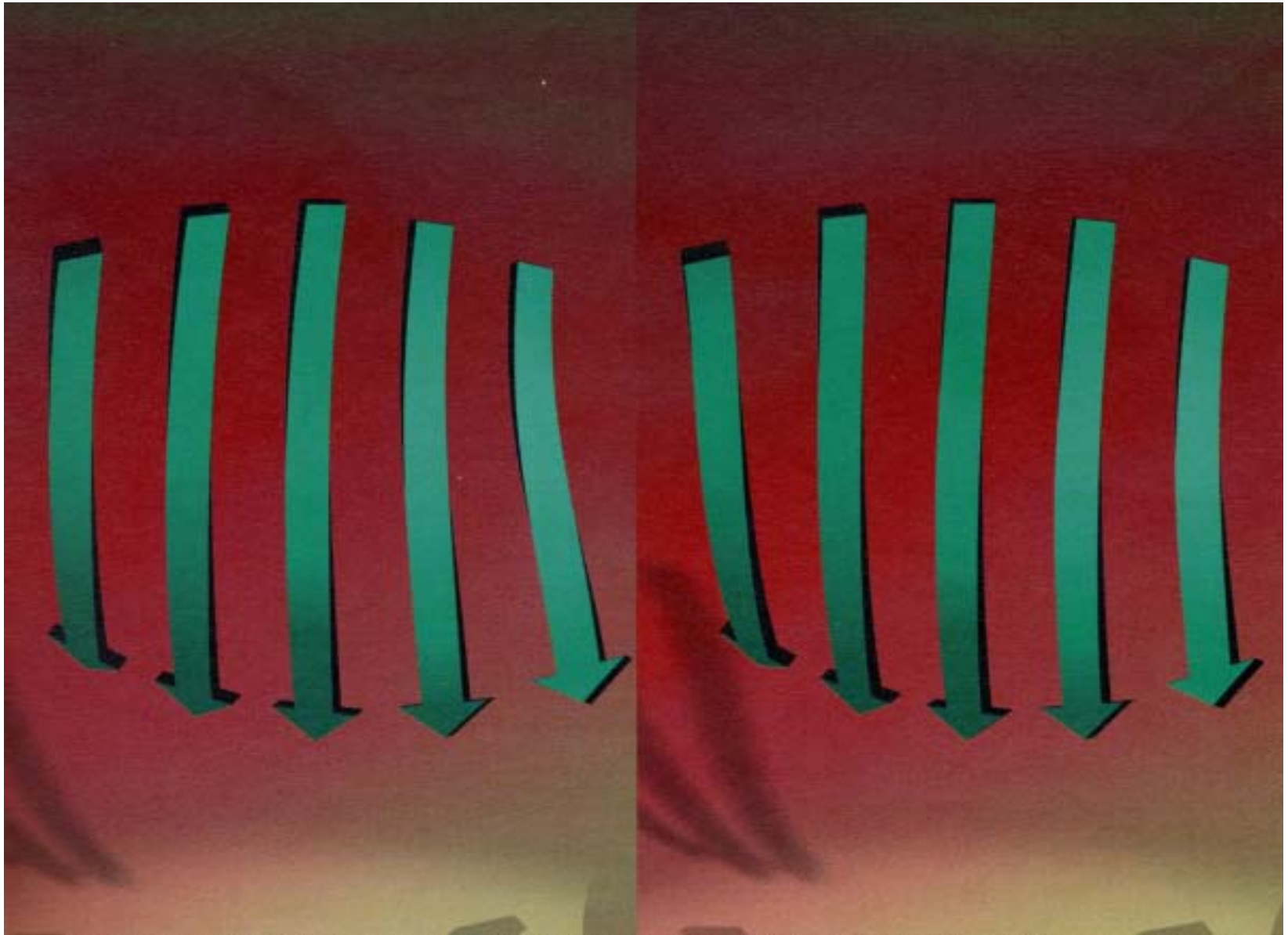
- secondary structure: **local** conformation of backbones
e.g.) α helix



e.g.) β sheet



β sheet (continued)

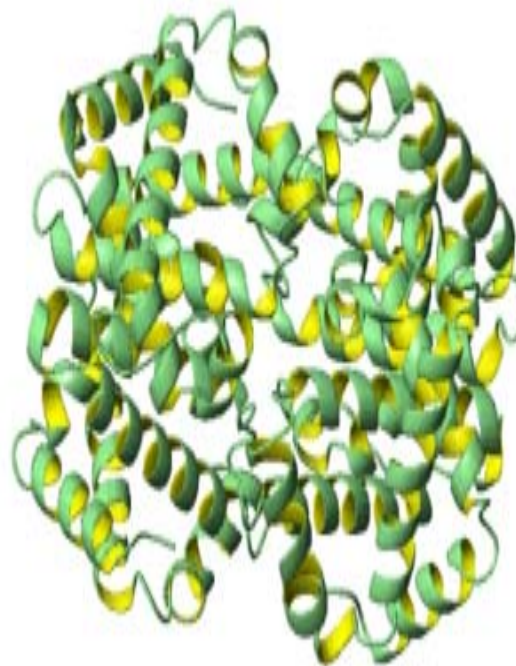


- 3 (Tertiary Structure):
overall topology of a folded protein



crambin (46 aa)

- 4 (Quaternary Structure)



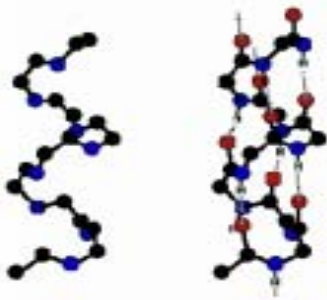
deoxy human
hemoglobin
(oxygen transport)
4 proteins
141 - 146 - 141 -
146 aa

- 1 (Primary Structure)

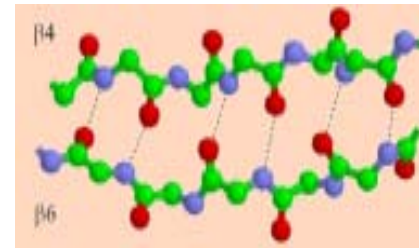
1D sequence of amino acids

e.g.) cys - gly - val - ala - ala

- 2 (Secondary Structure)



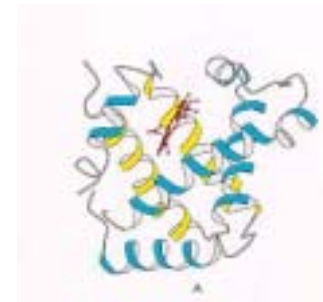
α helix



β sheet

- 3 (Tertiary Structure)

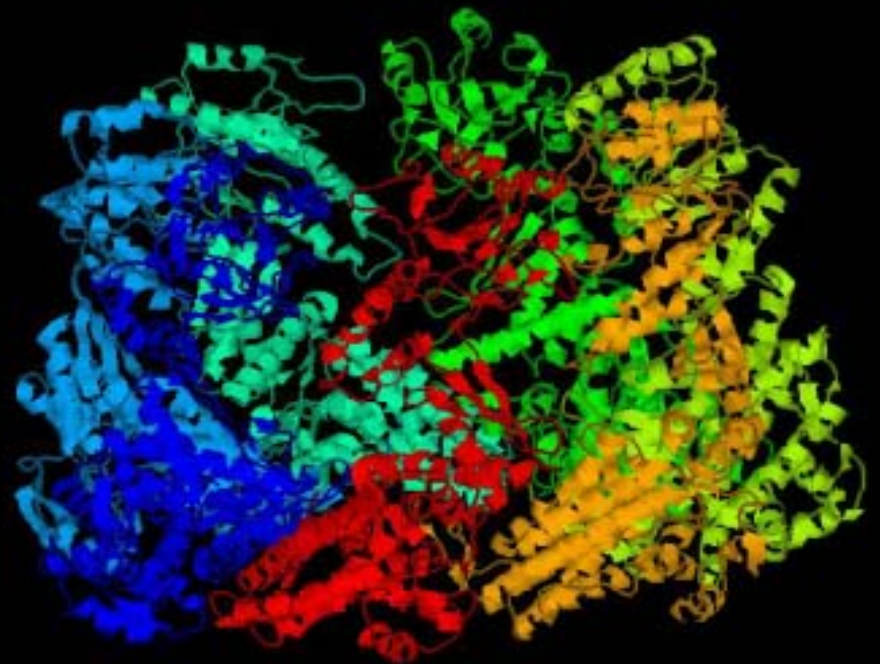
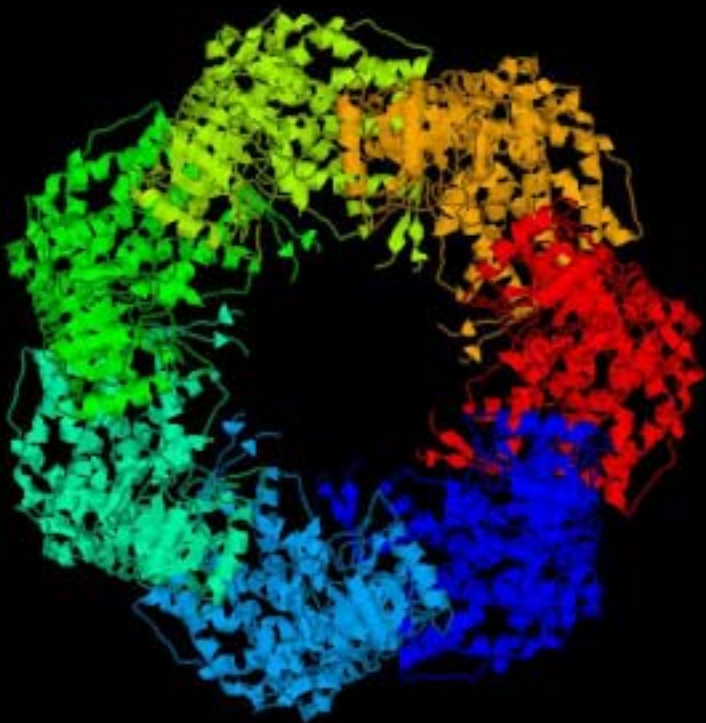
3D arrangement



GroEL

top view

side view



1. Protein Structure Determination

A. :

X-ray crystallography,

NMR

J. Kendrew (1957): myoglobin

M. Perutz (1959): hemoglobin

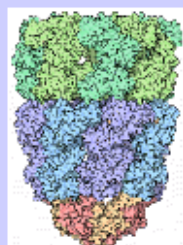
[DEPOSIT data](#)
[DOWNLOAD files](#)
[browse LINKS](#)
[BETA TEST new features](#)

Current Holdings

18488 Structures

Last Update: 20-Aug-2002

[PDB Statistics](#)



[Molecule of the Month:
Chaperones](#)

The Protein Data Bank (PDB) is operated by Rutgers, The State University of New Jersey; the San Diego Supercomputer Center at the University of California, San Diego; and the National Institute of Standards and Technology -- three members of the [Research Collaboratory for Structural Bioinformatics \(RCSB\)](#). The PDB is supported by funds from the [National Science Foundation](#), the [Department of Energy](#), and two units of the National Institutes of Health: the [National Institute of General Medical Sciences](#) and the [National Library of Medicine](#).

PROTEIN DATA BANK

Welcome to the PDB, the single worldwide repository for the processing and distribution of 3-D biological macromolecular structure data.

[ABOUT PDB](#) | [DATA UNIFORMITY](#) | [RECENT FEATURES](#) | [USER GUIDES](#) | [FILE FORMATS](#) | [EDUCATION](#) | [STRUCTURAL GENOMICS](#) | [PUBLICATIONS](#) | [SOFTWARE](#)

Search the Archive



Enter a [PDB ID](#) or keyword

[Query Tutorial](#)

Find a structure

- ☐ query by PDB id only ☐ match exact word
☐ remove sequence homologues

[SearchLite](#) keyword search form with examples

[SearchFields](#) customizable search form

[Status Search](#) find entries awaiting release

News

[Complete News
Newsletter](#)

[pdb-1 Mailing List
Subscribe](#)

20-Aug-2002

[New mmCIF and Data Processing Software Available](#) RCSB-developed programs for mmCIF and data processing -- including MAXIT, PDB_EXTRACT, and an mmCIF database loader -- are now available from the PDB's software page for download...

[\[MORE...\]](#)

PDB Mirrors

Please bookmark a mirror site

[San Diego Supercomputer Center*](#)

[Rutgers University*](#)

[National Institute of Standards and Technology*](#)

[Cambridge Crystallographic Data Centre, UK](#)

[National University of Singapore](#)

[Osaka University, Japan](#)

[Universidade Federal de Minas Gerais, Brazil](#)

[OTHER SITES](#)

**RCSB partner*

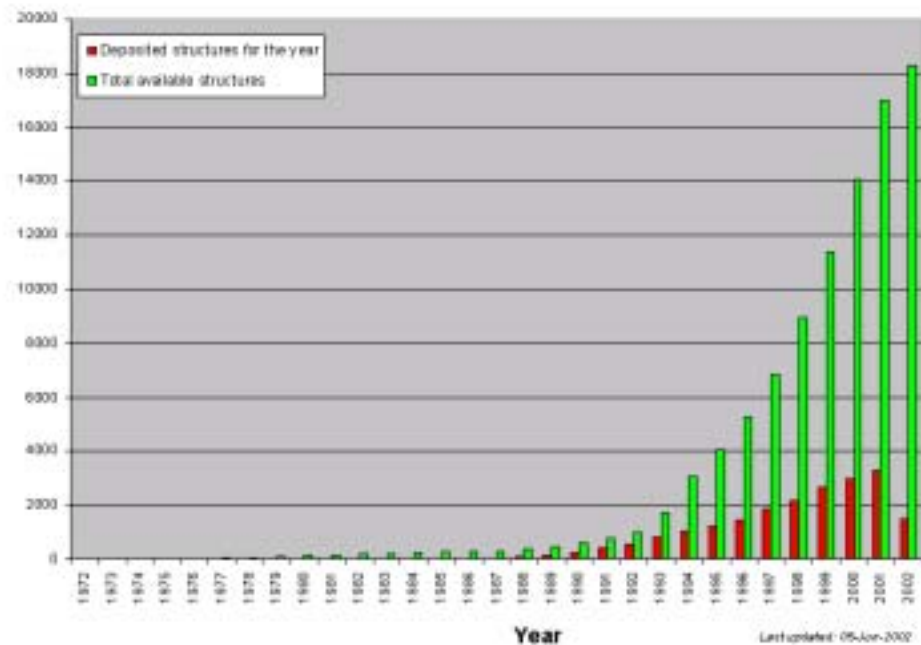
In citing the PDB please refer to:

H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Bhat, H. Weissig, I.N. Shindyalov, P.E. Bourne: [The Protein Data Bank](#). *Nucleic Acids Research*, **28** pp. 235-242 (2000)

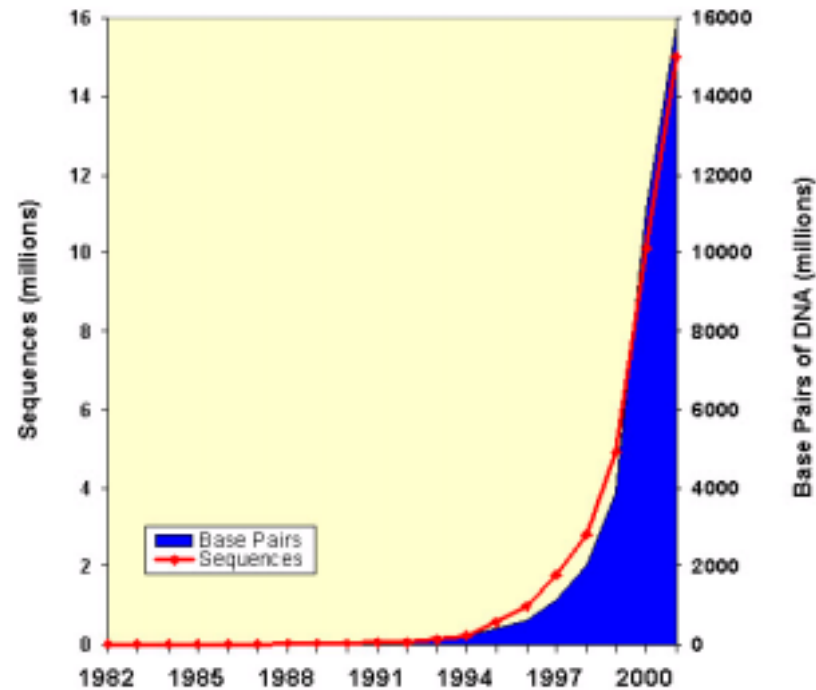
[ABOUT PDB](#) | [DATA UNIFORMITY](#) | [RECENT FEATURES](#) | [USER GUIDES](#) | [FILE FORMATS](#) | [EDUCATION](#) | [STRUCTURAL GENOMICS](#) | [PUBLICATIONS](#) | [SOFTWARE](#)

© RCSB

PDB Content Growth



Growth of GenBank



B. ()

knowledge-based methods

,
.(Protein Data Bank(PDB)
20000 가)

physics-based methods (ab initio prediction)

Knowledge-based method

- Comparative modeling (Homology modeling) :

- PDB

가 .

-

가 .

- PDB
similarity)

.(high sequence

- Fold recognition (Threading) :

- PDB

.

-

가 .

-

가

가

.

- PDB

.(medium

sequence similarity)

Sequences and Sequence alignment

- Two main kind of sequences
 - Sequence of base pairs in DNA molecules
 - (A+T+C+G)*
 - Sequence of aminoacids in a protein molecule
 - A(C+D+E+F+G+H+I+K+L+M+N+P+Q+R+S+T+V+W+X+Y)*Z
- Two main kind of sequence alignment
 - Global alignment
 - LGPSSKQT GKG S - -RIWDN
 - | | ||| | |
 - LN -ITKSA GKGAIMRLGDA
 - Local alignment
 - -----TGKG-----
 - |||
 - -----AGKG-----

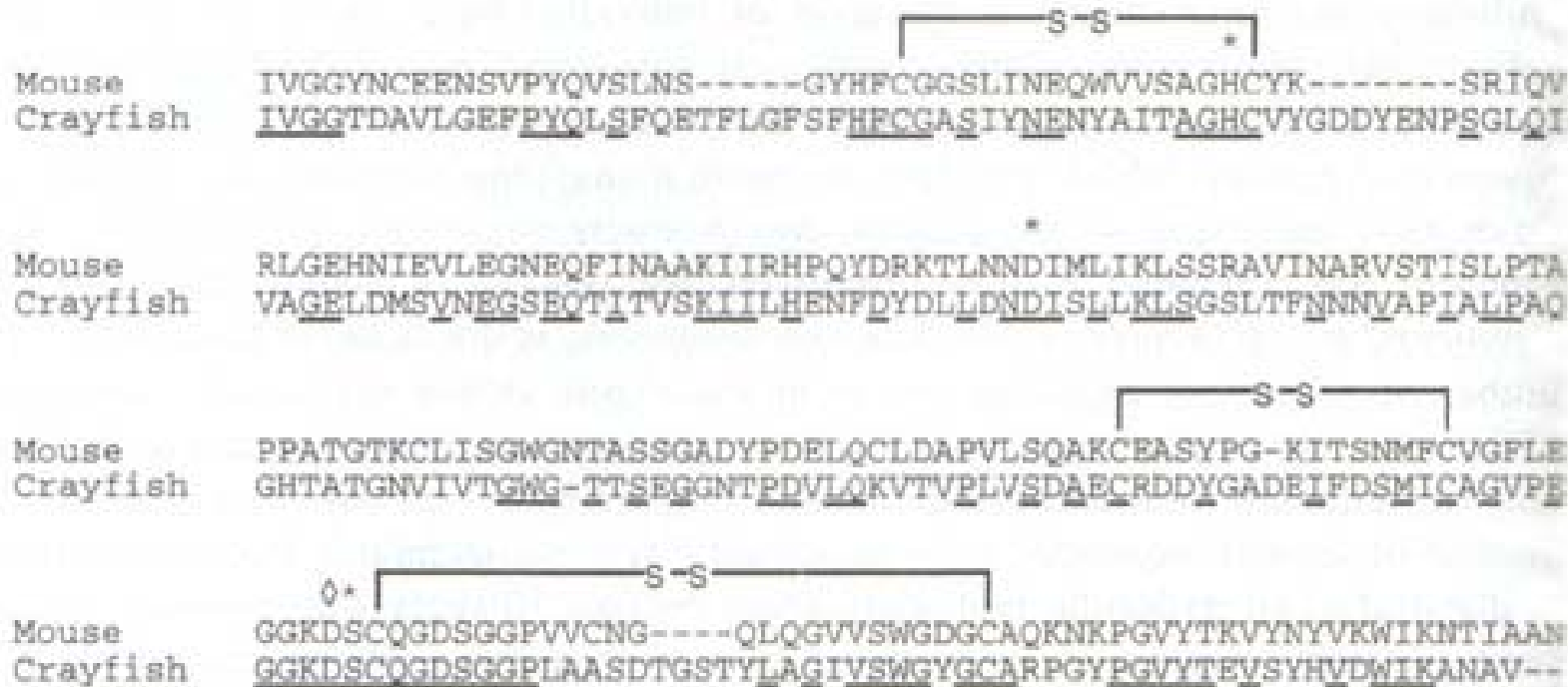
A BLOSUM62 scoring matrix

[illegible]

Basis of Basic Sequence Alignment

The objective of sequence alignment is to determine if 2 sequences are sufficiently similar to declare them homologous.

Here is an example:



The above illustration shows the alignment of trypsin proteins of mouse (SWISS-PROT P07146) and crayfish (SWISS-PROT P00765). Identical residues are underlined. Indicated above the alignments are 3 disulfide bonds (-S-S-) with participating cysteine residues conserved, amino acid side chains involved in the charge relay system (asterisk), and active site residue governing substrate specificity (diamond).

Ab-initio(energy-based,Physics-based) method

- - PDB
 -
- - (no sequence similarity)
 -

- There was a student who knew nothing about quantum mechanics
- This poor student took a “quantum mechanics” course.
- He had to take a take-home exam: Not to drop out of the class, he has two options to choose:

1. Examine last 40 year's problem set with answers. → Homology modeling / Threading.
2. Try to understand the problems and write down his own answers. → Ab initio (de novo, Energy-based,...).

Physics-based method ←

:

1

3

. (Anfinsen, 1973)



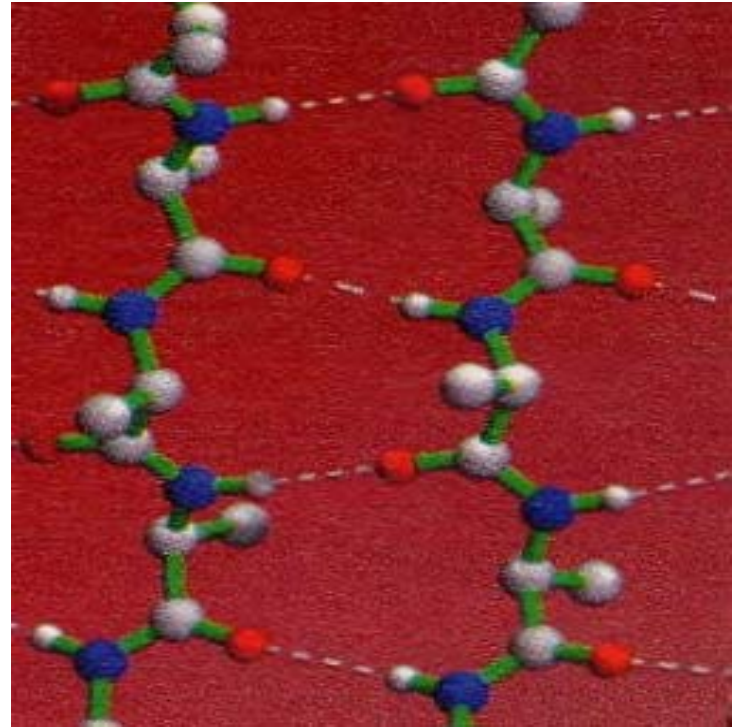
,

F가

.

$F=E-TS$ (E: energy, T: temperature,
S: entropy)

(Electrostatic interaction
between polar atoms,
Van der Waals force,
Disulfide bridge,
Etc.)



Various Potential Energy Functions

- All-atom Consideration
- Coarse-grained Function
Atom-centered, Fixed charge...
- Contacts only : Scoring Function

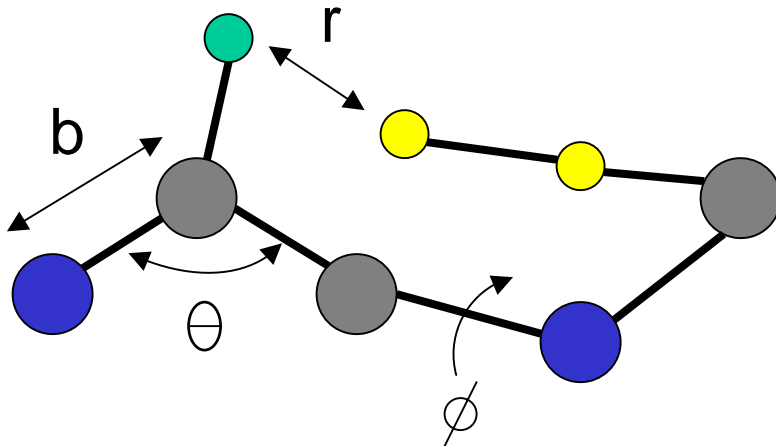
Potential energy functions

1. All atom, off-lattice potentials

- ECEPP, AMBER, CHARMM, and more
- E terms: vibrational, torsional, non-bonded, electrostatic

$$E = \frac{1}{2} \sum_{\text{bonds}} k_b (b - b_{\text{eq}})^2 + \frac{1}{2} \sum_{\substack{\text{bonds} \\ \text{angles}}} k_{\theta} (\theta - \theta_{\text{eq}})^2 +$$

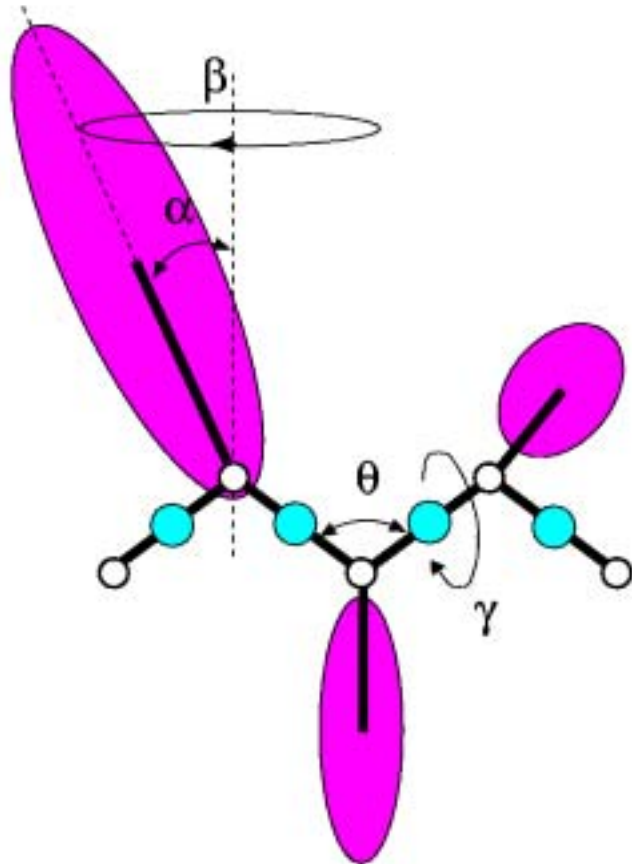
$$\sum_{\substack{\text{dihedral} \\ \text{angles}}} k_{\phi} \cos(n\phi - \delta) + \sum_{ij} \left(\frac{A_{ij}}{r_{ij}^6} + \frac{B_{ij}}{r_{ij}^{12}} + \frac{C_{ij}}{r_{ij}^{10}} + \frac{q_i q_j}{D r_{ij}} \right)$$



2. Coarse grained potentials:

e.g) UNRES(United-residue potential)

- Fixed bond (virtual bond) lengths
- Two interacting centers per residue
- Can treat larger molecules in reasonable CPU time
- Off-lattice model

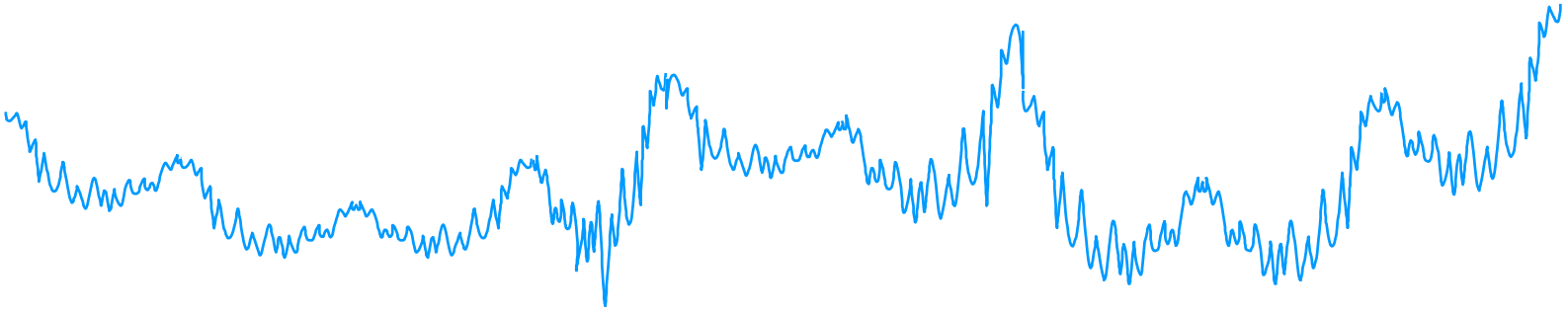


Definition of
degrees of
freedom for the
UNRES
representation of a
polypeptide chain

Side-chains are represented as
ellipsoids (Gay-Berne potential)

Interaction centers are marked
in colors

(Global optimization)



(global minimum)

●

(Global optimization)

- **Simulated Annealing (SA)**
- **Genetic Algorithms (GA)**
- **Hopfield Neural Network**
- **Monte Carlo with Minimization**

.

.

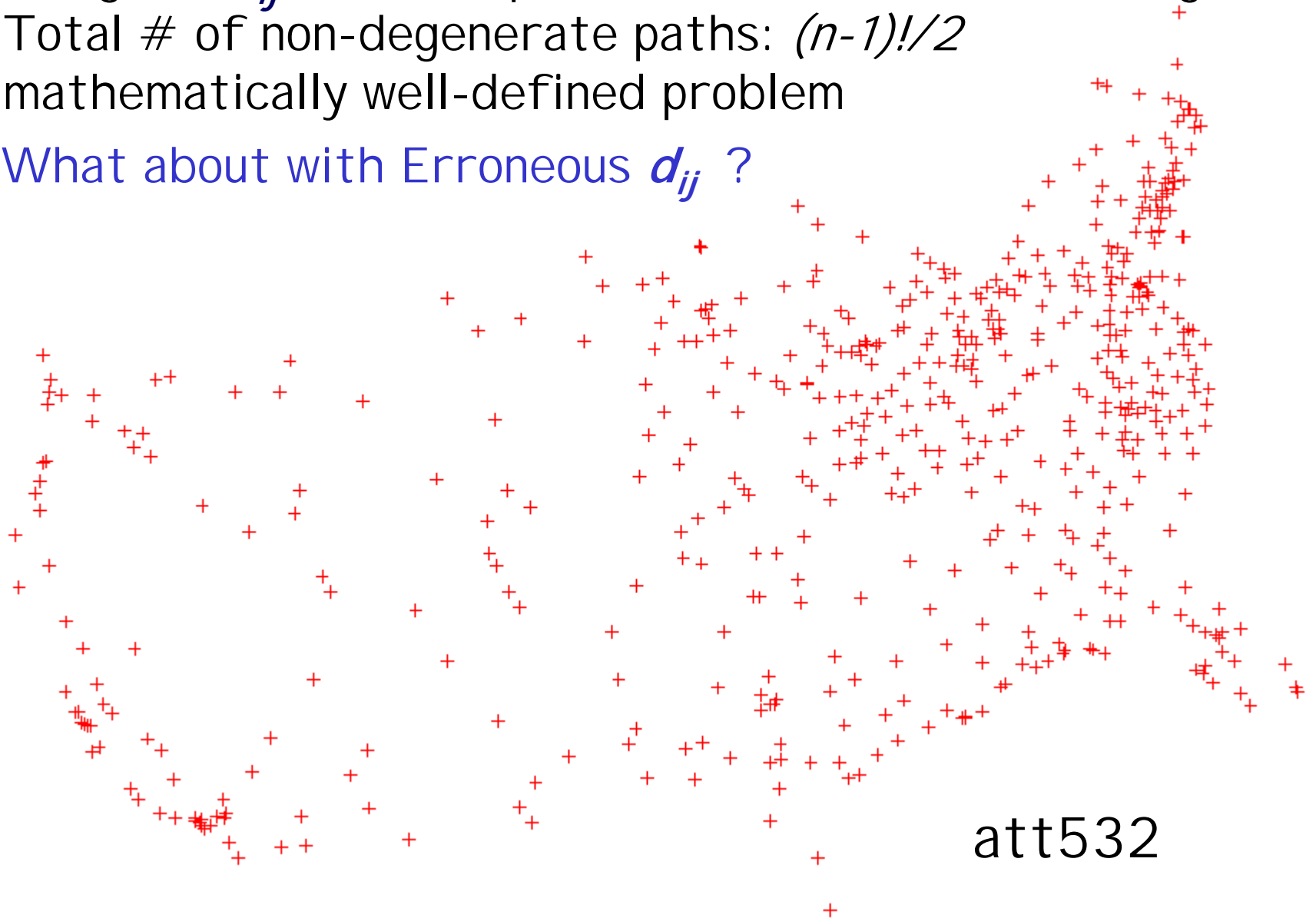
.

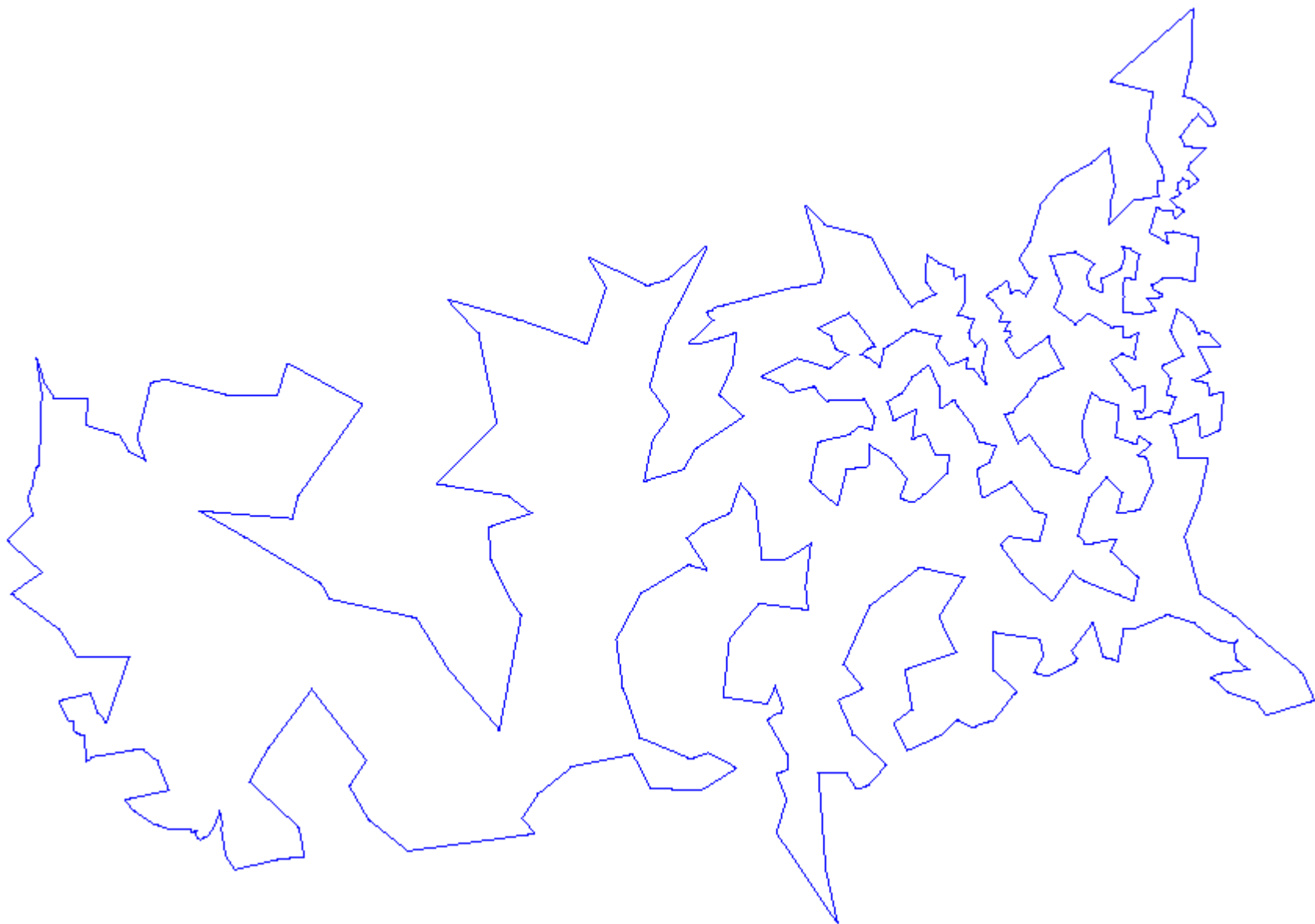
Protein Folding Problem

1. (protein) 가?
2. (protein synthesis) Central Dogma
- 3. Protein folding problem vs. Traveling salesman problem

Traveling Salesman Problem

- For given d_{ij} , find the path of the shortest tour length
- Total # of non-degenerate paths: $(n-1)!/2$
- mathematically well-defined problem
- What about with Erroneous d_{ij} ?







$E(x) = \dots$



$E(x) = \dots$

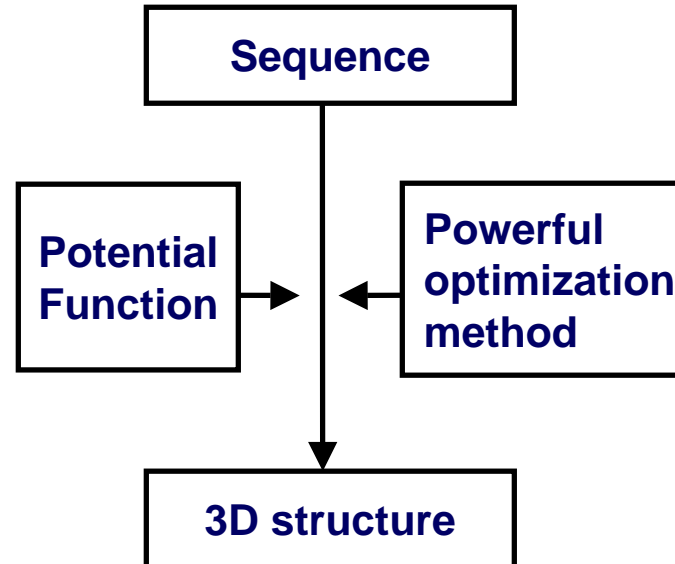


$E(x) = \dots$



$E(x) = \dots$

Ab initio protein folding by computer simulation



Protein Folding Problem

1. (protein) 가?
2. (protein synthesis) Central Dogma
3. Protein folding problem vs. Traveling salesman problem
- 4. **Ab. initio protein structure prediction**

[CASP1](#)[CASP2](#)[CASP3](#)[CASP4](#) ✓[Local services](#)[Other links](#)[People](#)[Website index](#)[Hide menu](#)

Protein Structure Prediction Center

Biology and Biotechnology Research Program
Lawrence Livermore National Laboratory, Livermore, California, USA



Welcome to the Protein Structure Prediction Center!

Our goal is to help advance the methods of identifying protein structure from sequence. The Center has been organized to provide the means of objective testing of these methods via the process of blind prediction. In addition to support of the CASP meetings our goal is to promote an objective evaluation of prediction methods on a continuing basis.

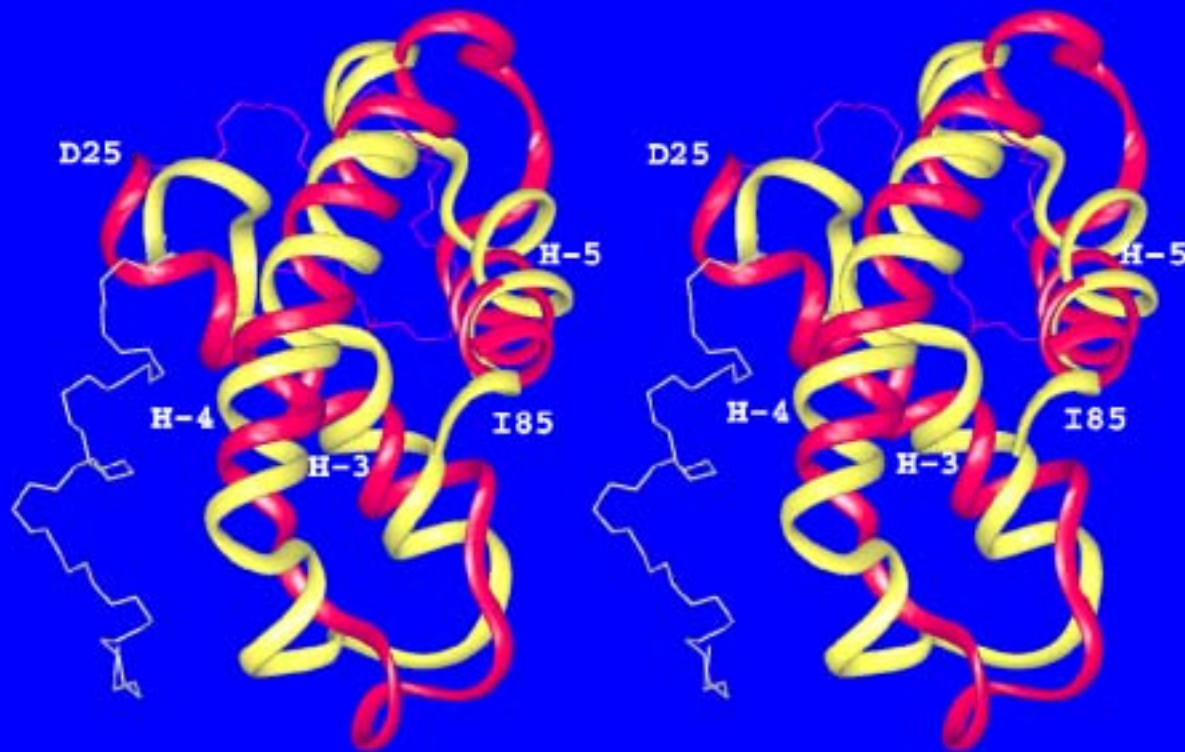
CASP experiment: [CASP1](#) | [CASP2](#) | [CASP3](#) | [CASP4](#)

Ten Most Wanted: [TMW](#)

The Center, supported by the U.S. Department of Energy, [Office of Biological and Environmental Research](#), is a part of the [Biology and Biotechnology Research Program](#) at the [Lawrence Livermore National Laboratory](#).

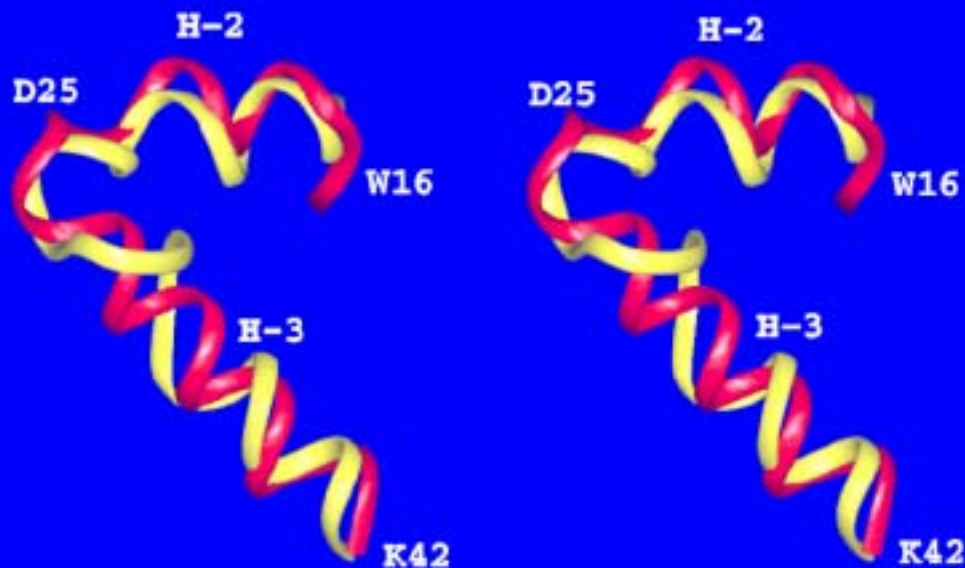
[Local services](#) | [Other links](#) | [People](#) | [Website index](#)

[CASP4 protein structure prediction evaluation data](#)



HDEA

RMSD=4.2 Å for 61
residues (80%,
residues 25-85)



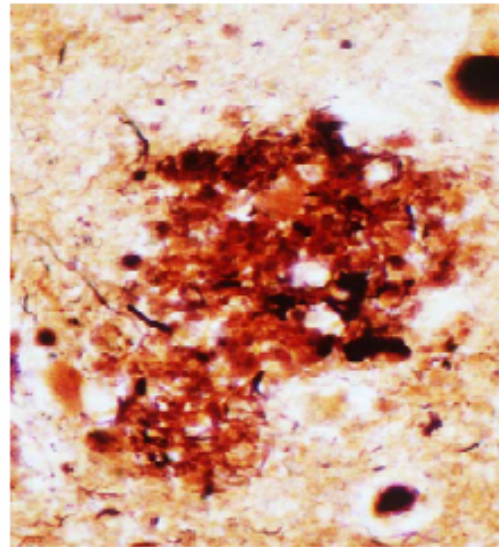
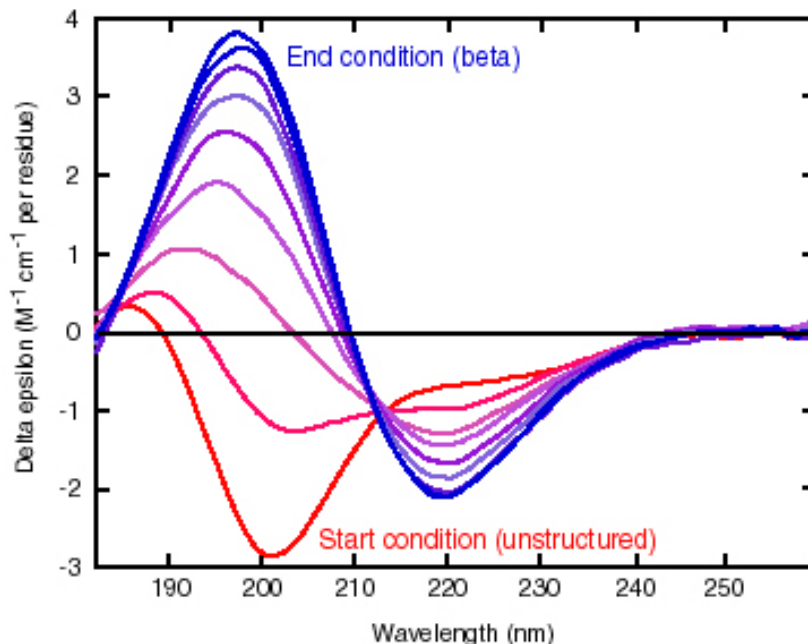
HDEA Segment

RMSD=2.9 Å for 27
residues (36%, residues
16-42)

Amyloid diseases

A number of diseases (e.g. Alzheimer's, CJD, BSE) involve the folding of proteins and peptides into beta-sheet structures which can polymerise, forming insoluble plaques in nerve tissue (below right).

A model for the Alzheimer's peptide is LRRN, which forms spontaneously into gels with a β -sheet structure.



SRCD spectra* (left) taken during the polymerisation of LRRN peptide show that the rate of polymerisation varies with substitution of a single amino acid residue.

*Collaboration with N.Gay and M. Symmons, Cambridge University

The SRCD data provide important information about the processes involved in polymerisation, and may lead to the development of drugs to treat these diseases.

Protein folding research topics:

- Protein structure prediction
- Protein folding mechanism: MD,MC simulations
- Docking problems
- Secondary structure prediction
- Contact prediction
- Order/disorder prediction
- Multiple sequence alignment
- Potential parameter optimization
- Global optimization of various systems: TSP, molecular clusters, etc.