

10/11/2024

Bioinformatics

↳ biological information → used to extract useful inferences

Biological Data/Information

→ use of comp. sci concepts

1D Data

genome sequences
(A,T,G,C)

Structural Data
3D Data

information about structure
(coordinates of macromolecules)

DBBS

Program-ming

3D Data:

generated from
some experiments

Proteins → structural → scaffolds that make
functional up organisms

Polymer of
amino acid
residues

enzymes

↳ should be folded
into a 3D structure → should perform/exhibit
some biological function

atoms in a protein: C, O, N, S, H

heavy atoms

light atom

coordinates
in PDB

coordinates
not in PDB

structures solved by X-ray Crystallography

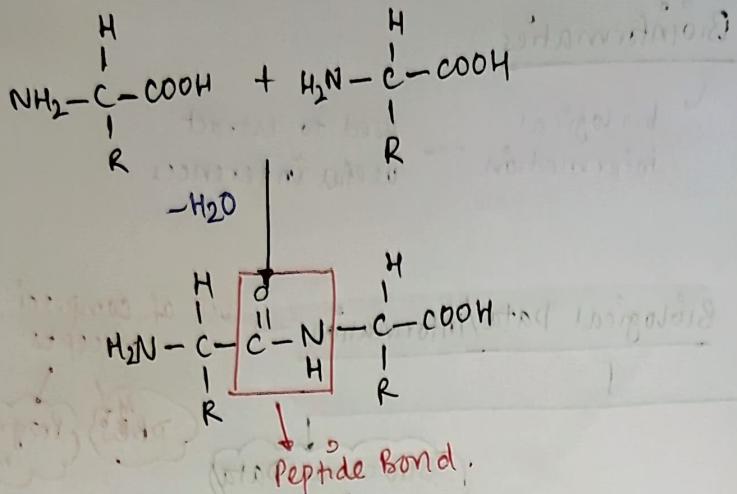
H does not
diffract X-ray → as only one

molecule
in solⁿ state

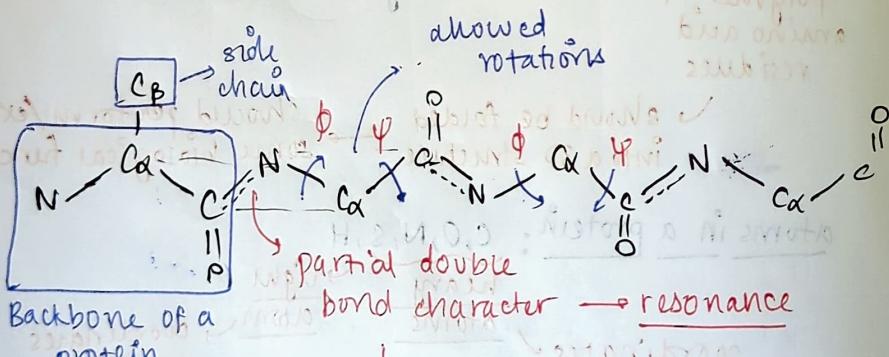
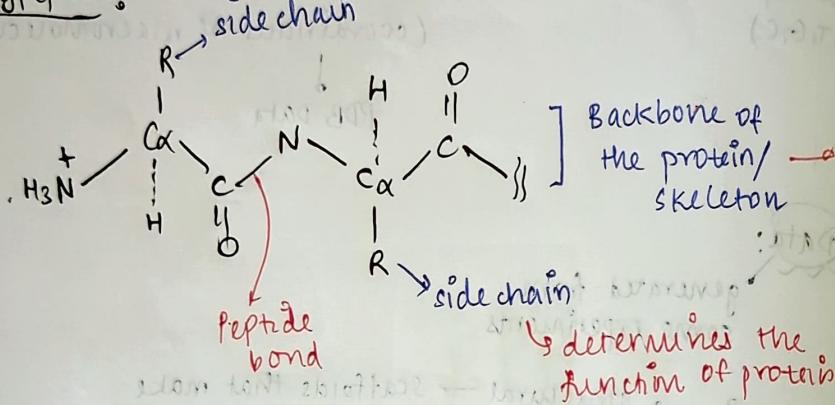
NMR

we get H
coordinates

we cannot
solve the
phase problem



Polypeptide:



more energy than a single bond

restricted rotation

found in protein structures

preferred!

cis conformation

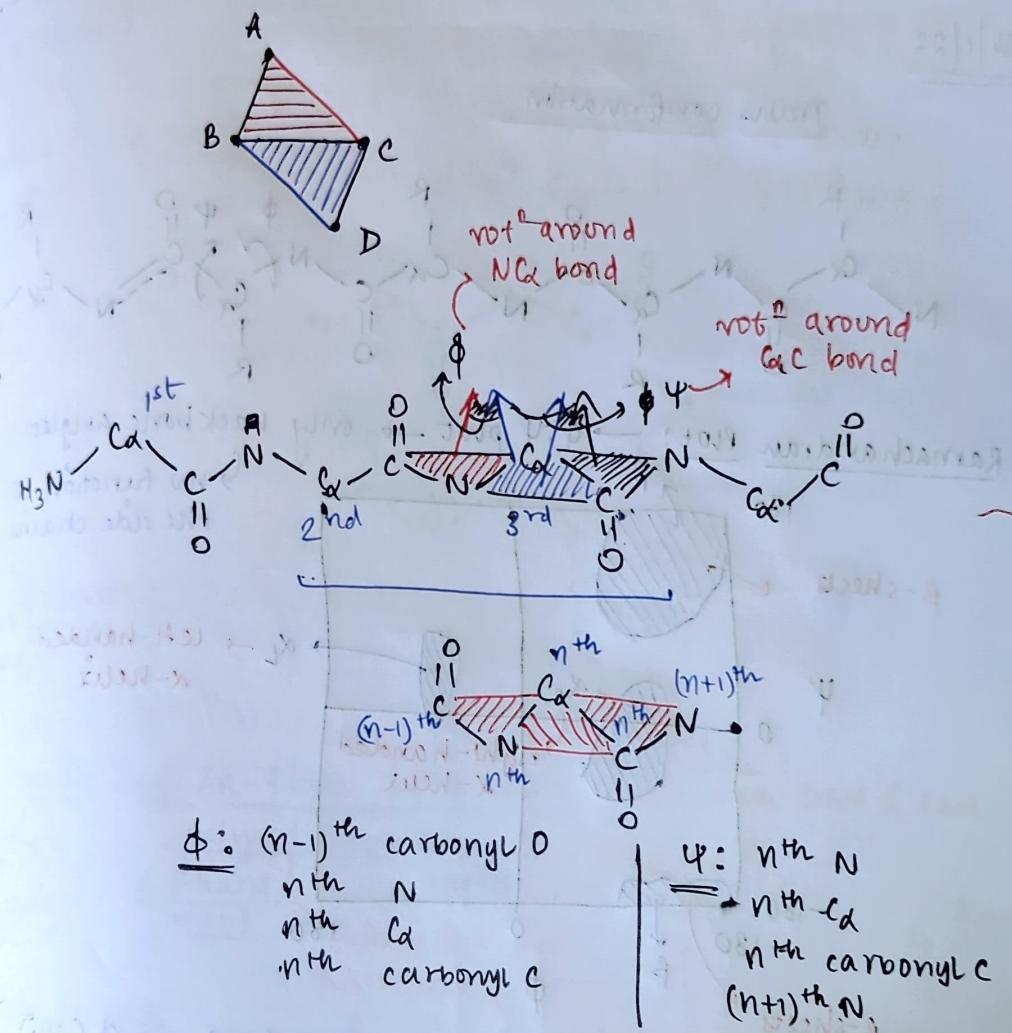
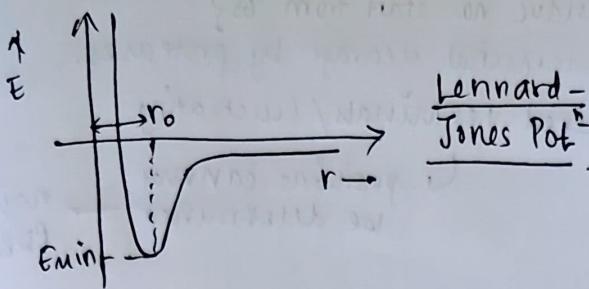
R group & carbonyl oxygen on opp. sides.

trans conformation

R group & carbonyl oxygen on the same side

steric hindrance

← interference of vdw radii



molecular visualization software

- PyMOL → most famous
- VMD
- RasMOL

protein crystallization → water molecules get trapped in the crystal → can see in a visualization → HETATM software

HETATM

Q. Read PDB Manual

Q. Why does residue no. start from 68?

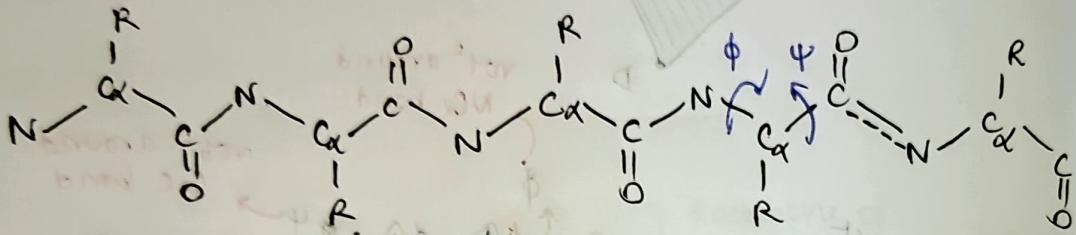
- ① N- or C-terminal cleavage by proteases.
- ② disordered terminals / fluctuating

↳ positions cannot be determined

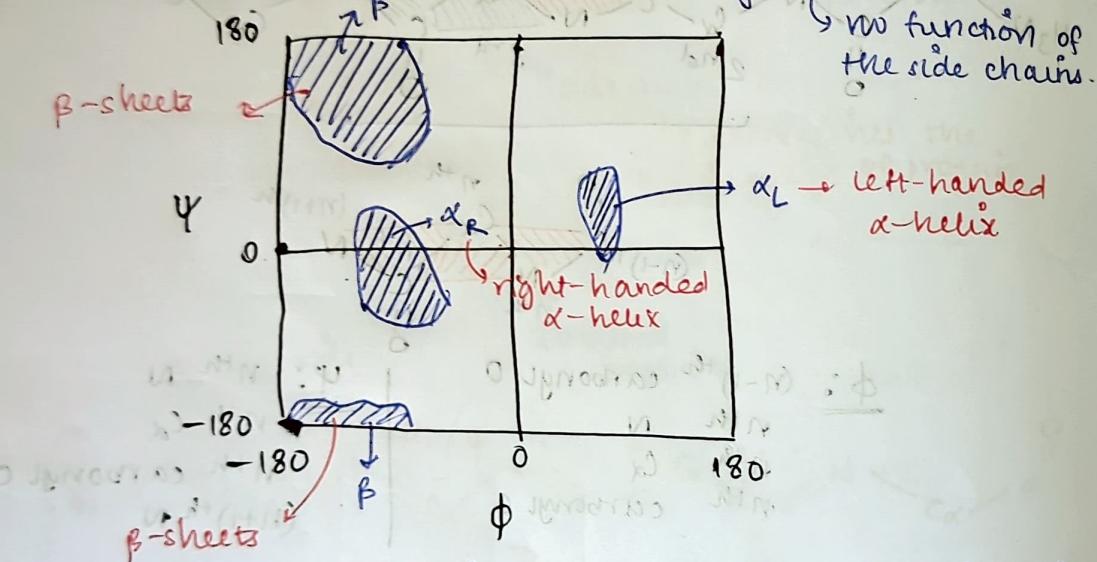
→ not in the PDB file.

24/1/22

Trans-conformation

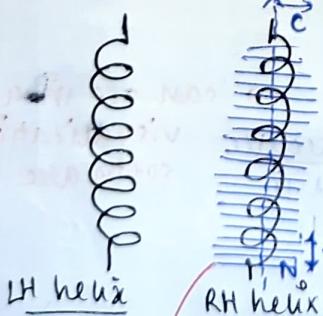


Ramachandran Plot: → ϕ, ψ plot → only backbone angles



rotation around N- α bond → ϕ (phi)

rotation around C- α bond → ψ (psi)



sliced at every a.a. residue

α -helix
per turn = 3.6 residues
rise per residue = 1.5 \AA
rise per turn = 5.4 \AA
Helices found in protein structures

- ↳ α -helices
- ↳ π -helices
- ↳ β_{10} -helices

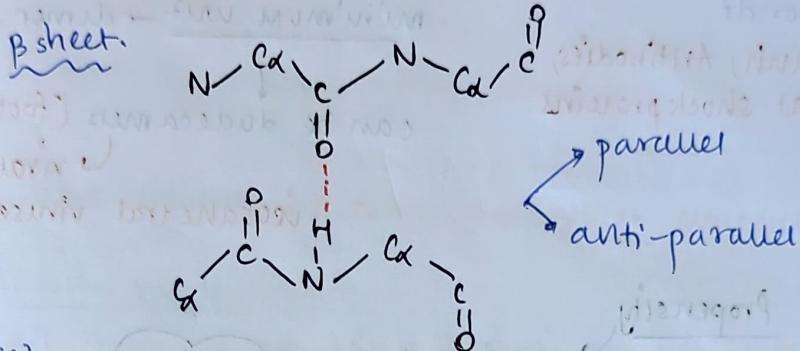
can be used to calculate length of helix

β -helix: if we compress α -helix
per turn = 3 residues

π -helix: if we extend the α -helix.
per turn = 5 residues

α -helix } \rightarrow 2 structures \rightarrow stabilized by H-bonds \rightarrow supports that hold helical structure
 β -sheet }
regular structure

α -helix H-bond pattern: i^{th} $N-H \cdots C=O$ } $(i+4)^{\text{th}}$ $C=O \cdots N-H$ } (or) i^{th} $C=O \cdots N-H$



irregular structure
loops

Motifs e.g. Zn-finger motifs \rightarrow interact with DNA & RNA

Calmodulin \rightarrow binding motif \rightarrow bind Ca^{2+}
E-F hand motif \rightarrow E helix & F-helix connected by a loop

E & H helices \rightarrow Hb where porphyrin ring is present

can be used to perform biotin formic acid analyses

e.g. presence of EF hand motif can predict Ca^{2+} binding ability of protein at hand.

RNA-binding proteins \rightarrow Arg-rich motifs.

Gamma-crystallin \rightarrow detect photons from EM waves in our eyes & convert it into nerve signals.

Tertiary structure

→ single polypeptide chain.
fold to form func² protein

e.g. Myoglobin (muscles)
Fibronectin (tear)

Quaternary structure

assembly of
more than one polypeptide
chain & perform functions

e.g. Hemoglobin

↓
some can be
covalently
linked
↓
by disulphide
bonds

↓
assembled by non-
covalent interactions

mostly hydrophobic
& vdw interactions,
electrostatic

e.g. insulin, Antibodies,
heat-shock proteins

minimum unit → dimer

↓
can be dodecamers (football-like)
viral capsids
icosahedral viruses

Propensity

ability to form
 α -helix

helical
residues

more likely to
form helices

proline → helix capping residue / helix breaker
if we want to terminate the helix

same protein

300 residues
↓
50 leucine
in total

	helix 1	helix 2
L		
I		
R		
H		
K		
L		
H		
R		
A		
G		

LRR structure (leucine-rich)

(say) Leucine (L)

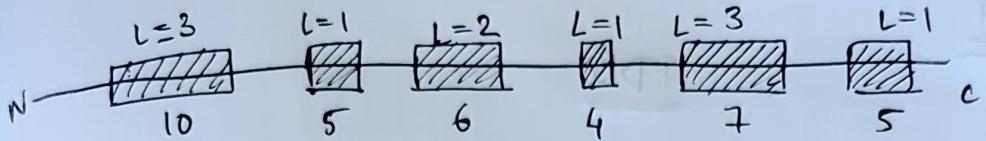
$$\text{freq. } (L) = 2/10$$

5 leucines in
19 helical residues

$$\text{Propensity} = \frac{2/10}{5/19}$$

$$5/19$$

y if we don't divide by this,
we cannot take into account
the fact that only 5 L's



$$\frac{3+1+2+1+3+1}{10+5+6+4+7+5} \quad \left. \begin{array}{l} f(h) \\ f_{ss}(h) \end{array} \right\}$$

$$P = \frac{\frac{70}{300}}{\left. \begin{array}{l} f_p(h) \\ f_{total} \end{array} \right\}}$$

propensity

sec. struct. helical region

(helical) 300 total a.a.

To total L in seqn

full framework

1DFU

PyMOL:

S-show

cartoon

ribbon

sticks

- only backbone

PyMOL > remove solvent # to remove water molecules

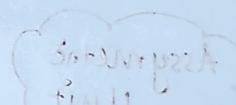
Viewing mode:

Edit mode \rightarrow select 2 atoms \rightarrow distance

 3 atoms \rightarrow angles

 4 atoms \rightarrow dihedral angle

PyMOL manual - 18 pages



PyMOL > select polymer, protein // to select protein only

PyMOL > select polymer, nucleic // to select nucleic acid only.

PyMOL > select chain M // to select chain M of nucleic acid

PyMOL > color red, chain M // coloring chain M to red

PyMOL > color blue, resi 91 // coloring residue 91

PyMOL > color orange, chain P and ALA // coloring all alanines of chain P.

31/1/24

PDB

Residues

2FBD

Asymmetric Unit

some protein \rightarrow monomeric
some protein \rightarrow oligomers \Rightarrow the functional unit of protein

homodimer:

(A A) two identical polypeptide chains
eg: tRNA synthetase (Asp)

heterodimer:

(A B) two different polypeptide chains
eg: light & heavy chains in Ab

Biological Assembly

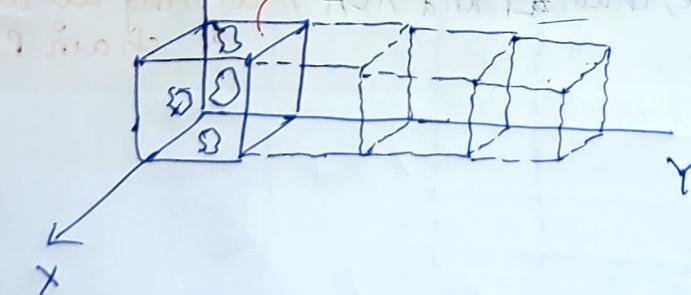
eg: $\alpha_2\beta_2$ of Hb
(all 4 chains are required for the function)

dimer of dimer

Asymmetric Unit

property of the crystal

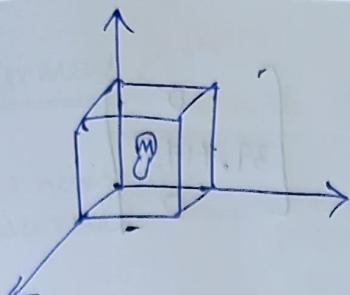
an unit cell may contain more than one molecule



crystallization is not a deterministic process.

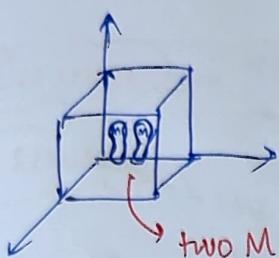
we cannot determine which lattice/group we are going to obtain

Myoglobin (say)



Myoglobin

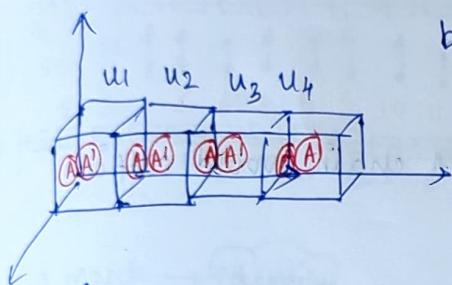
bio. assembly = monomeric
asy mm. unit = monomeric



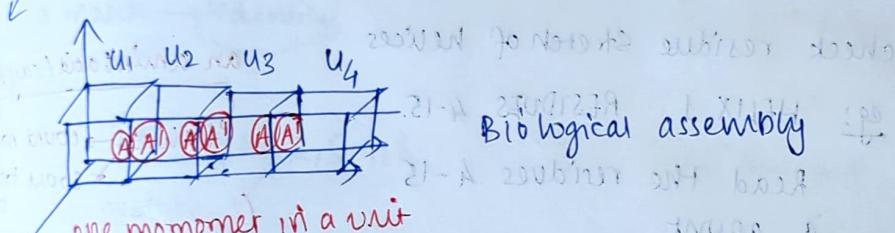
bio. assembly = monomer.
asy mm. unit = dimer

two M come together due
to forces during crystallization

Asp tRNA synthetase



biological assembly = (AA) = dimer
asy mm. unit = monomer



one monomer in a unit
cell & in adj. unit cell,
we have the other
monomer

Unit cell parameters in a PDB file:

CRYST1 36.523 79.435 45.203 90.0 102.97 90.00 P 1 211
 a b c α β γ

↓
2 fold screw
axis in y-
direction

transformation
matrix

$$T' = [R] \begin{pmatrix} X \\ Y \\ Z \end{pmatrix}_{3 \times 3} [T]_{3 \times 1}$$

Symmetry
Operator

$$\begin{pmatrix} X & Y & Z \\ -X & \frac{Y+1}{2} & -Z \end{pmatrix}$$

- Point group symmetry
- Space group symmetry

half unit cell
transl. along y-axis
 $\frac{79.435}{2} = 39.7175$

$$\begin{bmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{bmatrix} \times \begin{bmatrix} 39.7175 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} x_1' y_1' z_1' \\ x_2' y_2' z_2' \\ \vdots \vdots \vdots \\ x_n' y_n' z_n' \end{bmatrix} = \begin{bmatrix} x_1 y_1 z_1 \\ x_2 y_2 z_2 \\ \vdots \\ x_n y_n z_n \end{bmatrix}$$



Assignment:

↳ propensity of only A chain would suffice.

check residue stretch of helices

e.g.: HELIX 1 RESIDUES 4-15

Read the residues 4-15

& count

an amino acid (say Asp)

Helix 1 → count no. of Asp
Helix 2 → count total no. of Asp
...
Helix 8 → count # of Asp
→ # of aa

total Asp → total Asp
total aa → total aa

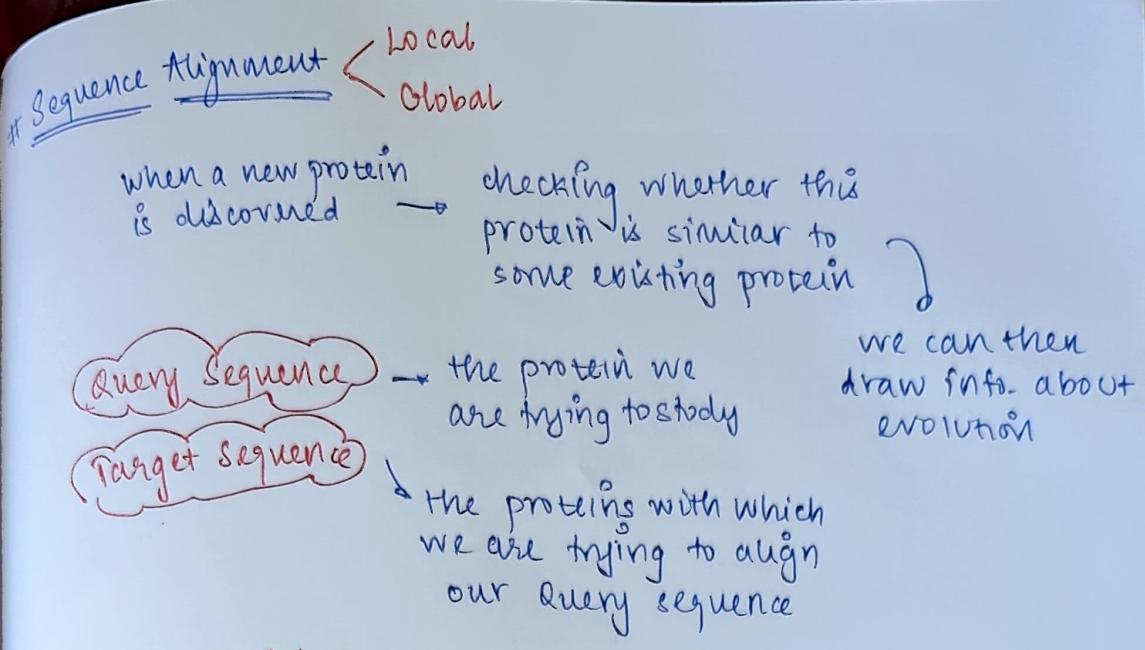
$$A_{\text{Asp}} = \frac{\# \text{ Asp in all } (8) \text{ helices}}{\# \text{ all aa in } 8 \text{ helices}}$$

$$\left(\frac{\text{total } \# \text{ Asp present in protein chain A}}{\text{total } \# \text{ aa in protein chain A}} \right)$$

Ala = $\frac{\# \text{ Ala in all } (8) \text{ helices}}{\# \text{ all aa in } 8 \text{ helices}}$

Asp = $\frac{\# \text{ Asp in all } (8) \text{ helices}}{\# \text{ all aa in } 8 \text{ helices}}$

Proline = $\frac{\# \text{ Proline in all } (8) \text{ helices}}{\# \text{ all aa in } 8 \text{ helices}}$



$\alpha =$

1	2	3	4	5	6	7	8	9	10	11	12	13	14
A	R	K	V	P	T	V	G	A	P	I	V	F	A
N	↑	↑	1	↑	↓	↑	↓	↑	↓	↑	↓	↑	C

$\beta =$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19
A	R	R	D	I	P	T	L	V	G	A	P	I	G	H	I	K		

a match → **Identity**

Identity: 1, 2, 5, 7, 8

position w.r.t. Q

Different: 6

both R & K are +ve charged

Similar: 3, 4

- ✓ if the aligned position has same amino acid
- ✓ if the aligned position has same type of amino acid residue (eg: hydrophobic, aromatic)

$$\% \text{ identity} = \frac{\text{no. of identical positions}}{\text{total no. of aligned positions}} \times 100$$

Identity
Similarity
eg: A R & K
D & E
T, D, F, W
N & Q.

$$\% \text{ similarity} = \frac{\text{no. of similar positions}}{\text{total no. of aligned positions}} \times 100$$

Similarity in nts:

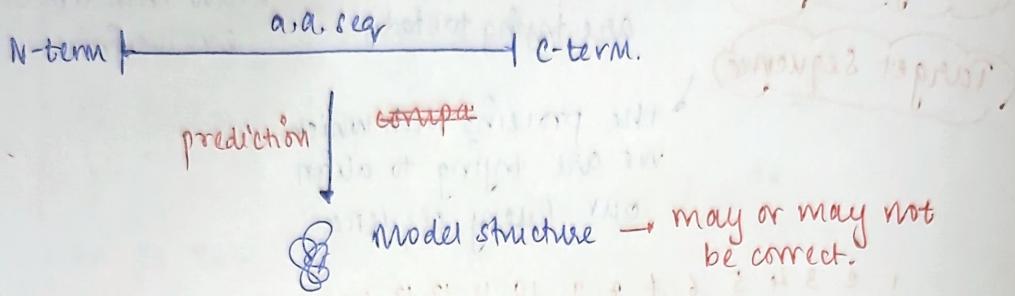
A & A → Identical

A & G → Similar (both purines)

7/2/24

Homology Modelling of Proteins

Prediction of 3D structure of a target protein from the amino acid seq (1st structure) of a homologous (template) protein for which an X-ray or NMR structure is available.

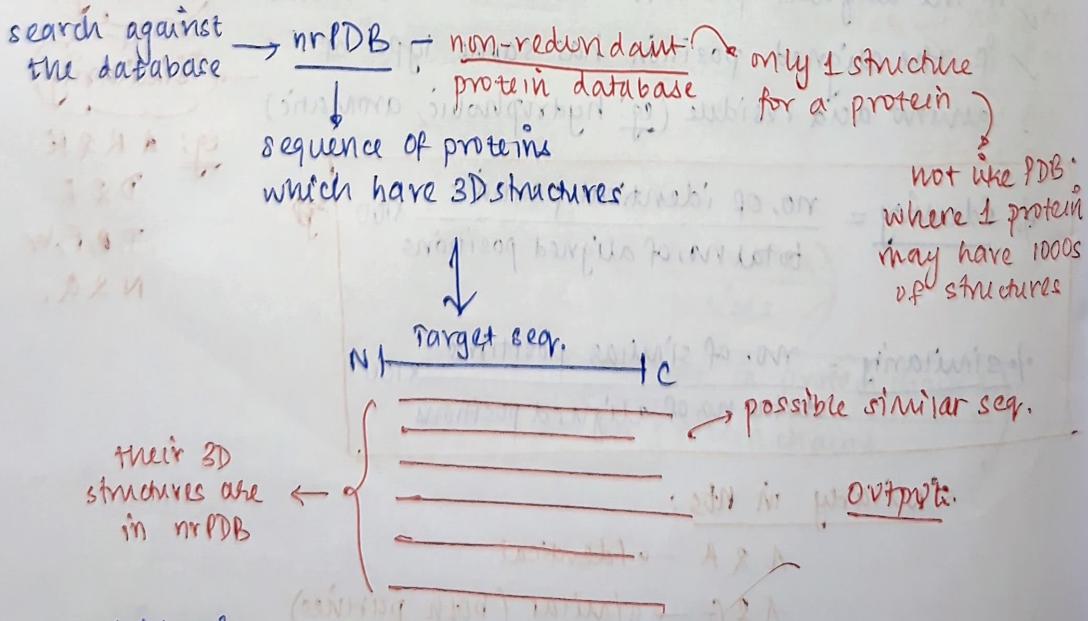
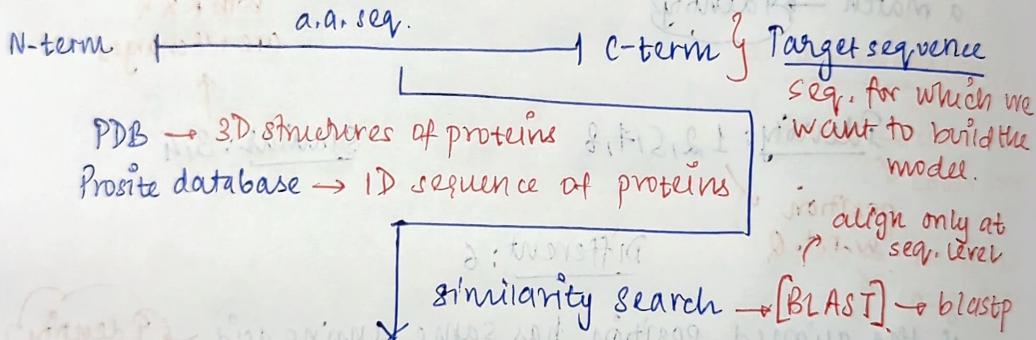


Methods:

① homology modelling / comparative modelling

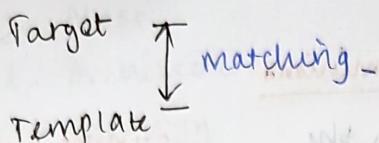
② Threading.

① Comparative/homology modelling:



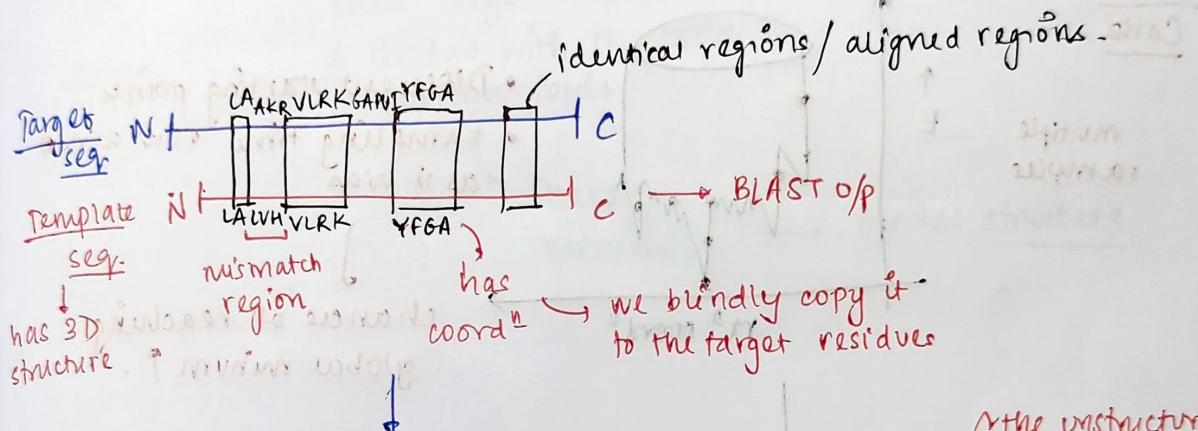
% identity, % similarity, Query Coverage → how many aa. are covered → poor query coverage makes homology modeling difficult

selection of one structure
that would act as
template structure → used as a reference
for homology modelling.



Conditions:

- ① Sequence identity b/w tg and tm > 30% 20-30% Twilight zone
- ② Query coverage → how much % of tg is aligned with tm > 95%



N \square $\overline{\square}$ \square \square \square \square \square \square now we know the structures of these fragments

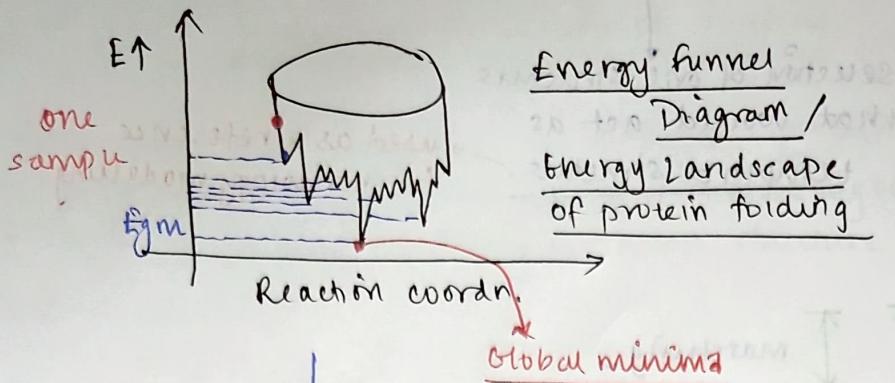
the unstructured parts
↓
fragmented parts of homology modeling

\square \square \square \square \square coordinates for each & every fragment take each fragment and do the BLAST search again on nrPDB

next step is to connect/stick the fragments together

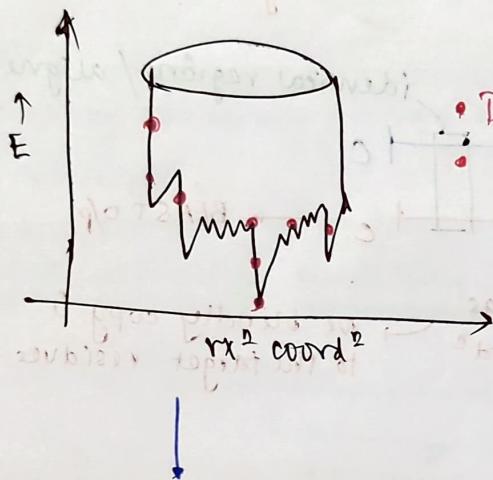
eq: A KR we will get similarity for such a small fragment

Energy minimization of the entire structure



more sampling → to obtain the global minima. → sample might get stuck to any local minima if it is deep enough

solution: 10 different starting structures



- Different starting points
- sampling time remains as it was

chances of reaching global minum ↑

once model is selected,
check Ramachandran plot

Homology Modelling Software

- ① SWISS-MODEL
- ② Modeller

CASP Critical assessment of structure prediction

John Hopkins University

✓ However, for homology modelling, template structure is a must.

② Threading: → performed when template structure is not available

3D structure space

what/how many different types of fold are possible in protein universe?

C - Class

A - Architecture

T - Topology

H - Homologous Superfamily

CATH Database

unit

1200 folds are possible in protein universe

fold space

any protein we take from any species, their architecture will be 1 of these 1200 unit-folds

take target sequence & thread with the 1200 possible folds

Energy minimizⁿ

→ Model structure.