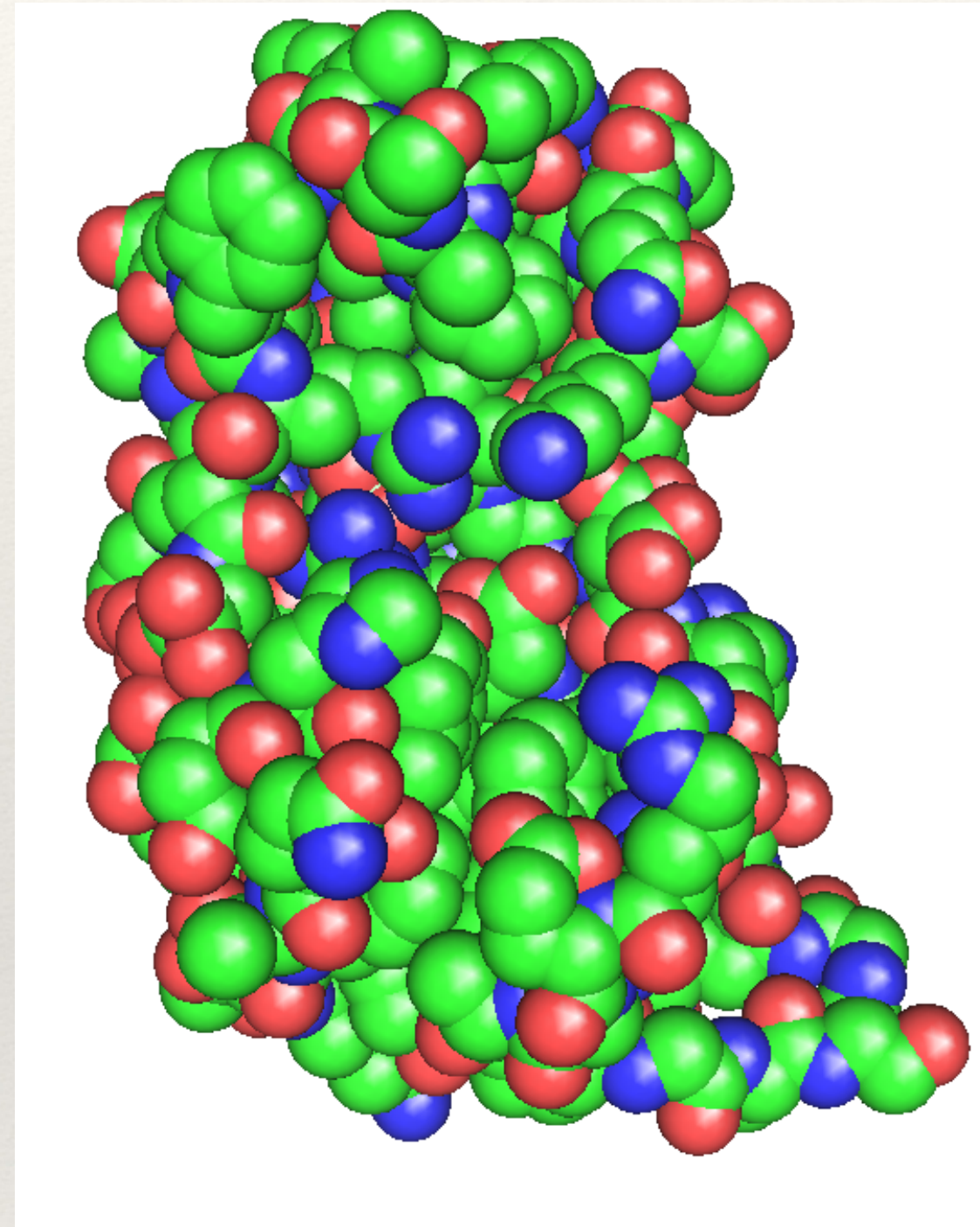
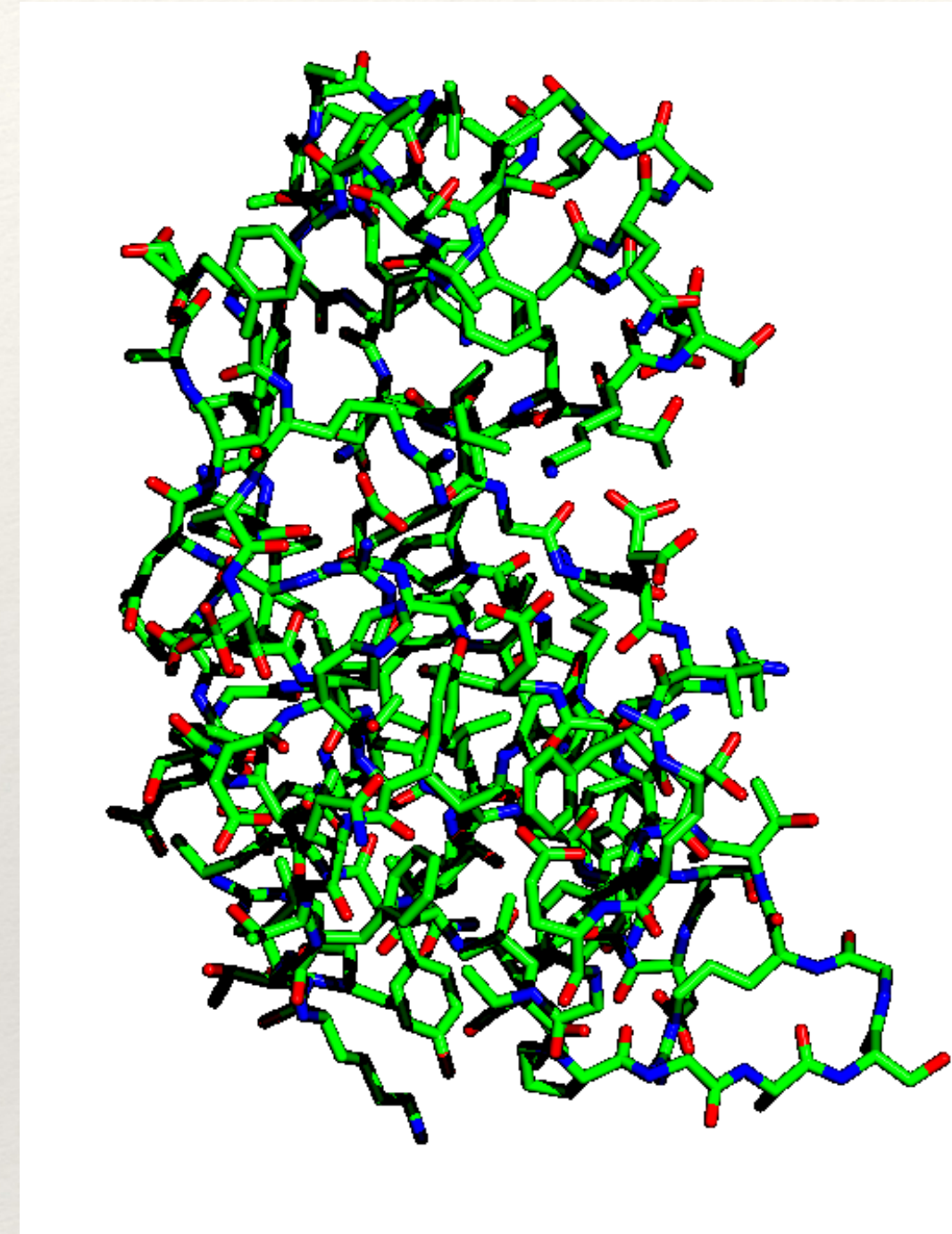

Protein Structure Prediction

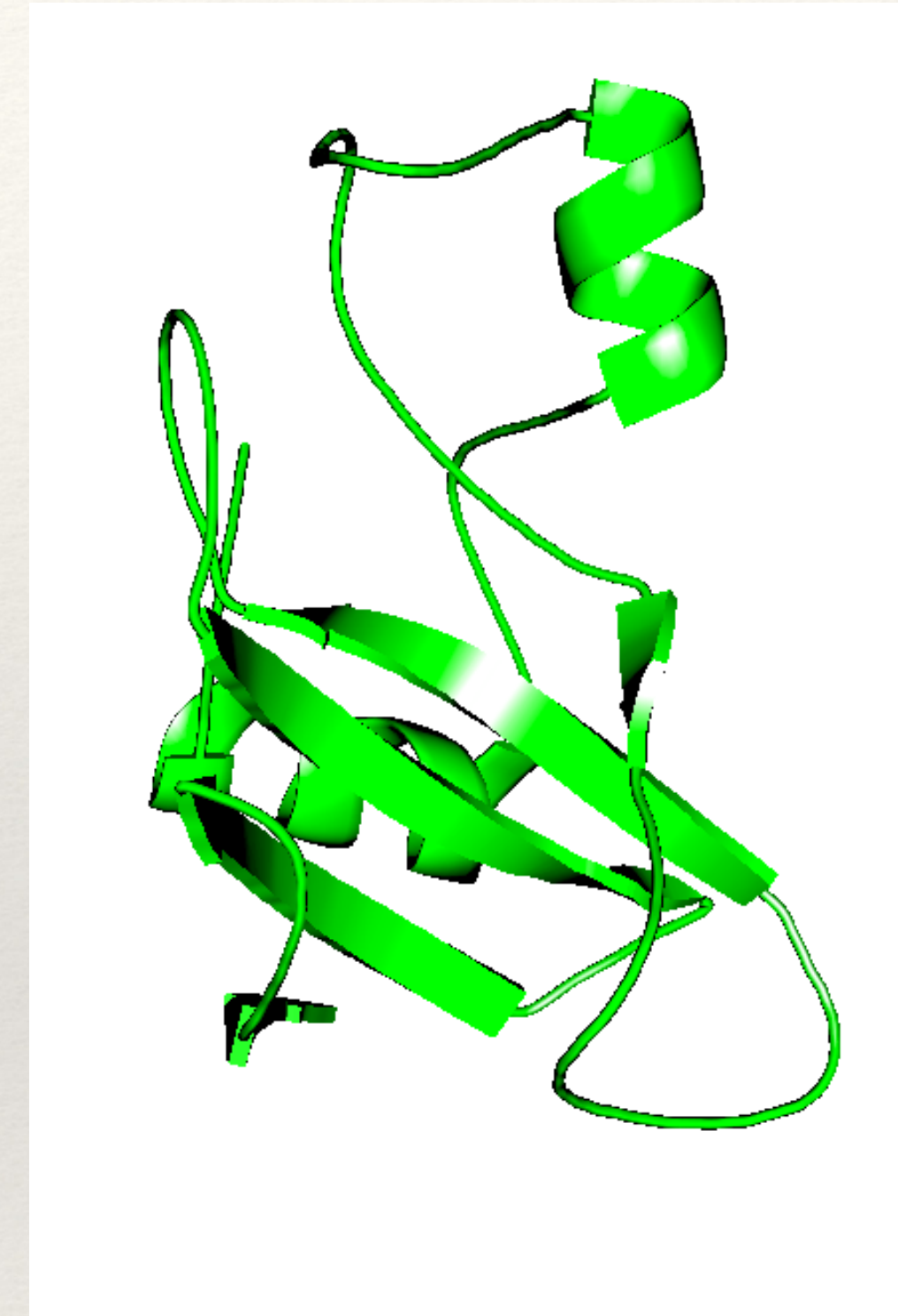
Protein Structure Representation



CPK: hard sphere model



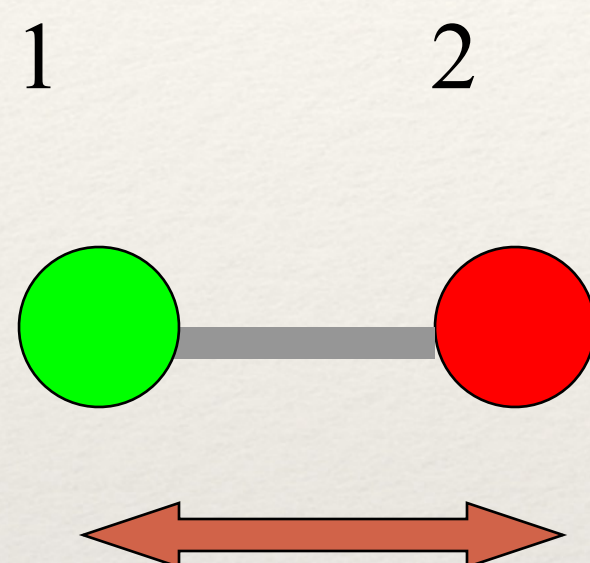
Ball-and-stick



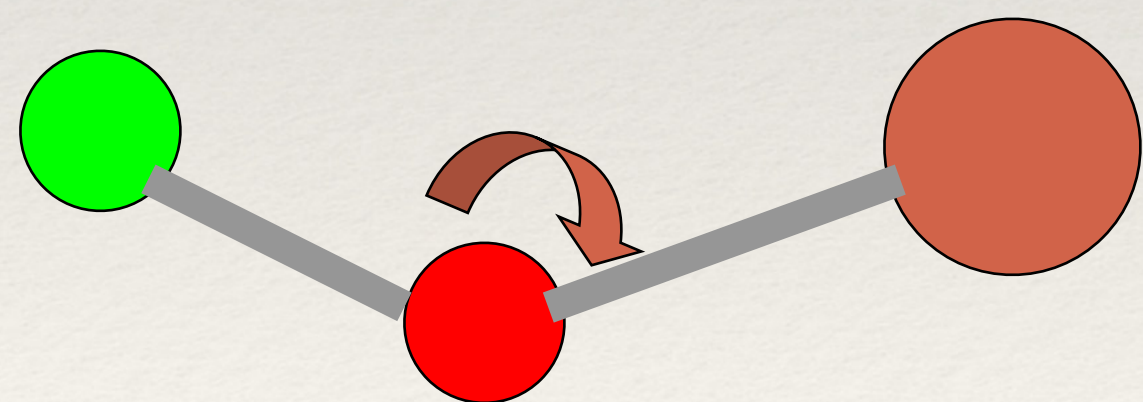
Cartoon

Degrees of Freedom in Proteins

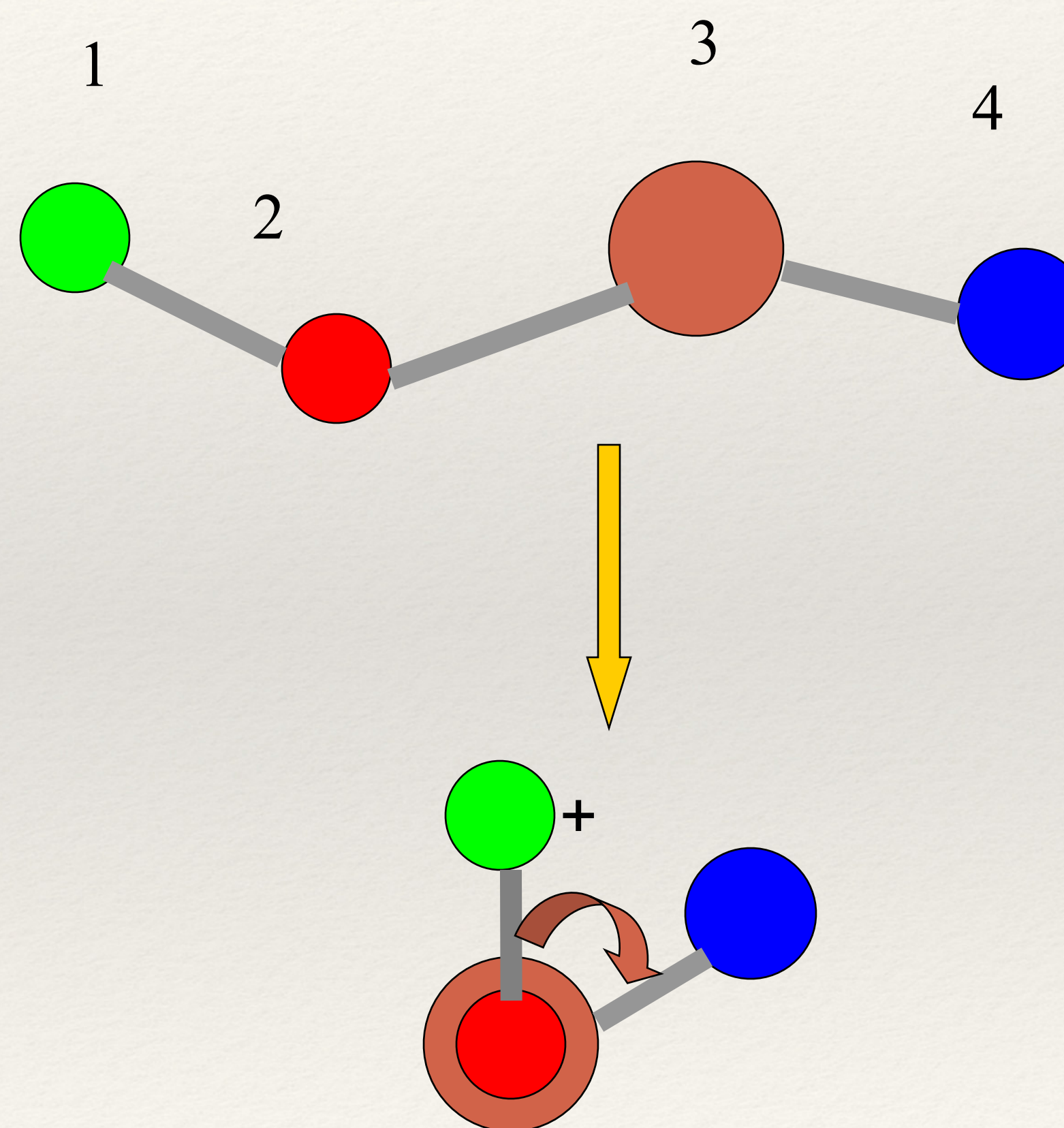
Bond length



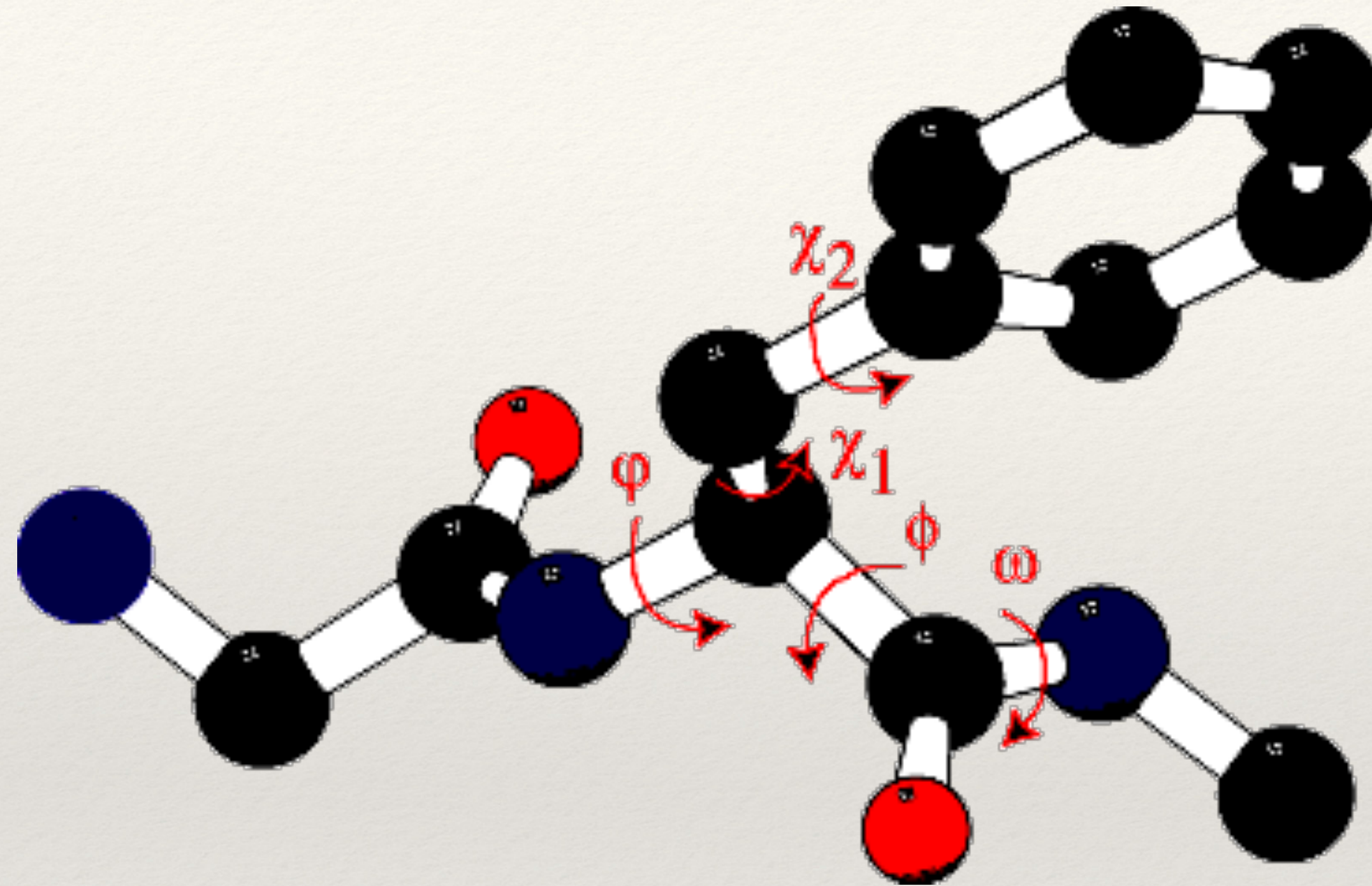
Bond angle



Dihedral angle



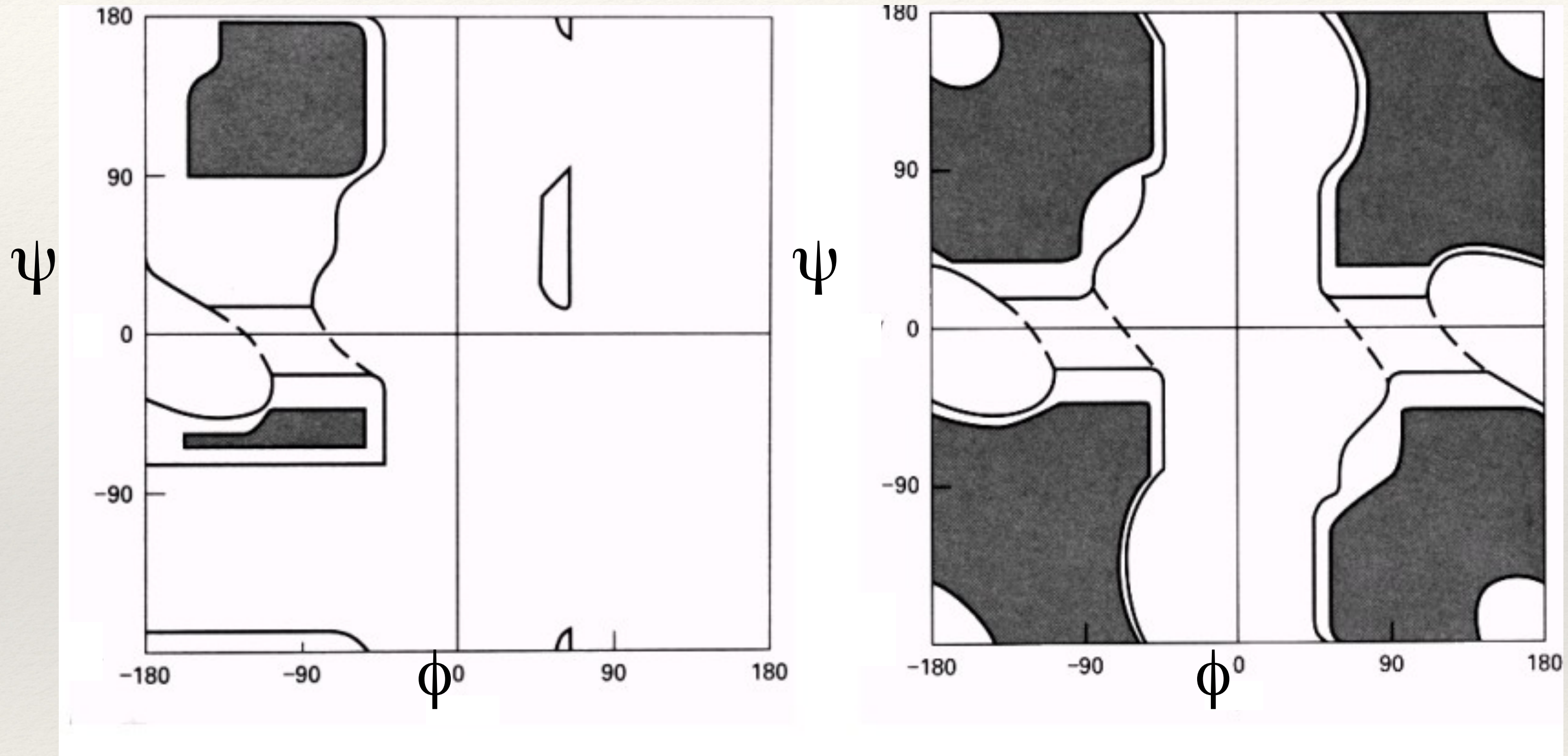
Protein Structure Representation



Backbone: 3 angles per residue : φ , ψ and ω

Sidechain: 1 to 7 angles, χ ; each χ has 3 favored values: 60° , -60° , 180° .

Ramachandran Plots



All residues, but glycine

Glycine

Root Mean Square Distance (RMSD)

To compare two sets of points (atoms) $A=\{a_1, a_2, \dots, a_N\}$ and $B=\{b_1, b_2, \dots, b_N\}$:

-Define a 1-to-1 correspondence between A and B

for example, a_i corresponds to b_i , for all i in $[1, N]$

-Compute RMS as:

$$RMS(A, B) = \sqrt{\frac{1}{N} \sum_{i=1}^N d(a_i, b_i)^2}$$

$d(A_i, B_i)$ is the Euclidian distance between a_i and b_i .

Root Mean Square Distance (RMSD)

- Simplified problem: we know the correspondence between set A and set B
- We wish to compute the rigid transformation T that best align a_1 with b_1 , a_2 with b_2 , ..., a_N with b_N
- The error to minimize is defined as:

*Old problem, solved in Statistics,
Robotics, Medical Image Analysis,
...*

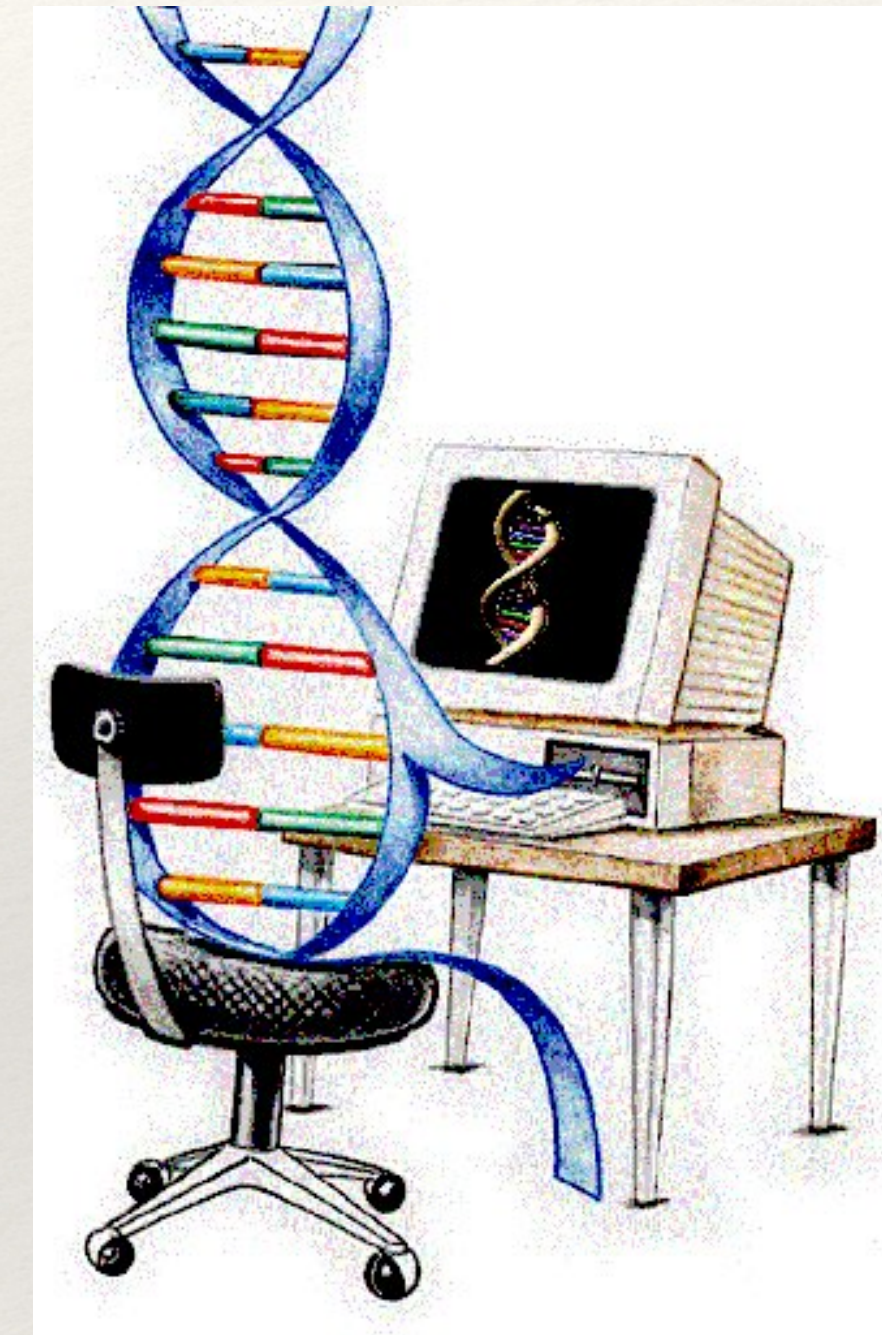
$$\varepsilon = \min_T \sum_{i=1}^N \|T(a_i) - b_i\|^2$$

Structural Bioinformatics: Proteins

Proteins: Sources of Structure Information

Proteins: Homology Modeling

Proteins: Secondary structure prediction



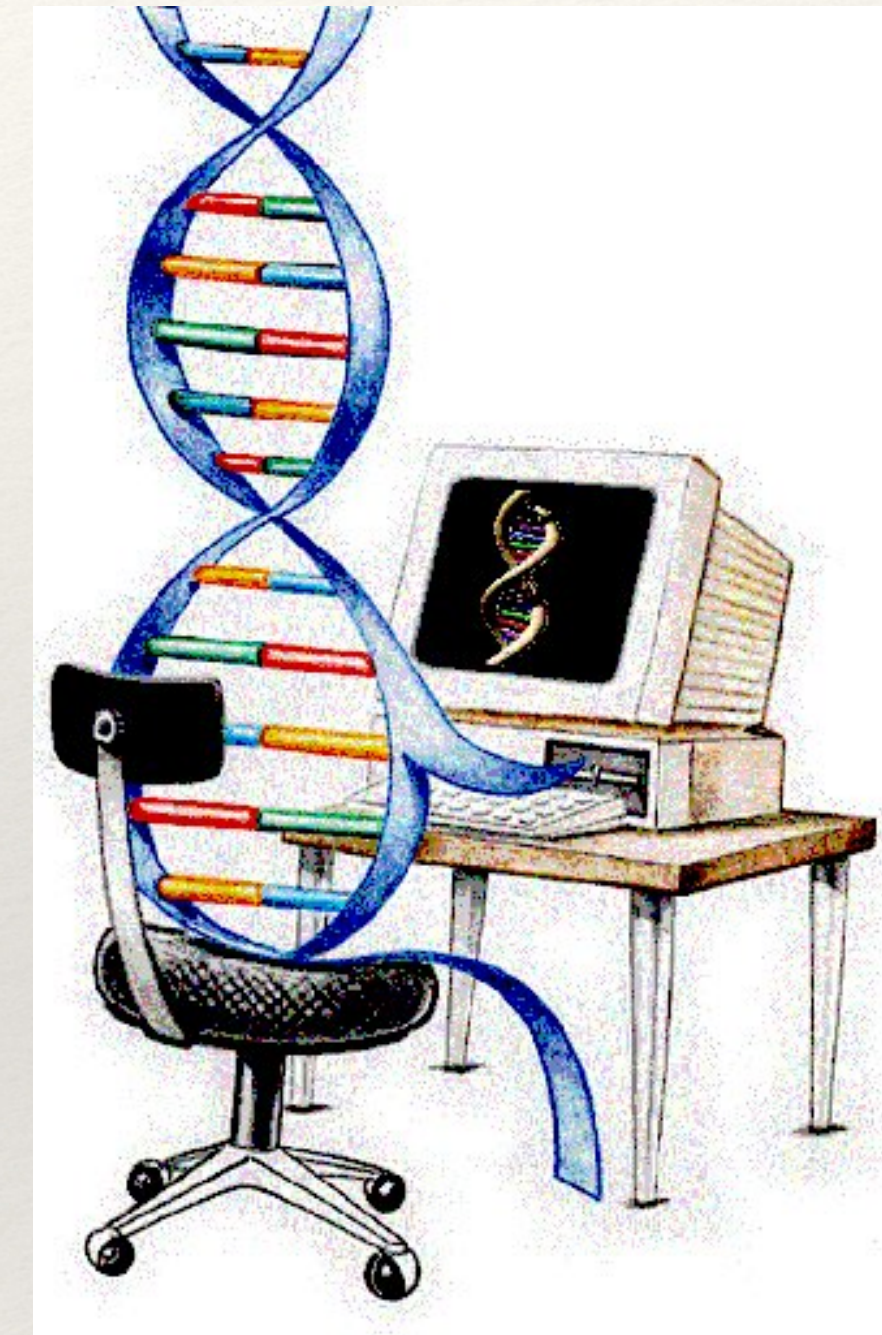
Structural Bioinformatics: Proteins

Proteins: Sources of Structure Information

Proteins: Homology Modeling

Proteins: Secondary structure prediction

Proteins: Ab initio prediction



Proteins: Finding the Primary Structure

Methods for finding the sequence of a protein:

-Translating gene sequence

- For proteins from prokaryotes, direct translation
- For proteins from eukaryotes, we need the sequence of mRNA or cDNA

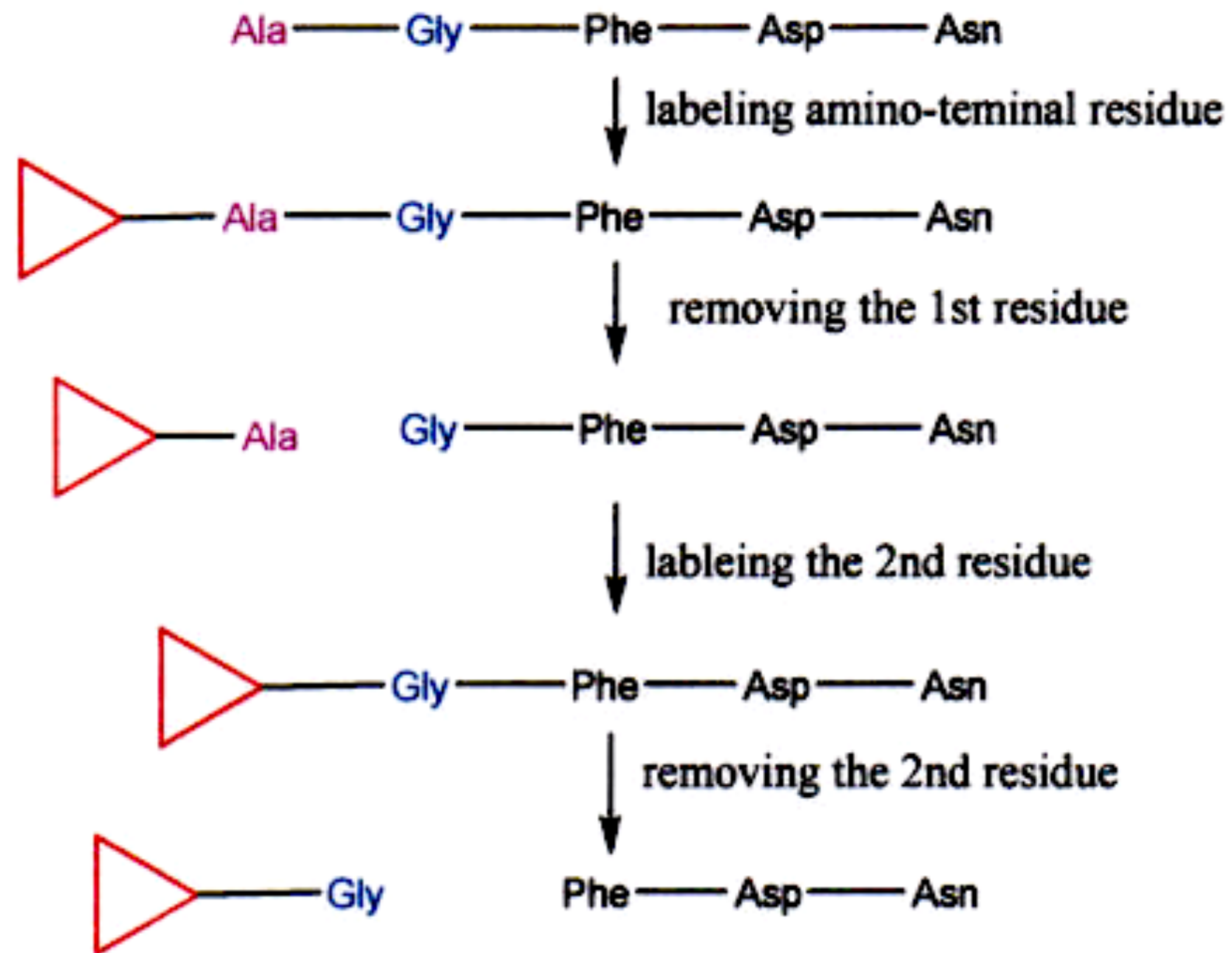
-Edman degradation

limited to “small” proteins, up to 50 amino acids
for automated sequencer

-Mass spectrometry

Proteins: Finding the Primary Structure

EDMAN DEGRADATION



(Phenyl isothiocyanate or PITC)

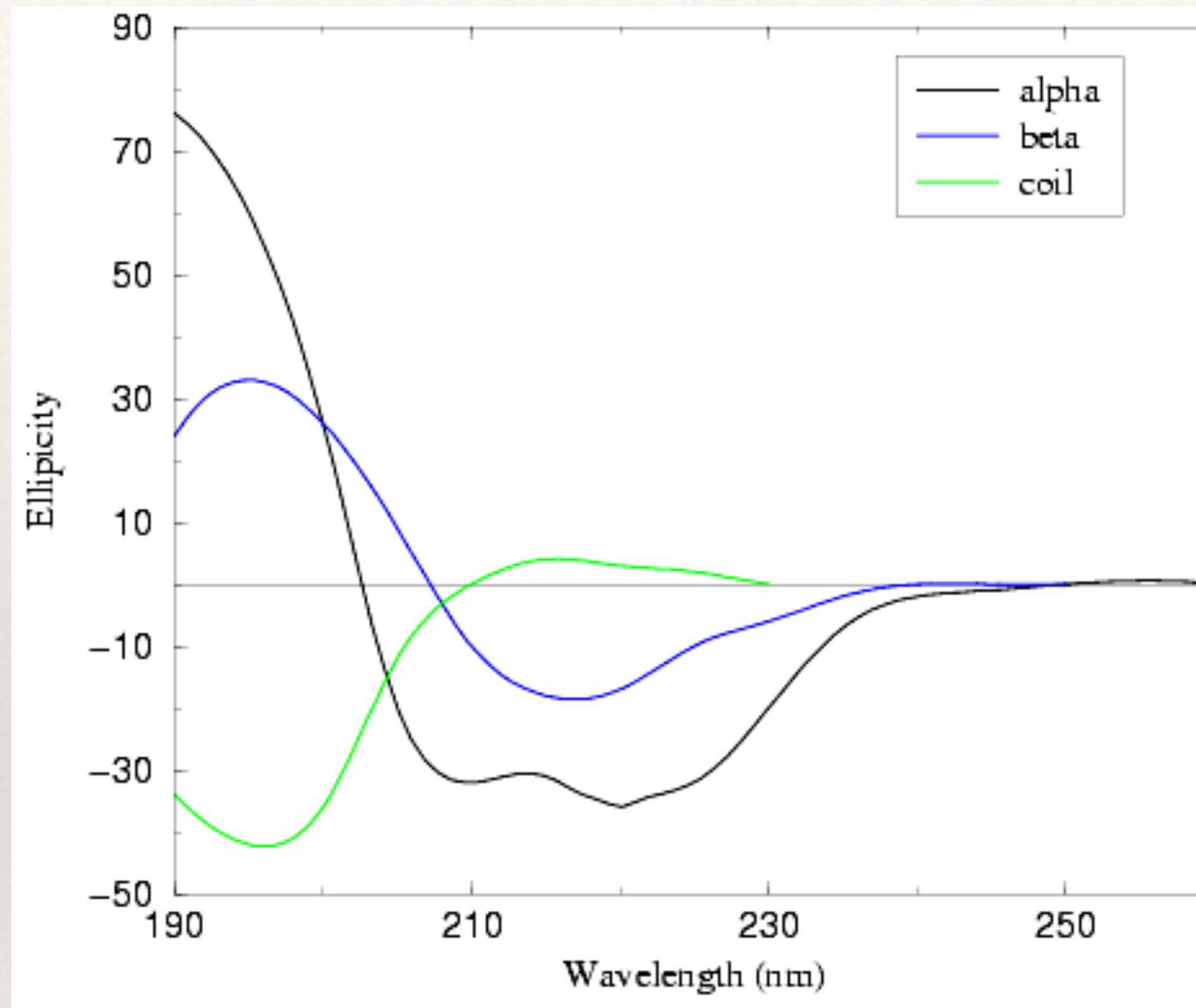
(Trifluoroacetic acid, TFA)

Proteins: Finding the Tertiary Structure

Methods for finding the 3D structure of a protein:

- **Circular Dichroism**
(low resolution; provides information on secondary structure)
- **X-ray crystallography**
(high resolution; finds structure of a protein in a crystal)
- **NMR spectroscopy**
(high resolution; finds structure of a protein in solution)

Proteins: Circular Dichroism

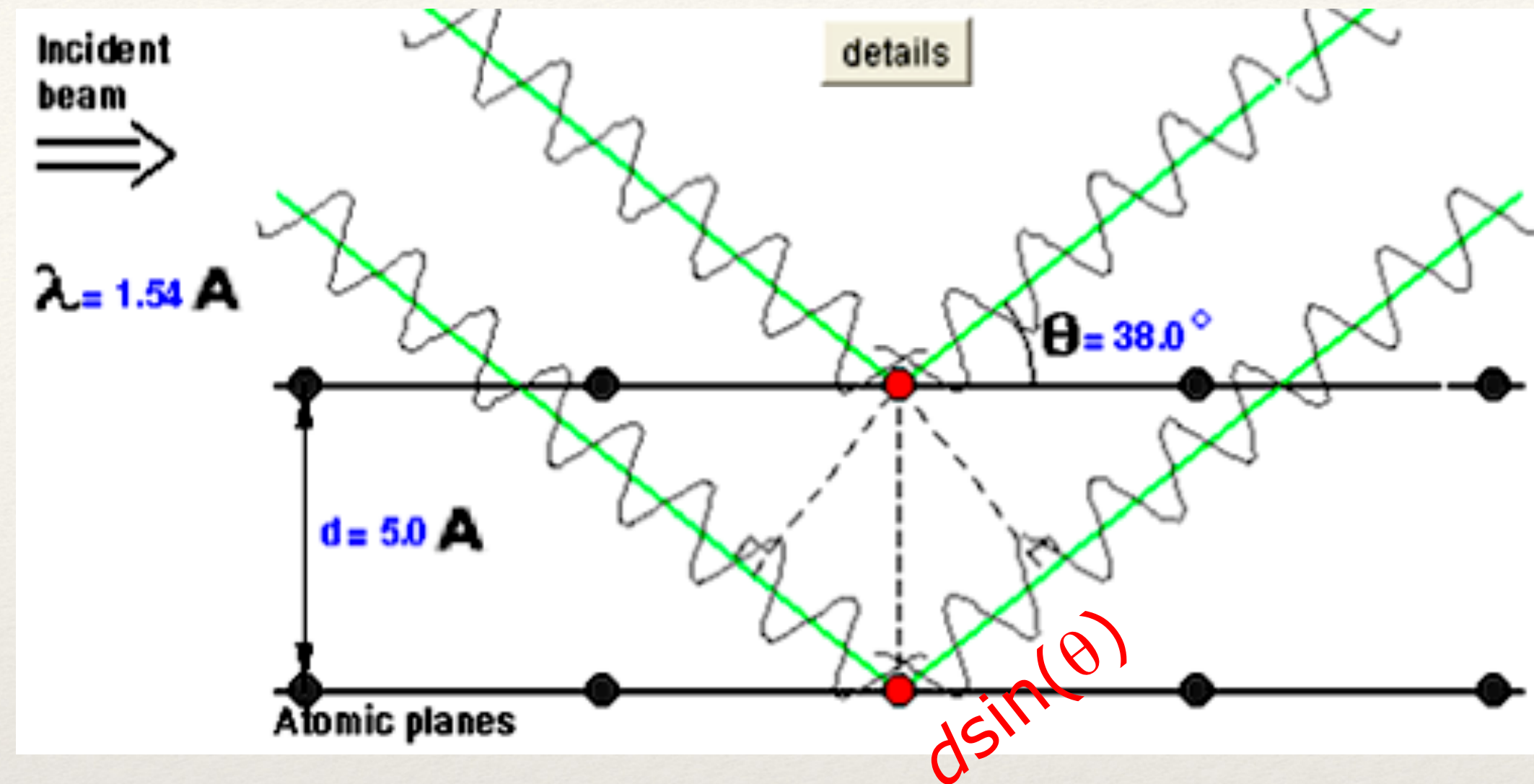


Circular dichroism (CD) spectroscopy measures differences in the absorption of left-handed polarized light versus right-handed polarized light which arise due to structural asymmetry.

Different secondary structures in proteins have different CD spectra as they have different asymmetry.

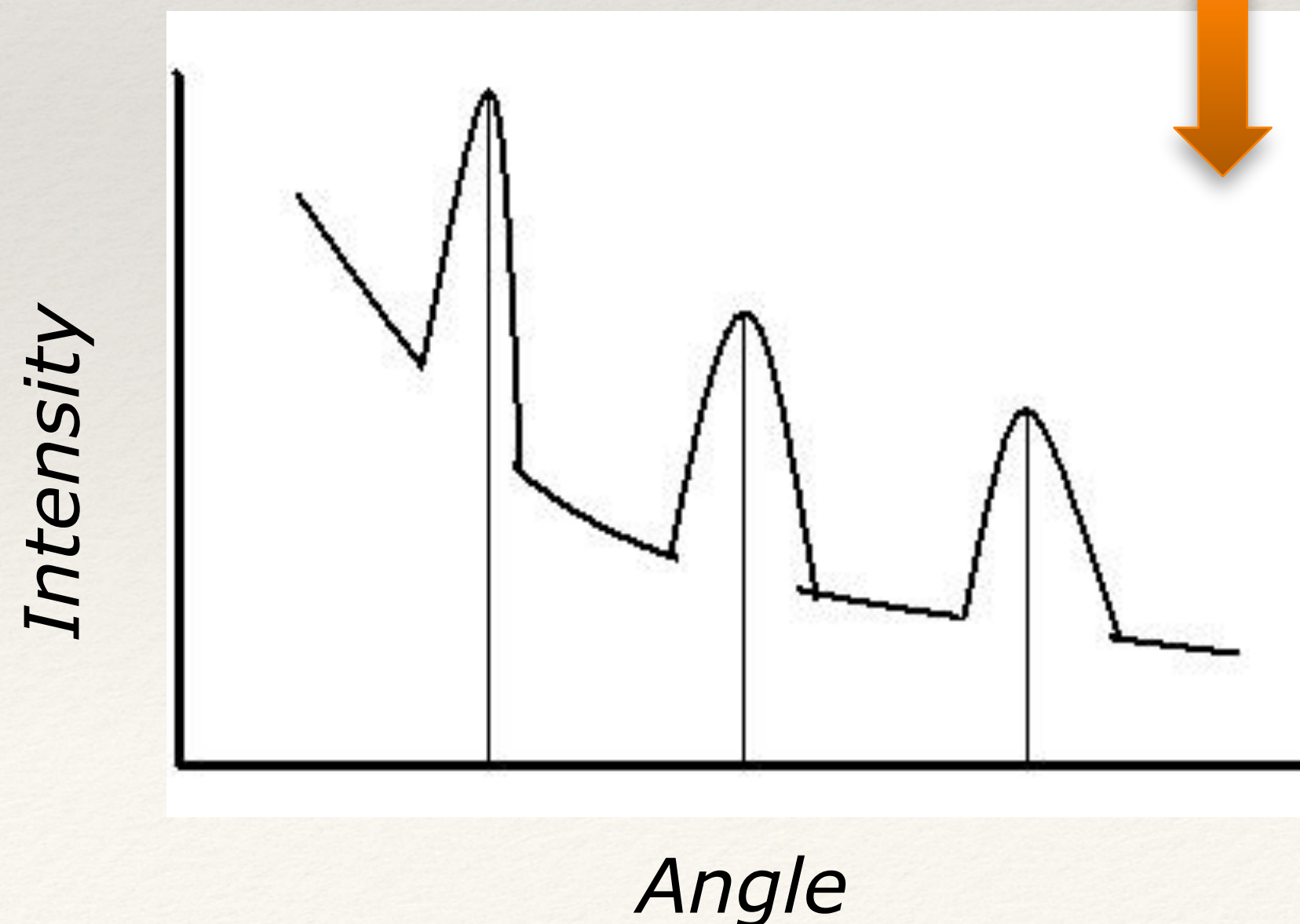
CD therefore can detect secondary structures in protein

Proteins: X-ray Crystallography

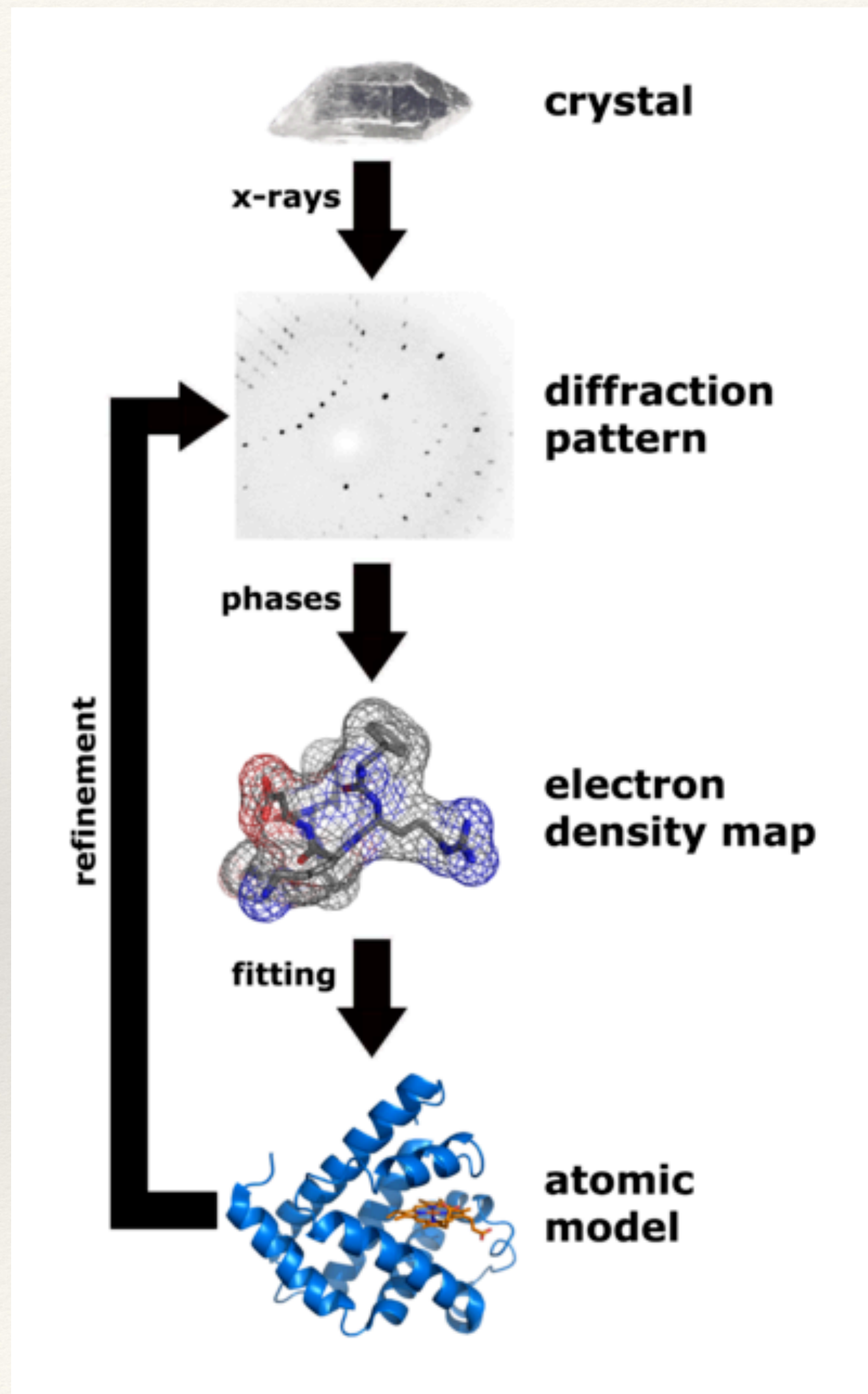


Bragg's Law:

$$2d \sin(\theta) = n\lambda$$



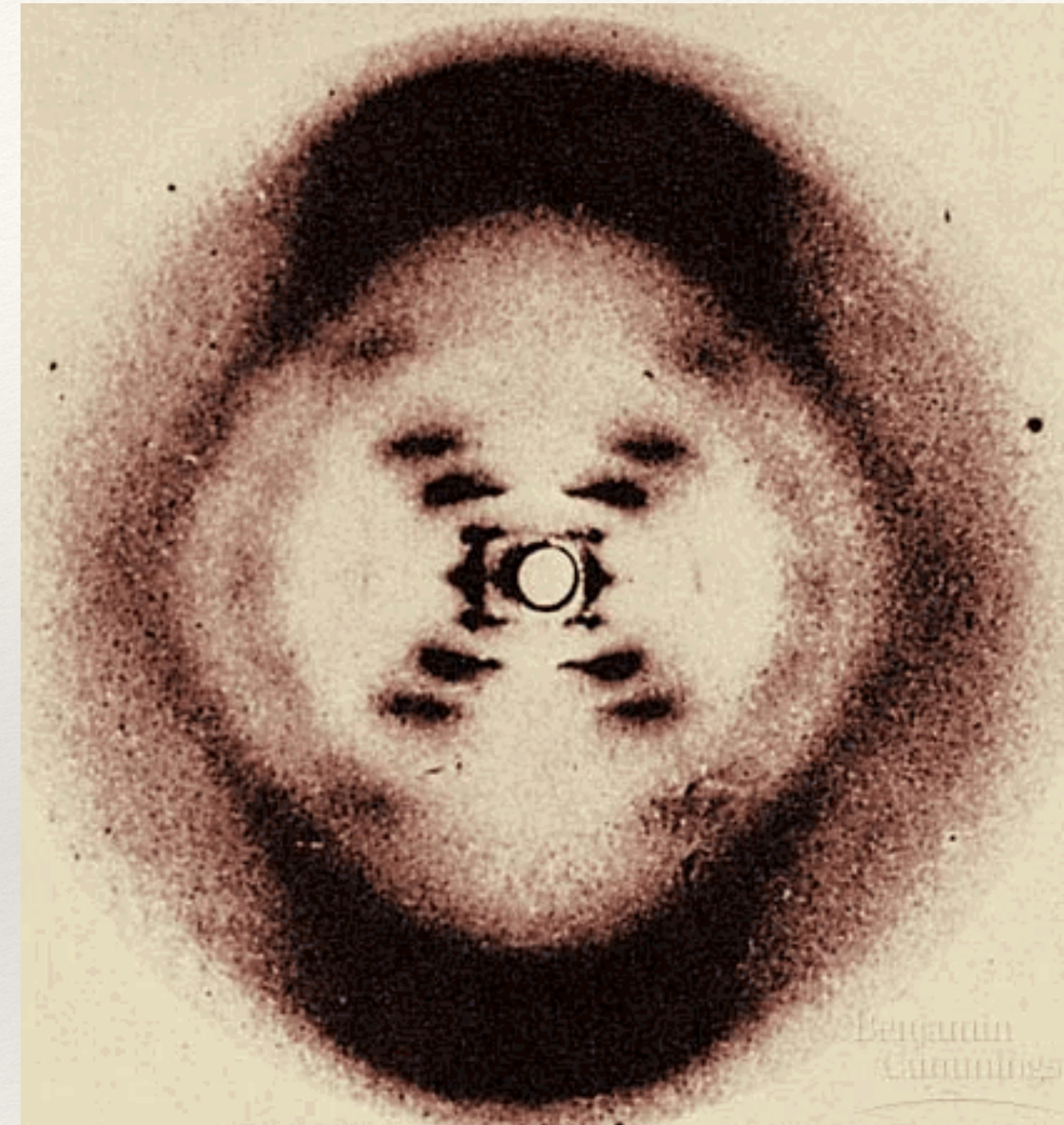
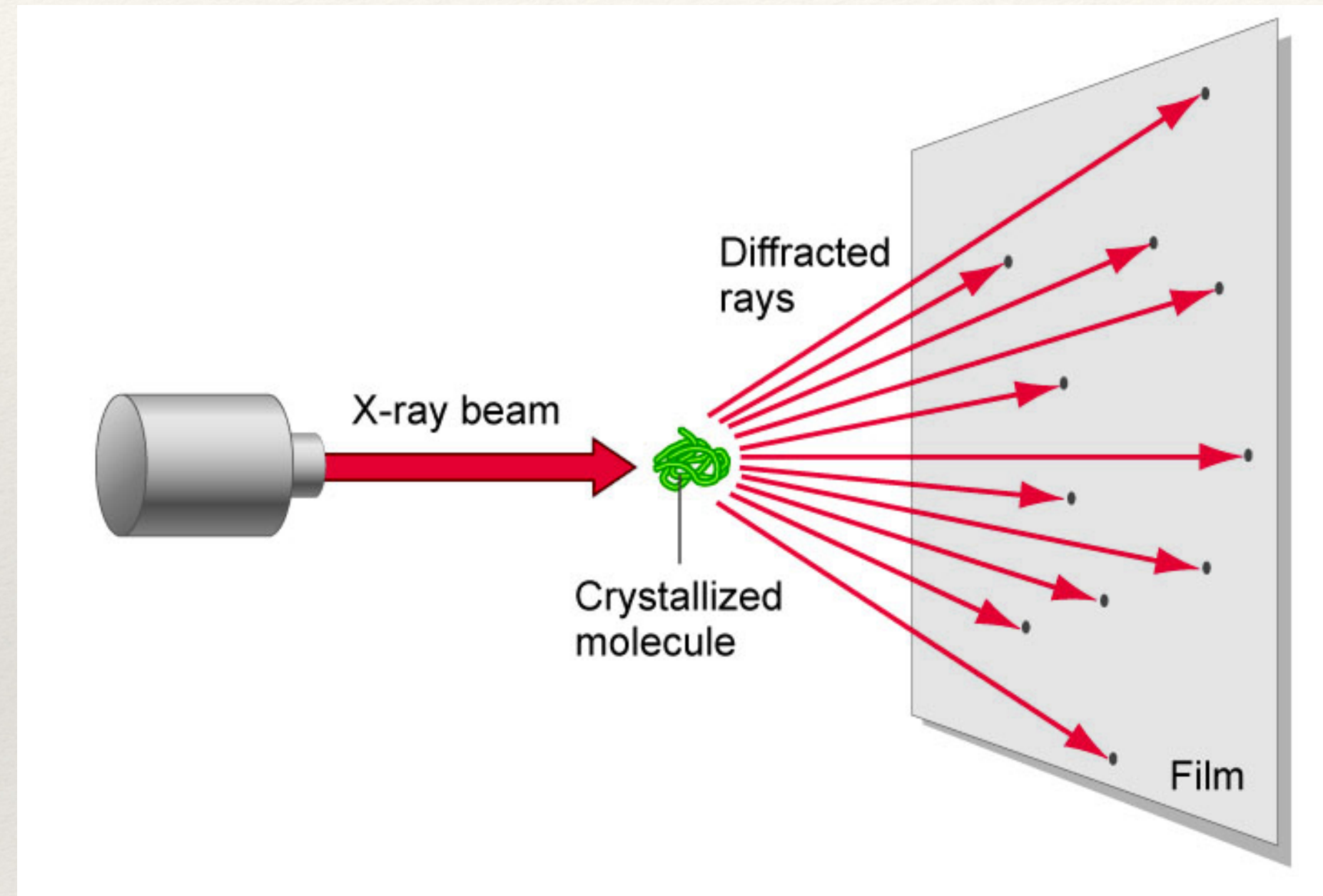
From the "pattern of diffraction", i.e. the maximum of intensities observed, we can find the angles of diffraction and for each angle we get the corresponding d using Bragg's law



General principle of X-ray crystallography applied to proteins:

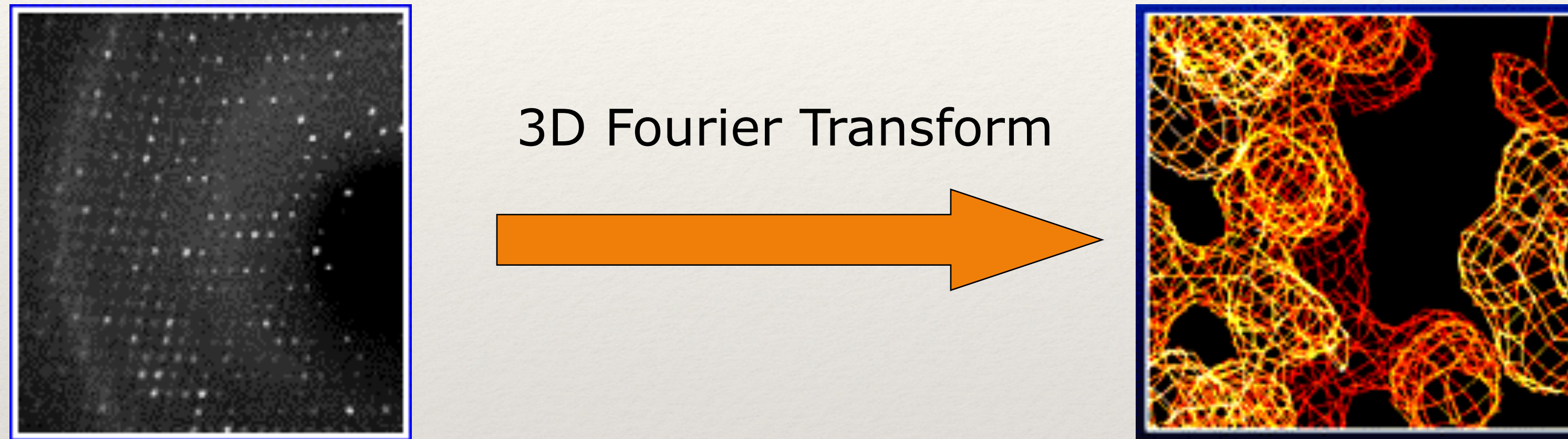
- 1) We need a crystal
- 2) From the diffraction pattern, we get the crystal organization
- 3) From the diffraction intensities, we get the electron densities
- 4) Once the electron density map we fit a structure that matches with this density
- 5) From the atomic model, we can compute a theoretical diffraction map; if it matches with the experimental one, we are done; otherwise refine

Getting the Diffraction Pattern



*Rosalyn Franklin:
Diffraction pattern for DNA*

From Diffraction to Electron Density Map



One hidden problem: diffraction patterns provide intensities; for Fourier transform, need intensity and phase. A significant step in X-ray crystallography is the solve the “phase problem”.

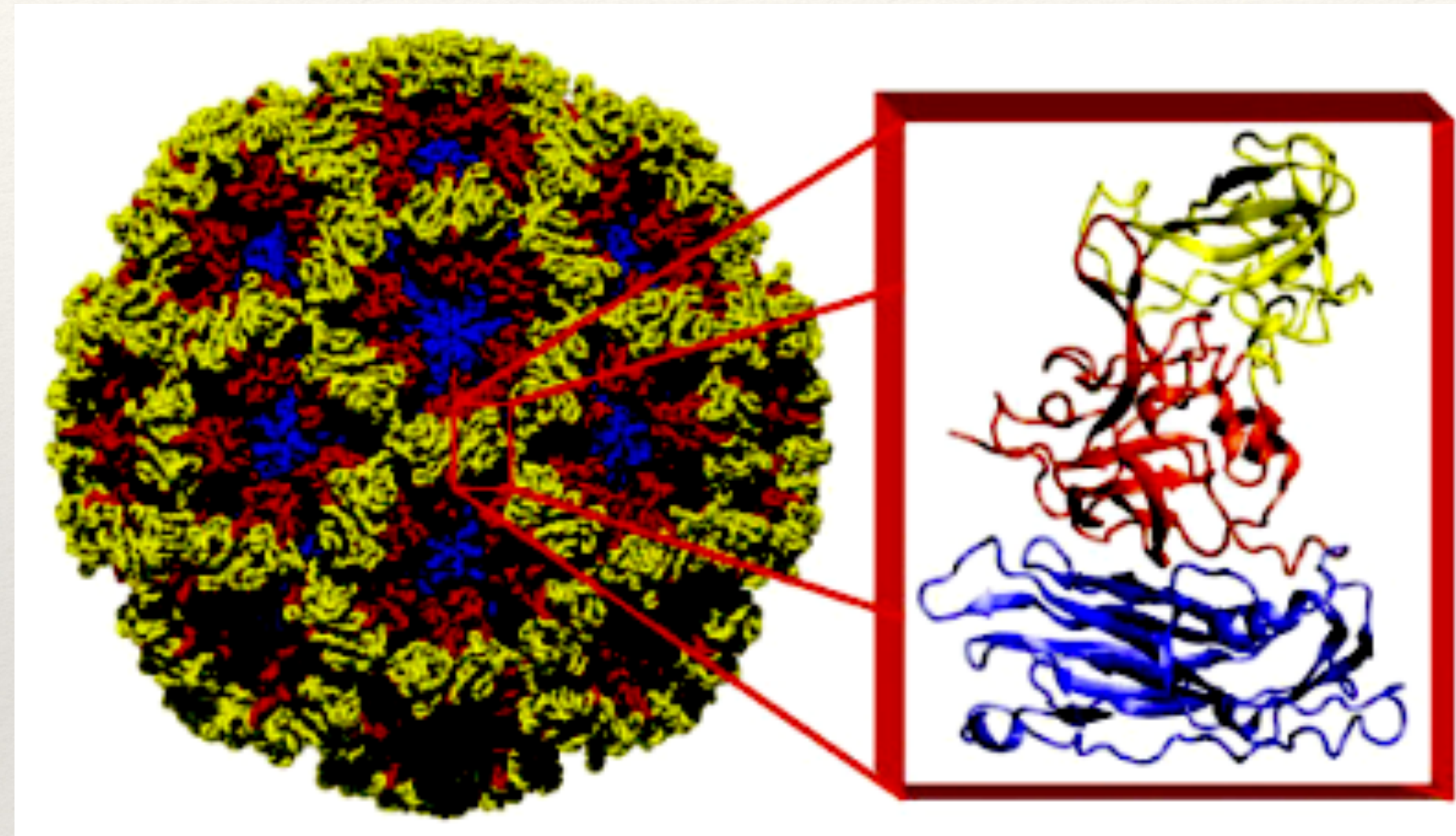
Resolution of X-ray structures

Resolution (Å)	Meaning
>4.0	Individual coordinates meaningless
3.0 - 4.0	Fold possibly correct, but errors are very likely.
2.5 - 3.0	Fold likely correct except that some surface loops might be mismodelled.
2.0 - 2.5	Many small errors can normally be detected. Fold normally correct and number of errors in surface loops is small.
1.5 - 2.0	Water molecules and small ligands become visible. Many small errors can normally be detected. Folds are extremely rarely incorrect, even in surface loops.
0.5 - 1.5	In general, structures have almost no errors at this resolution. geometry studies are made from these structures.

Large molecular assemblies: X-ray crystallography and Cryo-EM

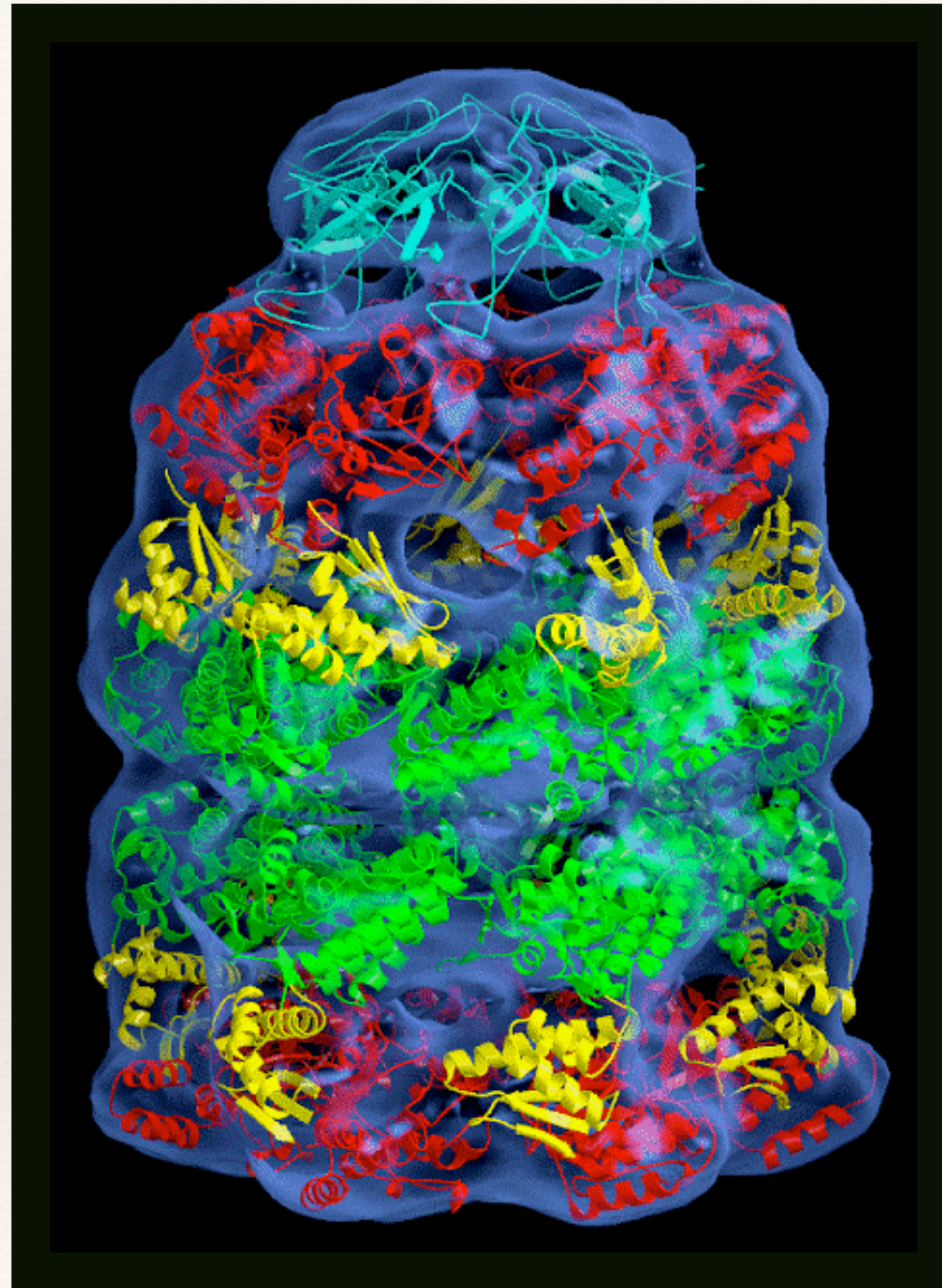
X-ray structure

(180 copies of
the same protein)



(Norwalk virus: <http://www.bcm.edu/molvir/norovirus>)

Large molecular assemblies: X-ray crystallography and Cryo-EM



Cryo-EM:

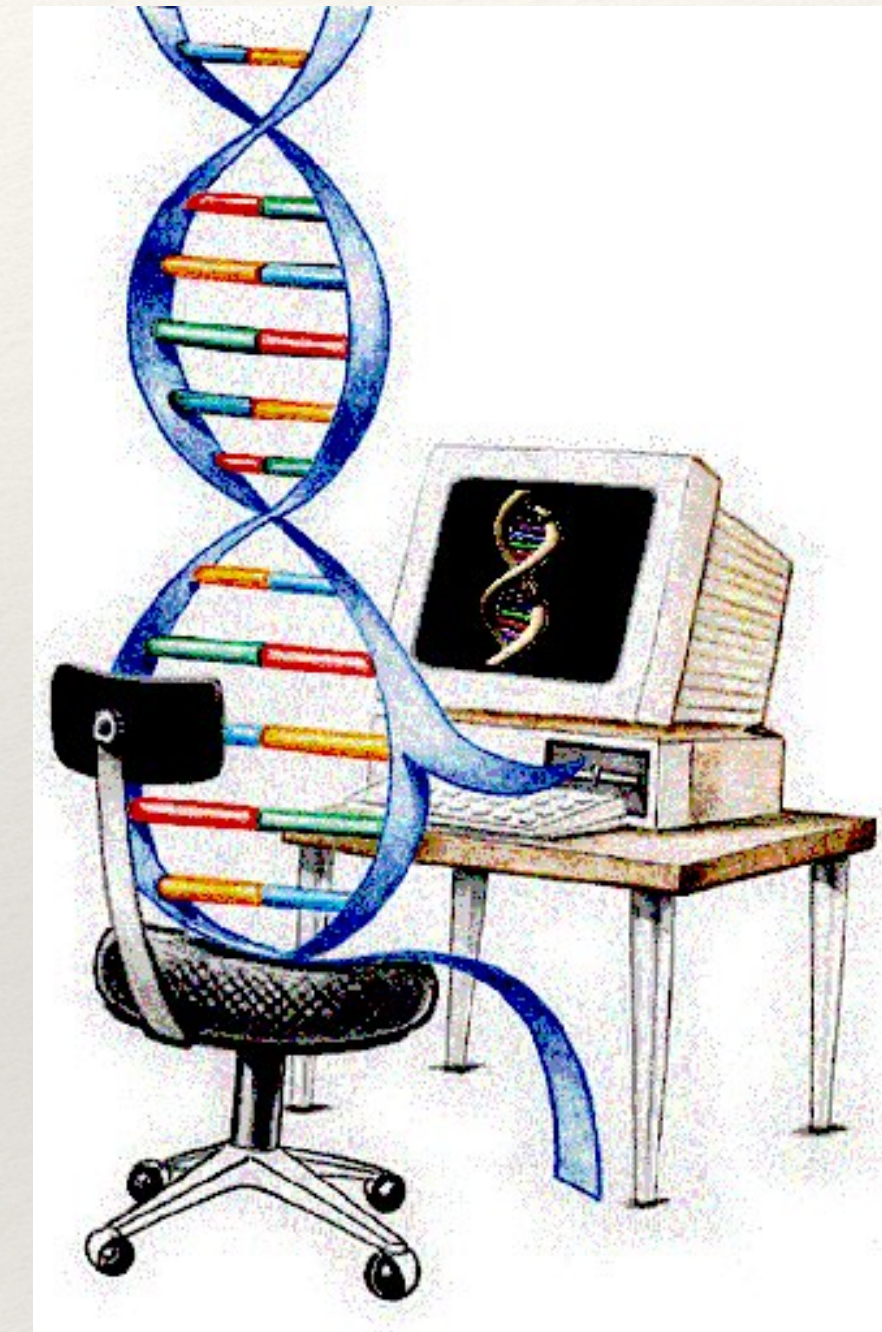
- Microscopy technique; as such, do not need crystal (closer to physiological conditions)
- Not high-resolution enough to provide atomic details; used in combination with modeling

Structural Bioinformatics: Proteins

Proteins: Sources of Structure Information

Proteins: Homology Modeling

Proteins: Secondary structure prediction

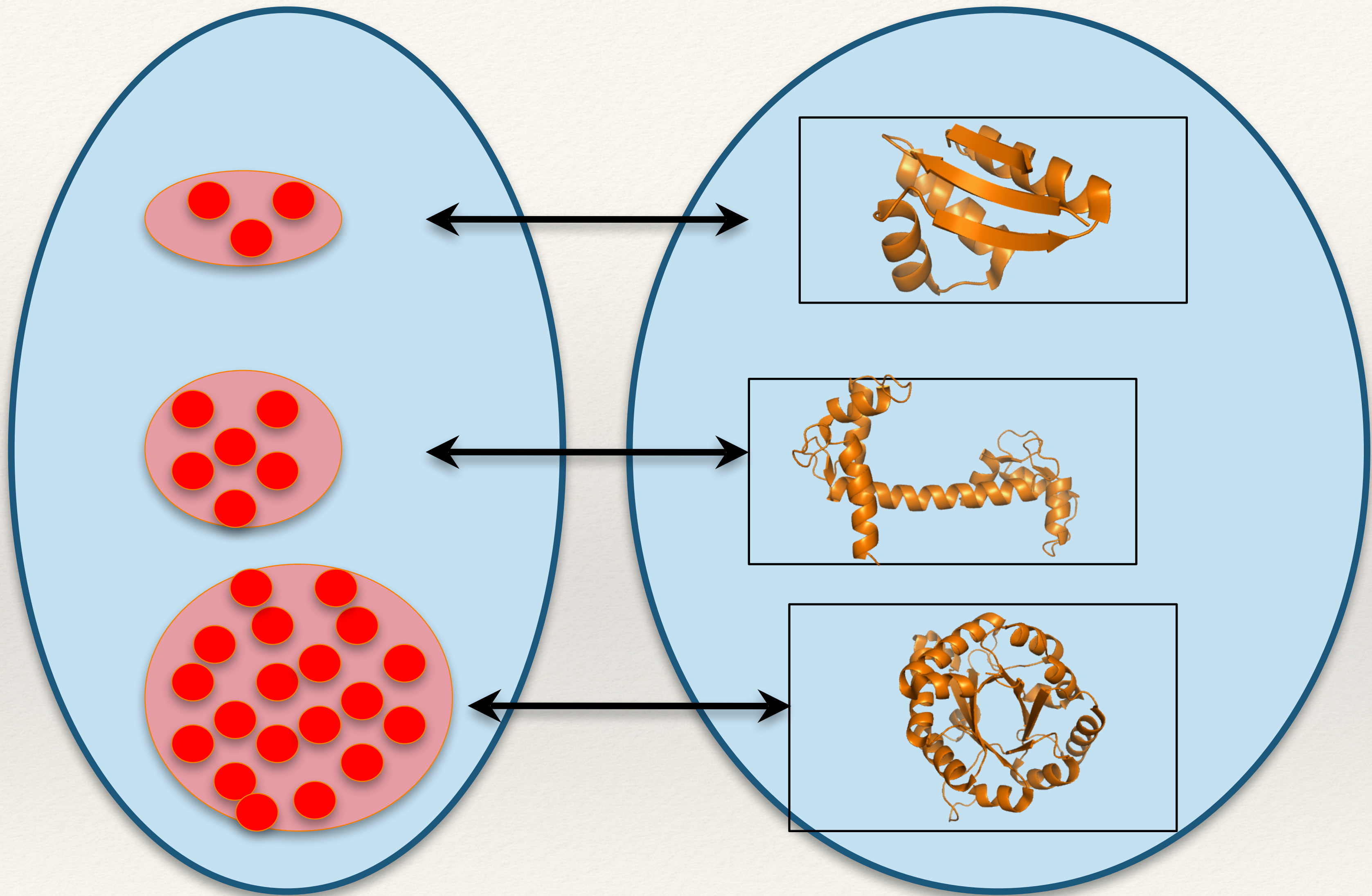


Why we need Homology Modeling?

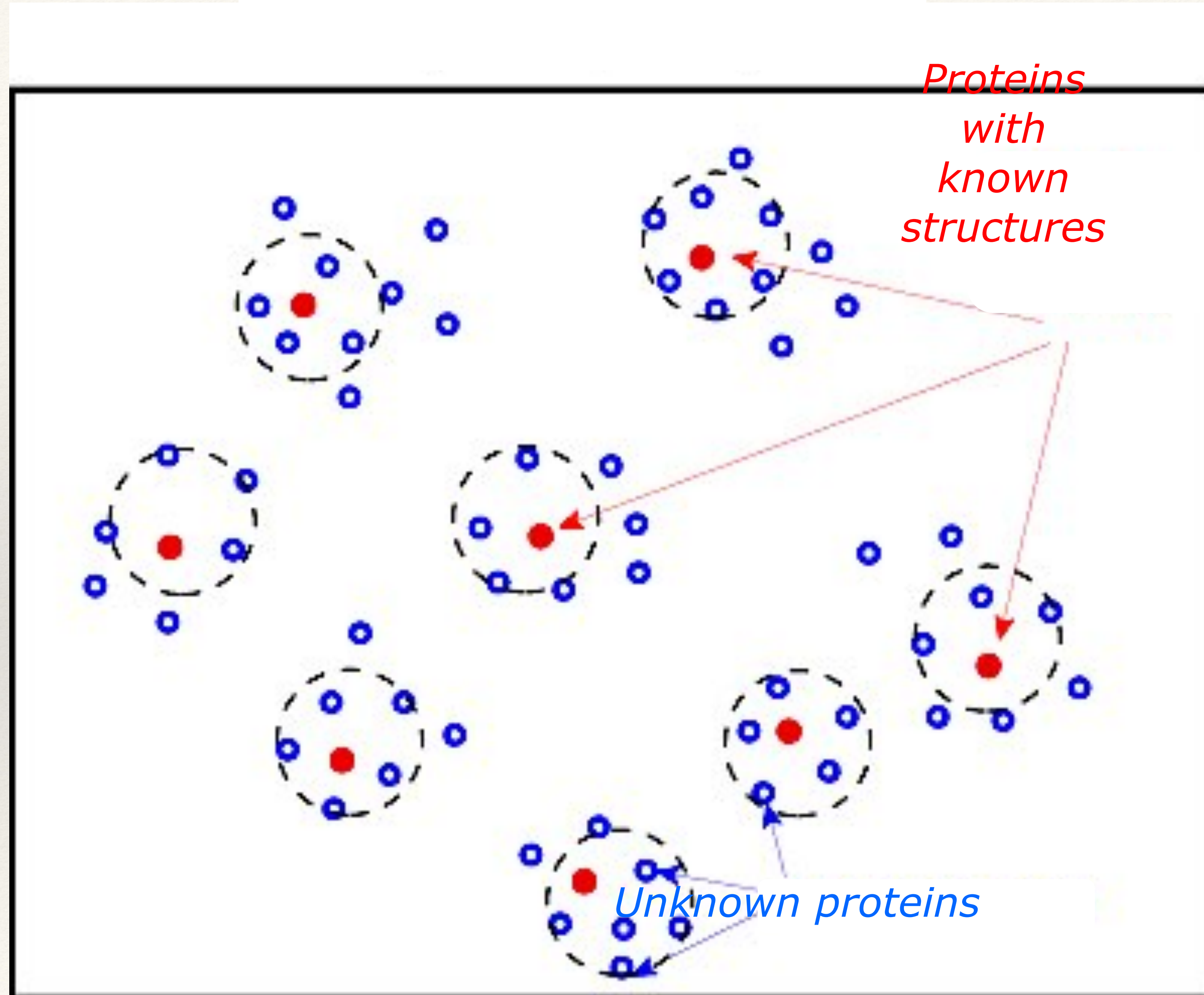
- ❖ Aim to solve the structure of all proteins: this is too much work experimentally!
- ❖ Solve enough structures so that the remaining structures can be inferred from those experimental structures
- ❖ The number of experimental structures needed depend on our abilities to generate a model.

Sequence Space

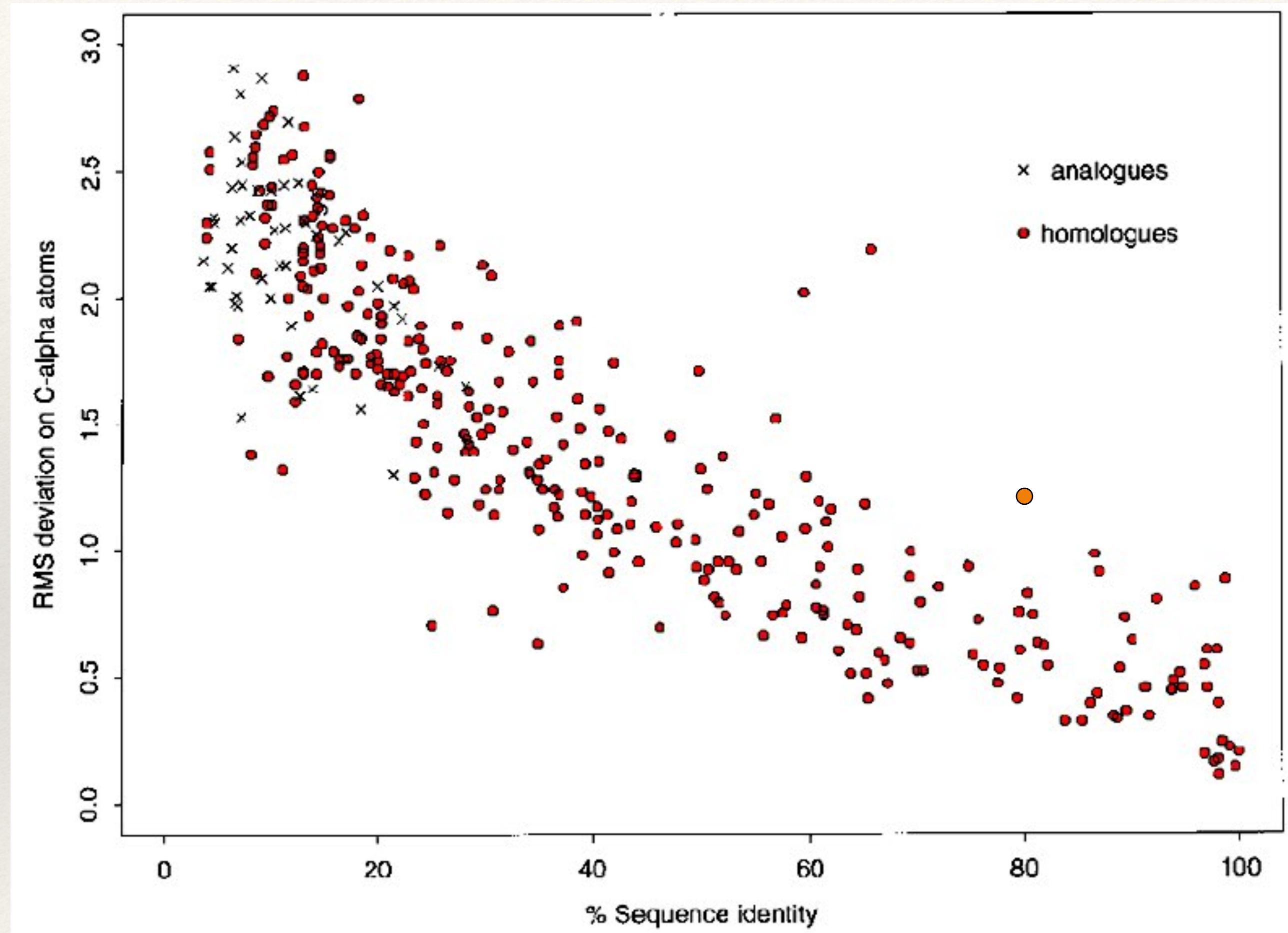
Structure Space



Why we need Homology Modeling?



Why does Homology Modeling Work?



*High sequence
identity*

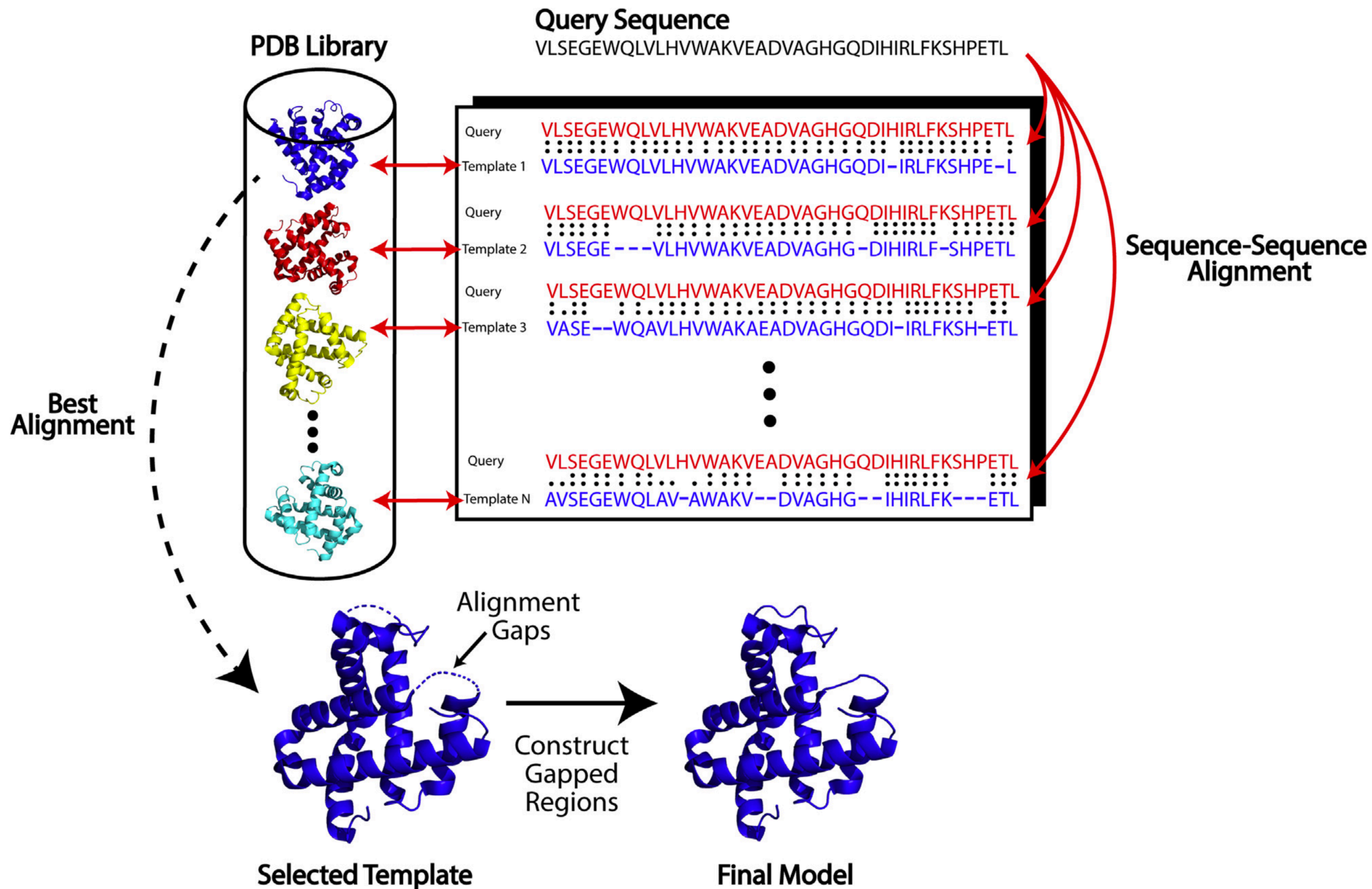


*High structure
similarity*

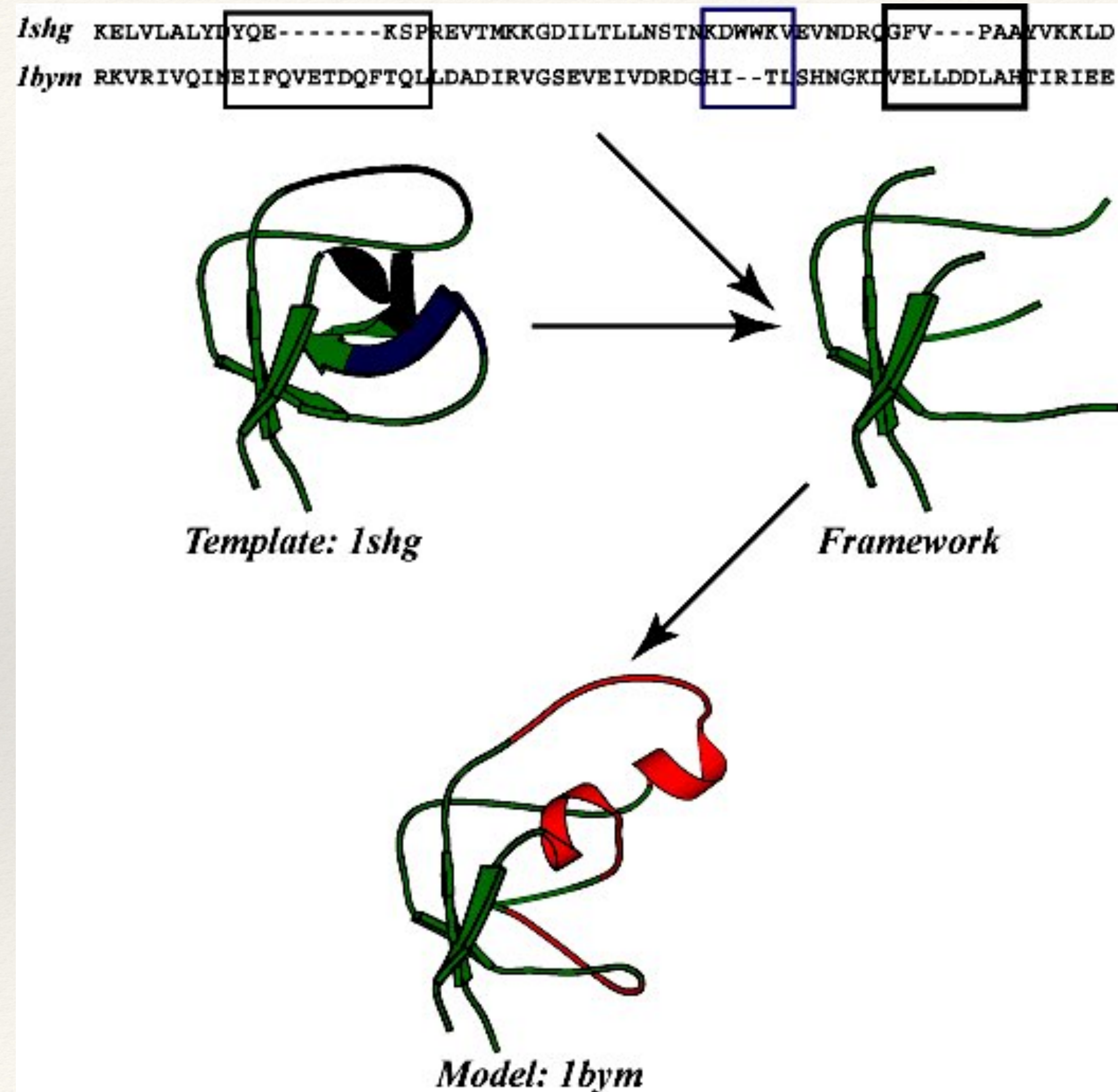
Russell et al. (1997) J Mol Biol 269: 423-439

Homology
Modeling:

How it
works



Homology Modeling: How it works



- Find template
- Align target sequence with template
- Generate model:
 - add loops
 - add sidechains
- Refine model

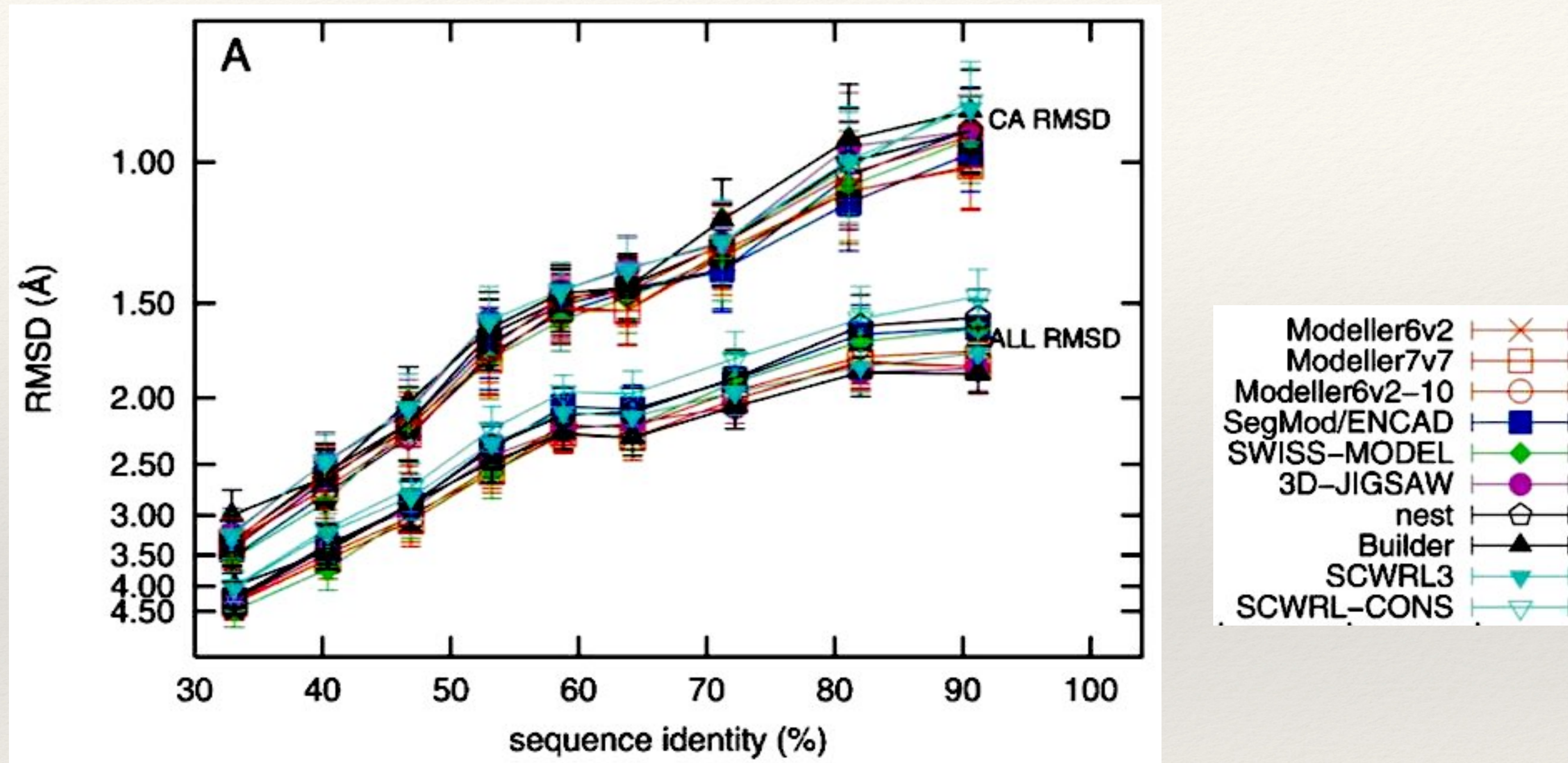
Homology Modeling: Input

The query, also called target sequence

The template structure and sequence
(need high quality structure)

The sequence alignment between query and template sequence
(Probably the most important input!)

Homology Modeling: Which program to use?



Wallner B, Elofsson A.

All are not equal: a benchmark of different homology modeling programs.

Protein Sci. 2005 14(5):1315-27.

Homology Modeling: Which program to use?

1) **Web service: SwissModel**

<http://swissmodel.expasy.org/>

3 modes:

- fully automatic
- "Alignment mode": you provide your own target-template alignment
- "Project mode": provides an environment to edit alignment

2) **Software: Modeller**

<http://www.salilab.org/modeller/>

Probably the best maintained software the homology modeling

Do not start a homology modeling project before checking...

– **Swiss-Model** repository

(<http://swissmodel.expasy.org/repository/>)

- Companion to the Swiss-Model tools – over 2.0 million models of protein domains

– **ModBase**

(<http://modbase.compbio.ucsf.edu/>)

- Companion to **Modeller**

ModBase Contents

Number of Models	38,284,888
Number of Unique Sequences modeled	6,492,640
Unique sequences attempted	8,072,845
Number of PDB chains	85,448

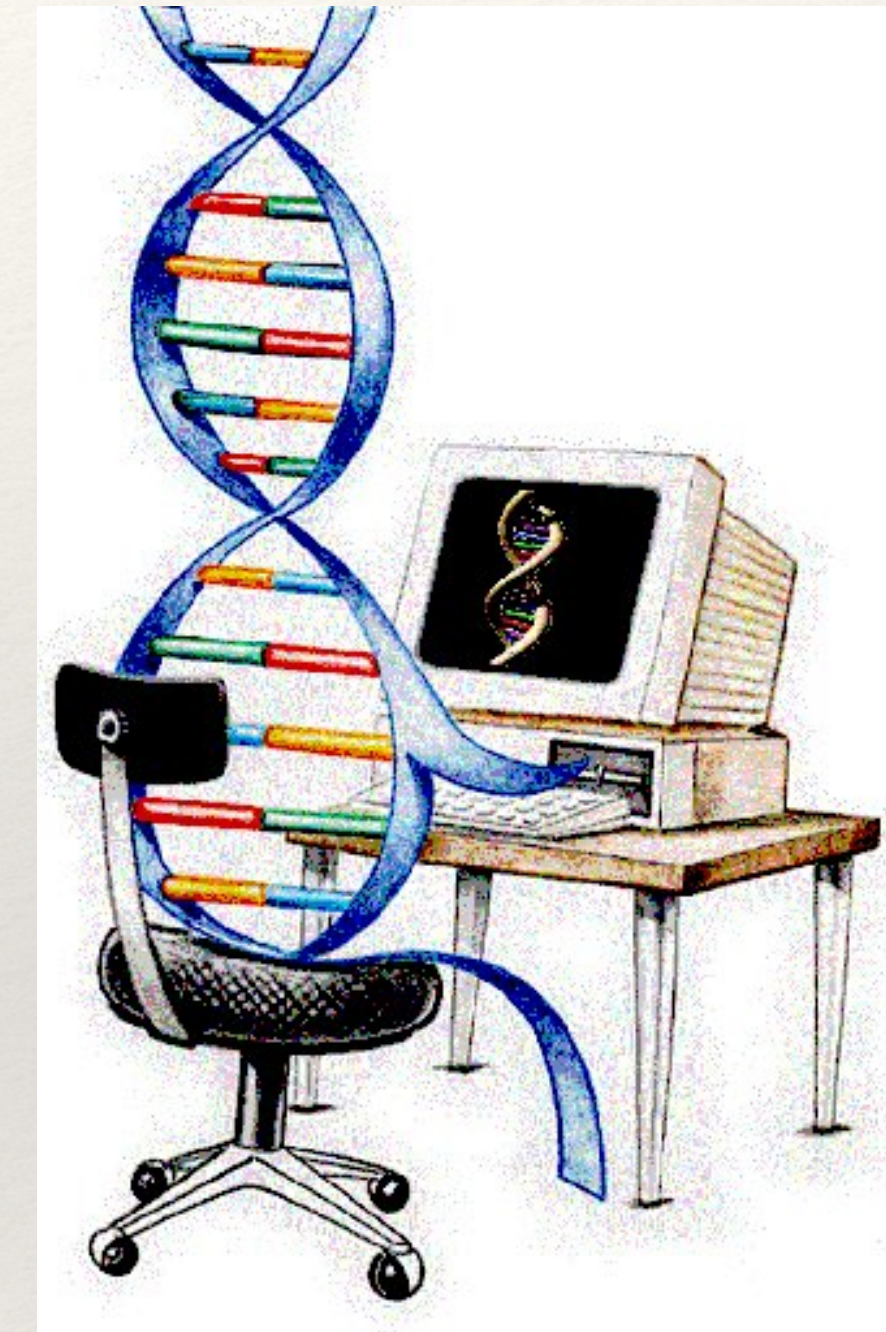
Statistics updated at Fri Jan 14 22:01:34 PST 2022

Structural Bioinformatics: Proteins

Proteins: Sources of Structure Information

Proteins: Homology Modeling

Proteins: Secondary structure prediction



Secondary Structure Prediction

- ❖ *Given a protein sequence $a_1a_2\dots a_N$, secondary structure prediction aims at defining the state of each amino acid a_i as being either H (helix), E (extended=strand), or O (other) (Some methods have 4 states: H, E, T for turns, and O for other).*
- ❖ *The quality of secondary structure prediction is measured with a “3-state accuracy” score, or Q_3 . Q_3 is the percent of residues that match “reality” (X-ray structure).*

Secondary Structure Assignment

Determine Secondary Structure positions in known protein structures using DSSP or STRIDE:

1. Kabsch and Sander. Dictionary of Secondary Structure in Proteins: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymer* 22: 2571-2637 (1983) (DSSP)
2. Frischman and Argos. Knowledge-based secondary structure assignments. *Proteins*, 23:566-571 (1995) (STRIDE)

Early methods for Secondary Structure Prediction

- ❖ *Chou and Fasman*

(Chou and Fasman. Prediction of protein conformation. Biochemistry, 13: 211-245, 1974)

- ❖ *GOR*

(Garnier, Osguthorpe and Robson. Analysis of the accuracy and implications of simple methods for predicting the secondary structure of globular proteins. J. Mol. Biol., 120:97- 120, 1978)

Chou and Fasman

- ❖ *Start by computing amino acids propensities to belong to a given type of secondary structure:*

$$\frac{P(i / \text{Helix})}{P(i)} \quad \frac{P(i / \text{Beta})}{P(i)} \quad \frac{P(i / \text{Turn})}{P(i)}$$

Propensities > 1 mean that the residue type i is likely to be found in the corresponding secondary structure type.

Chou and Fasman

Amino Acid	α -Helix	β -Sheet	Turn	
Ala	1.29	0.90	0.78	Favors α -Helix
Cys	1.11	0.74	0.80	
Leu	1.30	1.02	0.59	
Met	1.47	0.97	0.39	
Glu	1.44	0.75	1.00	
Gln	1.27	0.80	0.97	
His	1.22	1.08	0.69	
Lys	1.23	0.77	0.96	
Val	0.91	1.49	0.47	
Ile	0.97	1.45	0.51	Favors β -strand
Phe	1.07	1.32	0.58	
Tyr	0.72	1.25	1.05	
Trp	0.99	1.14	0.75	
Thr	0.82	1.21	1.03	
Gly	0.56	0.92	1.64	
Ser	0.82	0.95	1.33	Favors turn
Asp	1.04	0.72	1.41	
Asn	0.90	0.76	1.23	
Pro	0.52	0.64	1.91	
Arg	0.96	0.99	0.88	

Chou and Fasman

Predicting helices:

- find nucleation site: 4 out of 6 contiguous residues with $P(\alpha) > 1$
- extension: extend helix in both directions until a set of 4 contiguous residues has an average $P(\alpha) < 1$ (breaker)
- if average $P(\alpha)$ over whole region is > 1 , it is predicted to be helical

Predicting strands:

- find nucleation site: 3 out of 5 contiguous residues with $P(\beta) > 1$
- extension: extend strand in both directions until a set of 4 contiguous residues has an average $P(\beta) < 1$ (breaker)
- if average $P(\beta)$ over whole region is > 1 , it is predicted to be a strand

Chou and Fasman

Position-specific parameters for turn:

Each position has distinct amino acid preferences.

Examples:

-At position 2, Pro is highly preferred; Trp is disfavored

-At position 3, Asp, Asn and Gly are preferred

-At position 4, Trp, Gly and Cys preferred

f(i)	f(i+1)	f(i+2)	f(i+3)	
	i	i+1	i+2	i+3
Ala	0.060	0.076	0.035	0.058
Arg	0.070	0.106	0.099	0.085
Asp	0.147	0.110	0.179	0.081
Asn	0.161	0.083	0.191	0.091
Cys	0.149	0.050	0.117	0.128
Glu	0.056	0.060	0.077	0.064
Gln	0.074	0.098	0.037	0.098
Gly	0.102	0.085	0.190	0.152
His	0.140	0.047	0.093	0.054
Ile	0.043	0.034	0.013	0.056
Leu	0.061	0.025	0.036	0.070
Lys	0.055	0.115	0.072	0.095
Met	0.068	0.082	0.014	0.055
Phe	0.059	0.041	0.065	0.065
Pro	0.102	0.301	0.034	0.068
Ser	0.120	0.139	0.125	0.106
Thr	0.086	0.108	0.065	0.079
Trp	0.077	0.013	0.064	0.167
Tyr	0.082	0.065	0.114	0.125
Val	0.062	0.048	0.028	0.053

Chou and Fasman

Predicting turns:

- for each tetrapeptide starting at residue i , compute:
 - P_{Turn} (average propensity over all 4 residues)
 - $F = f(i) \cdot f(i+1) \cdot f(i+2) \cdot f(i+3)$
- if $P_{\text{Turn}} > P_{\alpha}$ and $P_{\text{Turn}} > P_{\beta}$ and $P_{\text{Turn}} > 1$ and $F > 0.000075$ tetrapeptide is considered a turn.

Chou and Fasman prediction:

http://fasta.bioch.virginia.edu/fasta_www/chofas.htm

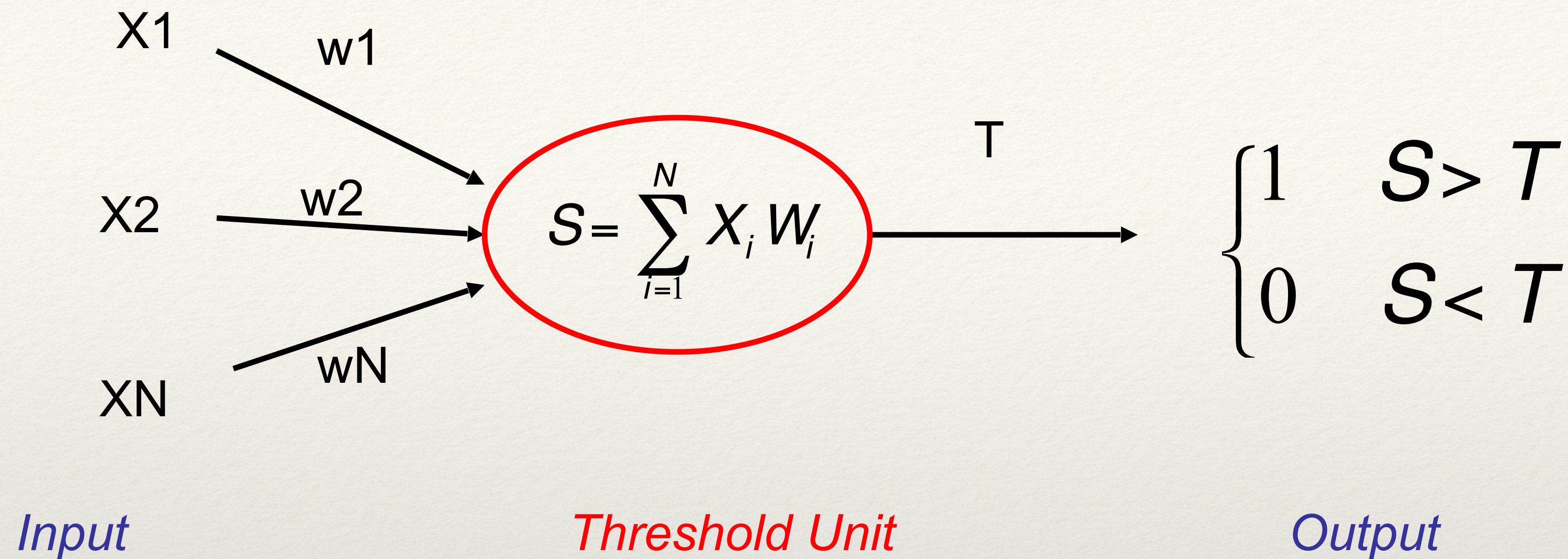
Neural Networks

The most successful methods for predicting secondary structure are based on neural networks. The overall idea is that neural networks can be trained to recognize amino acid patterns in known secondary structure units, and to use these patterns to distinguish between the different types of secondary structure.

Neural networks classify “input vectors” or “examples” into categories (2 or more).

They are loosely based on biological neurons.

The perceptron



The **perceptron** classifies the input vector X into two categories.

If the weights and threshold T are not known in advance, the perceptron must be **trained**. Ideally, the perceptron must be trained to return the correct answer on all training examples, and perform well on examples it has never seen.

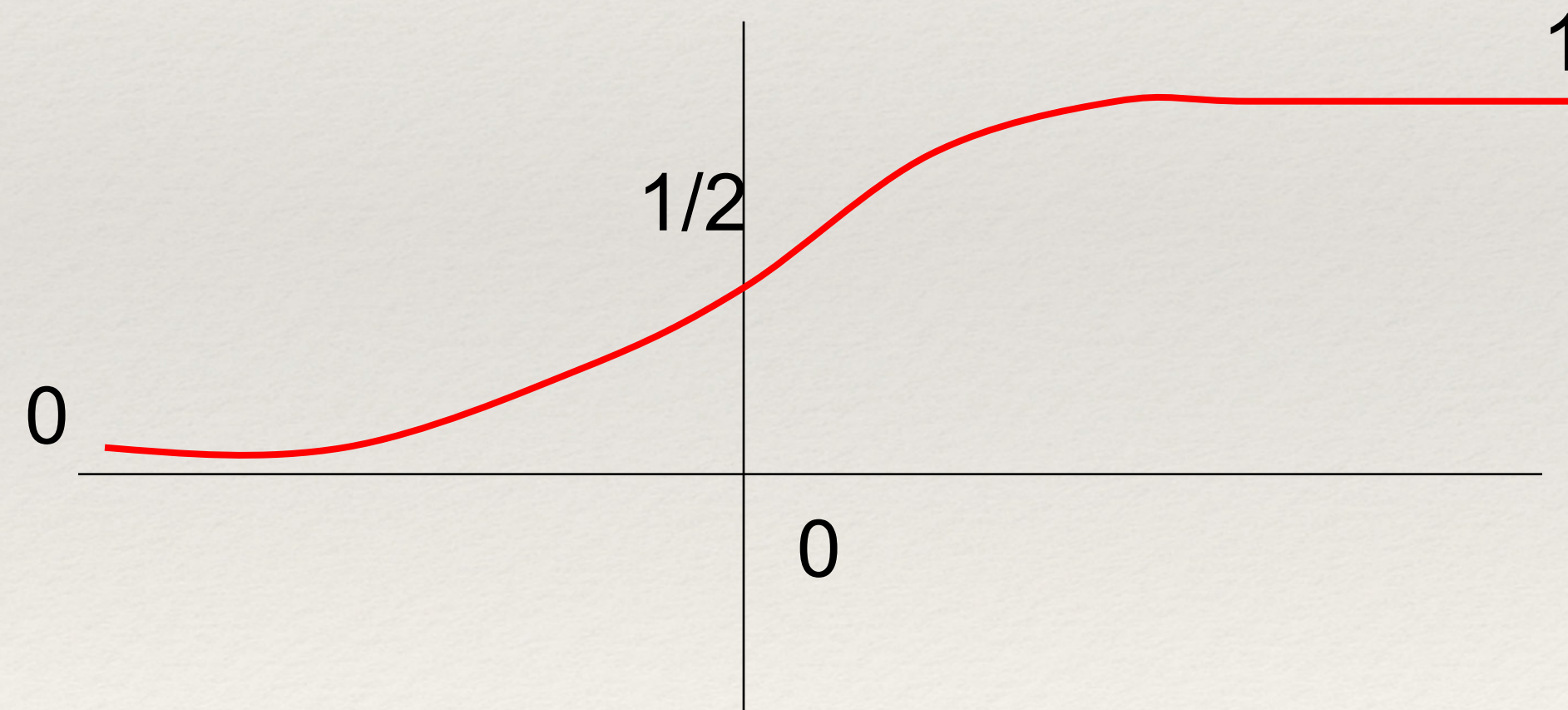
The training set must contain both type of data (i.e. with "1" and "0" output).

The perceptron

Notes:

- The input is a vector X and the weights can be stored in another vector W .
- the perceptron computes the dot product $S = X.W$
- the output F is a function of S : it is often set discrete (i.e. 1 or 0), in which case the function is the step function.
For continuous output, often use a sigmoid:

$$F(X) = \frac{1}{1 + e^{-X}}$$



- Not all perceptrons can be trained ! (famous example: XOR)

The perceptron

Training a perceptron:

Find the weights W that minimizes the error function:

$$E = \sum_{i=1}^P \left(F(X^i \cdot W) - t(X^i) \right)^2$$

P : number of training data

X^i : training vectors

$F(W \cdot X^i)$: output of the perceptron

$t(X^i)$: target value for X^i

Use steepest descent:

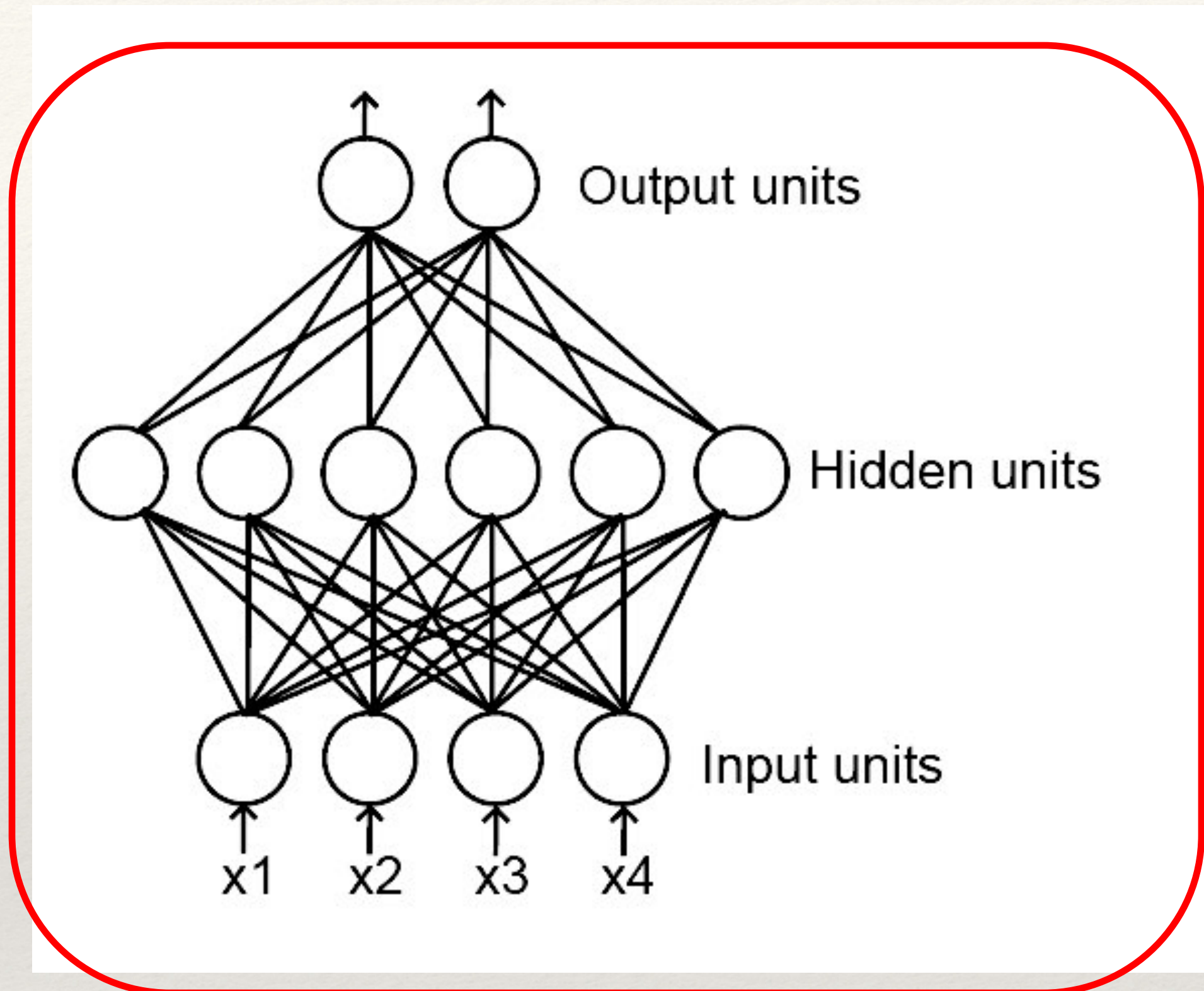
- compute gradient:
- update weight vector:
- iterate

$$\nabla E = \left(\frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \frac{\partial E}{\partial w_3}, \dots, \frac{\partial E}{\partial w_N} \right)$$

$$W_{new} = W_{old} - \varepsilon \nabla E$$

(ε : learning rate)

Neural Network



A complete neural network is a set of perceptrons interconnected such that the outputs of some units becomes the inputs of other units. Many topologies are possible!

Neural networks are trained just like perceptron, by minimizing an error function:

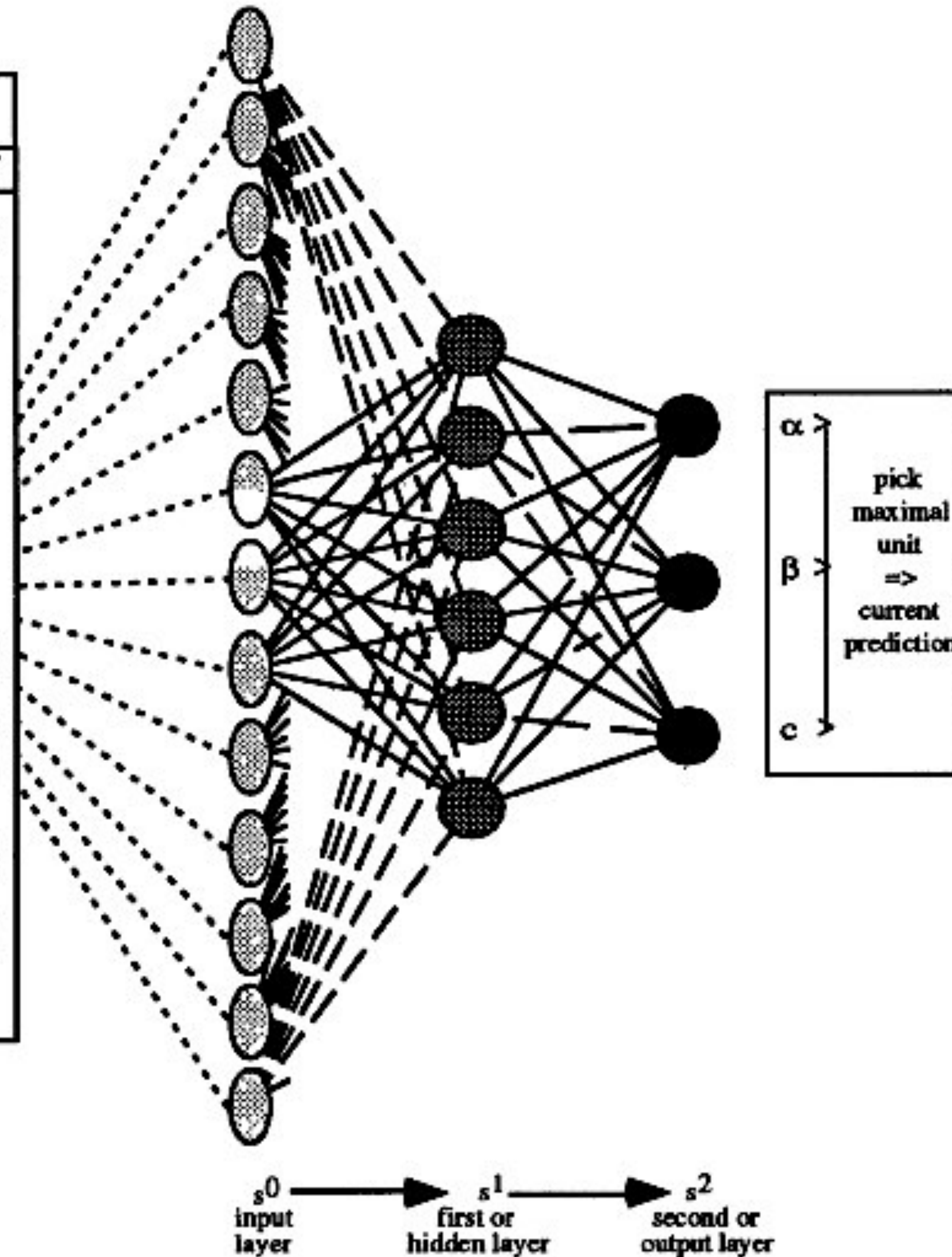
$$E = \sum_{i=1}^{N_{data}} \left(NN(X^i) - t(X^i) \right)^2$$

PHD: Secondary structure prediction using NN

Biophysics: Rost and Sander

Proc. Natl. Acad. Sci. USA 90 (1993) 7559

Protein	Alignments	profile table
		GSAPD NT EKQ C VH IR LM YFW
:	:: :: :: :	
G	GG GG	5
Y	YY YY 5 . .
I	II EE 2 3
Y	YY YY 5 . .
D	DD DD 5
P	PP PP 5
E	AE AA	. . 3 2
D	VVEE 1 . . 2 2
G	GG GG	5
D	DD DD 5
P	PP PP 5
D	DT DD 4 . . 1
D	NQ NN 1 3 1
G	GNGG	4 1
V	VI VV 4 . 1
N	EP KK 1 . . 1 . 1 2
P	PP PP 5
G	GG GG	5
T	TT TT 5
D	EK SA	. 1 1 . 1 1 1
F	FF FF 5
:	:: :: :: :	



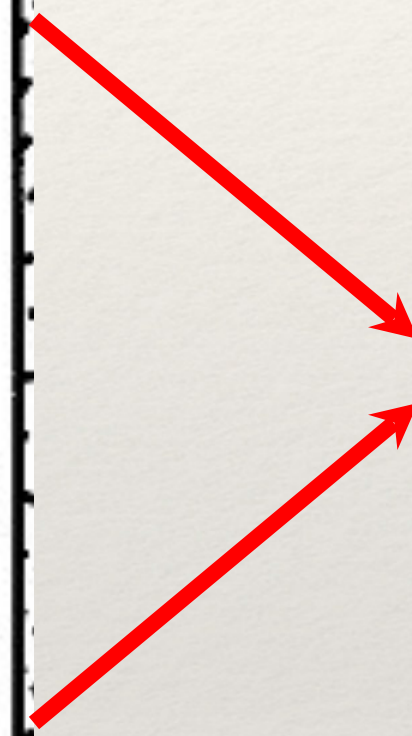
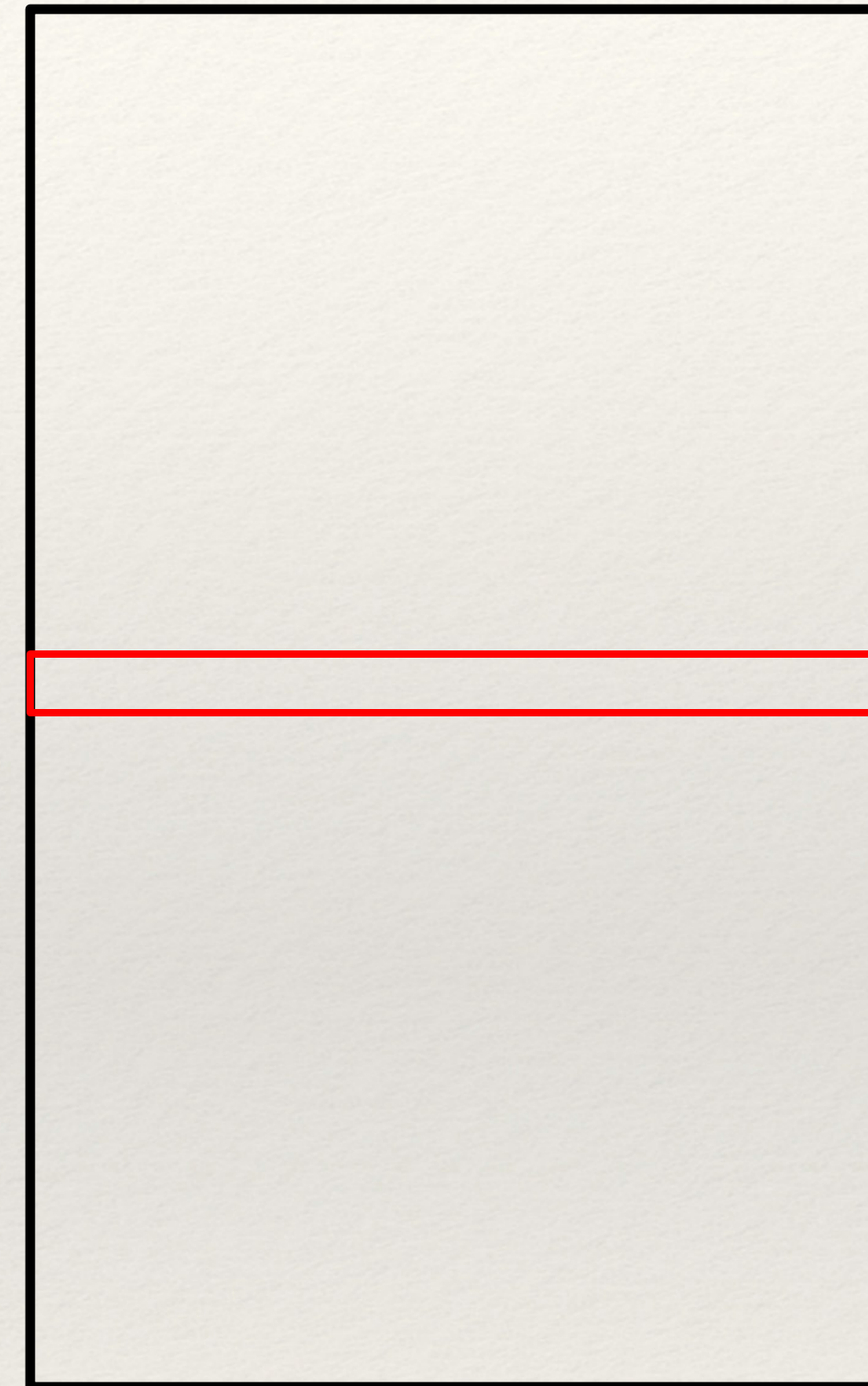
PHD: Input

For each residue, consider a window of size 13:

13x20=260 values

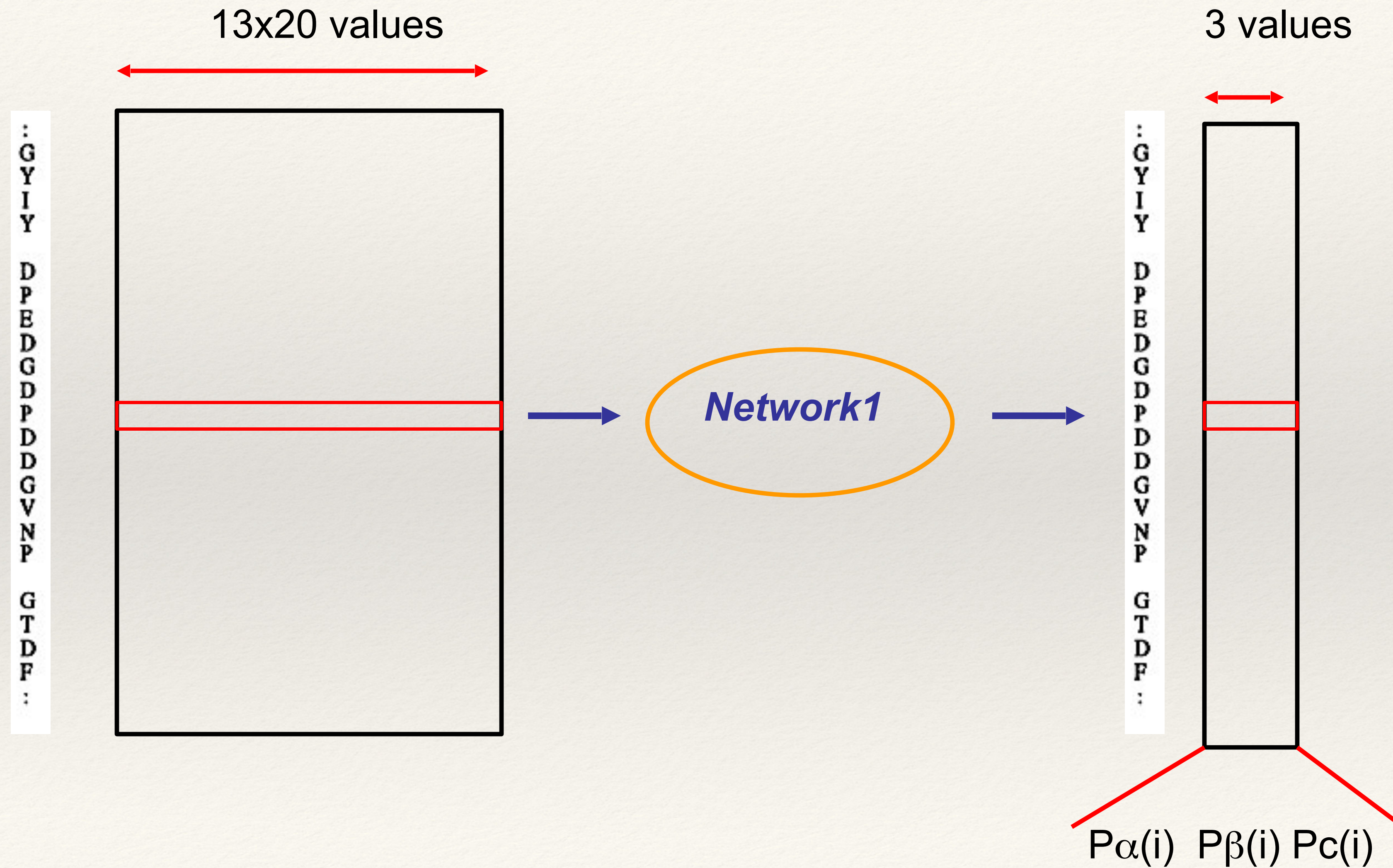
Protein	Alignments	profile table
:	:: :: :	GSAPD NT EKQ C VH IR LM YFW
G	GG GG	5.....
Y	YY YY5..
I	II EE2..3..
Y	YY YY5..
D	DD DD5
P	PP PP	...5.
E	AE AA	..3..2..
D	VV EE	...1..2..2..
G	GG GG	5.....
D	DD DD5
P	PP PP	...5.
D	DT DD	...4.1..
D	NQ NN	...13..1
G	GN GG	4...1..
V	VI VV4.1.
N	EP KK	..1.1.12.
P	PP PP	...5.
G	GG GG	5.....
T	TT TT5..
D	EK SA	.11.1..11.
F	FF FF5.
:	:: :: :	

: G Y I Y D P E D D G G D D P D D D G V N P G T D F :



PHD: Network 1

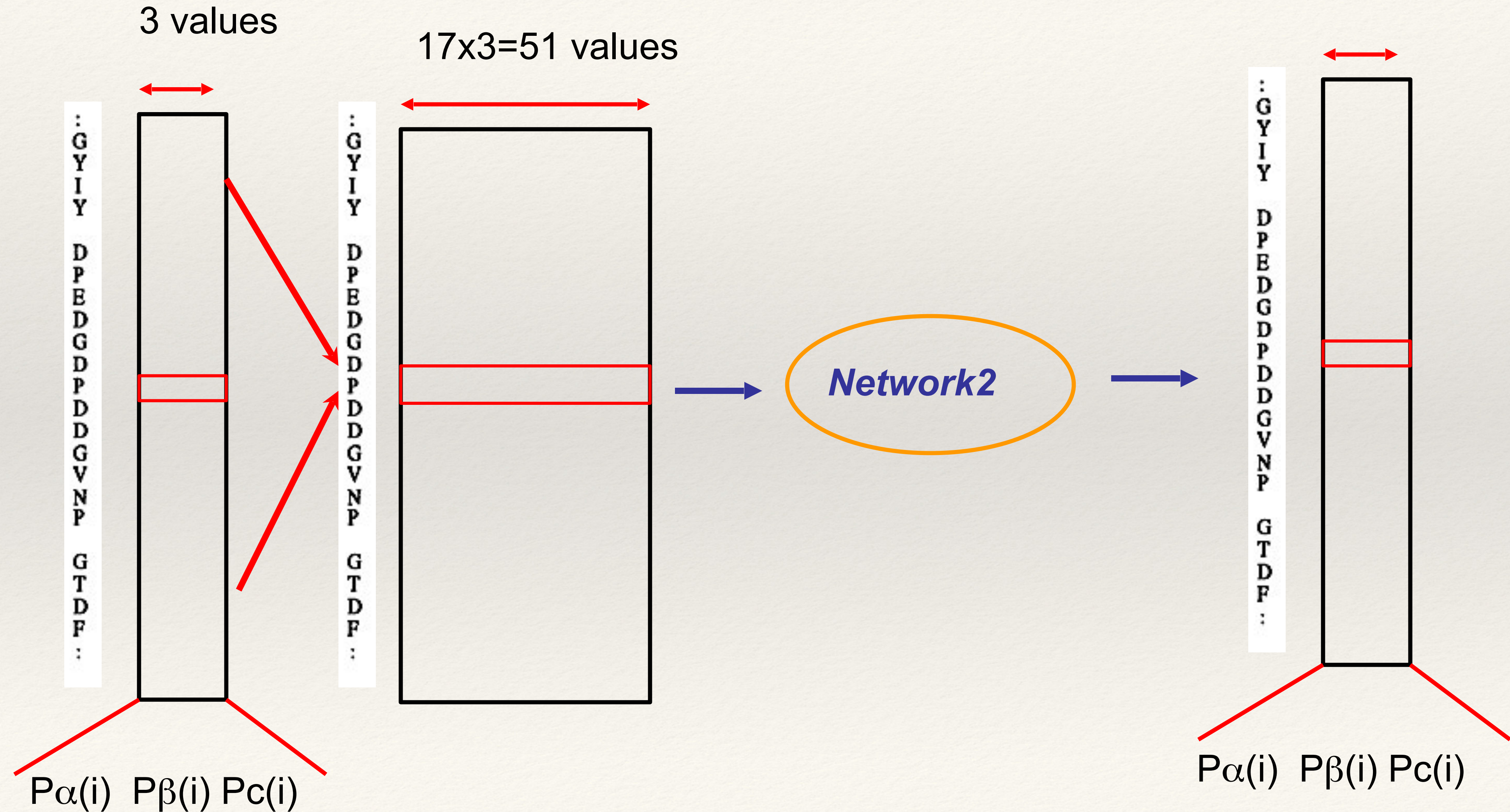
Sequence → Structure



PHD: Network 2

Structure \rightarrow Structure

For each residue, consider a window of size 17:



Protein Structure Prediction

One popular model for protein folding assumes a sequence of events:

- ❖ Hydrophobic collapse
- ❖ Local interactions stabilize secondary structures
- ❖ Secondary structures interact to form motifs
- ❖ Motifs aggregate to form tertiary structure

Protein Structure Prediction

A physics-based approach:

- ❖ - find conformation of protein corresponding to a thermodynamics minimum (free energy minimum)
- ❖ - cannot minimize internal energy alone! Needs to include solvent
- ❖ - simulate folding...a very long process!
- ❖ Folding time are in the ms to second time range; however, Folding simulations at best run 1 ns in one day...

PHD: Secondary structure prediction using NN

- **Sequence-Structure network**: for each amino acid a_j , a window of 13 residues $a_{j-6} \dots a_j \dots a_{j+6}$ is considered. The corresponding rows of the sequence profile are fed into the neural network, and the output is 3 probabilities for a_j : $P(a_j, \alpha)$, $P(a_j, \beta)$ and $P(a_j, \text{other})$
- **Structure-Structure network**: For each a_j , PHD considers now a window of 17 residues; the probabilities $P(a_k, \alpha)$, $P(a_k, \beta)$ and $P(a_k, \text{other})$ for k in $[j-8, j+8]$ are fed into the second layer neural network, which again produces probabilities that residue a_j is in each of the 3 possible conformations
- **Jury system**: PHD has trained several neural networks with different training sets; all neural networks are applied to the test sequence, and results are averaged
- **Prediction**: For each position, the secondary structure with the highest average score is output as the prediction