Lesson 11. Content

1. The protein folding problem.

2. Comparative (Homology) modelling. Choice of template and alignment.

3. Comparative(Homology) modelling. Modelling of the protein core and loops

4. Comparative(Homology) modelling. Optimization and quality check.



Why knowing the 3D structure of proteins? **Experimental Models (X-Ray, NMR, EM)** Theoretical Models (protein modelling)

- 1. Unraveling the function/s if not known
- 2. Study of the molecular mechanisms of interaction
- 3. Study of biochemical mechanisms
- 4. Rational drug design

Example of experimental 3D structure in the PDB



© 3D View: Structure | 1D-3D View | Electron Density | Validation Report | Ligand Interaction

Global Symmetry: Asymmetric - C1 Global Stoichiometry: Hetero 2-mer - A1B1

🖪 6M0J

Crystal structure of SARS-CoV-2 spike receptor-binding domain bound with ACE2

🖿 Display Files 🗸

Ownload Files -

PDB DOI: 10.2210/pdb6M0J/pdb

Classification: VIRAL PROTEIN/HYDROLASE Organism(s): Homo sapiens, Severe acute respiratory syndrome coronavirus 2 Expression System: Trichoplusia ni

Mutation(s): No

Deposited: 2020-02-21 Released: 2020-03-18 Deposition Author(s): Wang, X., Lan, J., Ge, J., Yu, J., Shan, S.

Experimental Data Snapshot	wwPDB Validation	3D Report Full Report
Method: X-RAY DIFFRACTION	Metric	Percentile Ranks Value
Resolution: 2.45 Å	Rfree	0.228
R-Value Free: 0.227	Clashscore	4
R-Value Work: 0.192	Ramachandran outliers	0.1%
R-Value Observed: 0.194	RSRZ outliers	2.7%
	Worse Percentile re Decembin re	Better stative to all X-ray structures slative to X-ray structures of similar mentation

Besides the coordinates of the macromolecule(s) structure, additional information is reported, including especially the experimenthal method (here X-ray diffraction) and relative quality parameters (Resolution and R-values, the lower the better)

Why predicting the 3D structure of proteins?



Why predicting the 3D structure of proteins?



Between 2012 and 2015, the number of protein families has more than doubled, then it remained pretty stable

The hundreds of millions of known protein sequences are grouped in a finite number (thousands) of protein families

The folding problem:

Predicting the fold (3D structure) of a protein starting from



The folding problem:

Predicting the fold (3D structure) of a protein starting from

ACTFGARTEADEASRTFCGABHI _ GFRLPMNHTYWPLYHMVCS...

HGRTDEPLPMNWQACVFRGHEF GPLMNSSFGHINV...

MKWSDFHITTLPEQCVNTHILSS TPLMYHVGVCQQTHLMS...

LIPREDSGHWQPLMTRFHSDAAS LKPLRTENMVCDERSTGHKL...



The folding problem:

The cell environment can affect the folding, e.g. through chaperones (proteins that assist the conformational folding and unfolding)

ACTFGARTEADEASRTFCGABHI GFRLPMNHTYWPLYHMVCS...

HGRTDEPLPMNWQACVFRGHEF GPLMNSSFGHINV...

MKWSDFHITTLPEQCVNTHILSS TPLMYHVGVCQQTHLMS...

LIPREDSGHWQPLMTRFHSDAAS LKPLRTENMVCDERSTGHKL...





The Anfinsen's experiment



Denatured protein (inactive)

Control
The Anfinsen's experiment

Native protein (active)

Protein in a non-native conformation (inactive)

"Now, pack the hydrophobic core, fold helix A along the dotted line, taking charged residues of A close to the ionic groups on the surface of helix B..."



The Levinthal's paradox



The protein backbone can rotate around the dihedral angles ϕ & ψ



Intervals of 120 degrees...

N amino acids: 3²⁽ⁿ⁻¹⁾ alternatives 100 amino acids: 10⁴⁷ alternatives



If each attempt would take 100 femtoseconds (10⁻¹³ sec) trying all the possible combinations would take:

 $10^{47} \times 10^{-13}$ seconds = 10^{34} seconds = 2×10^{27} years !!



N amino acids: 3 ²⁽ⁿ⁻¹⁾ alternatives **1**00 amino acids: 10 ⁴⁷ alternatives



If each attempt would take 100 femtoseconds (10^{-13} sec) trying all the possible combinations would take:

 $10^{47} \times 10^{-13}$ seconds = 10^{34} seconds = 2×10^{27} years !!

And we are considering only (a subset) of the backbone conformations



Most proteins fold on timescales on the order of a millisecond (ms), with a median of ~5 ms

The "speed limit" of protein folding has been proposed to be set by N/ 100 µs, where N is the length of the protein (note: this only works with single-domain proteins)

Potential energy surface





a cross-section of the energy landscape

A surface can be drawn that represents the variation of the potential energy as the conformation varies

It resembles a landscape where valleys and peaks represent energy minima and maxima, respectively (the most frequently occuring conformations lie at the bottom of the deapest valleys)

Peaks are energy barriers, therefore reaching the global minimum is not trivial

The difference in energy between the native and non-native conformation can be very small



Even if we were able to generate all the possible conformations (and we aren't) a small error in the energy estimate would lead



The problem of folding: a hand from evolution





Example of homologous proteins



AMINO ACIDS ARE THE BUILDING BLOCKS OF PROTEINS IN LIVING ORGANISMS. THERE ARE OVER 500 AMINO ACIDS FOUND IN NATURE - HOWEVER, THE HUMAN GENETIC COD ONLY DIRECTLY ENCODES 20. 'ESSENTIAL' AMINO ACIDS MUST BE OBTAINED FROM THE DIET, WHILST NON-ESSENTIAL AMINO ACIDS CAN BE SYNTHESISED IN THE BODY.



Note: This chart only shows those amino acids for which the human genetic code directly codes for. Selenocysteine is often referred to as the 21st amino acid, but is encoded in a special manner. In some cases, distinguishing between asparagine/aspartic acid and glutamine/glutamic acid is difficult. In these cases, the codes asx (B) and glx (Z) are respectively used.

Example of homologous proteins





Function

structure



sequence

ACCEAAGAAGTCAGLVTCCTCGADGCTGAAGQRNDFCTGTQQGCTCLIAVRCCTACAACTRGACAAGTDFHAATGCAACCFHTTTGCTILMGAGGAAAAGGAGTLIGCNDEGAGGRCTTCTGAGFRCGGCAAFHRLCCATTQFCTGACAGDYGHTWTTTACYWACTTGCFTRNQKICCTCTGAWHY...

Similarity in sequence implies similarity in structure!



structure

TCCTCGADGC TGAAGQRNDF CTGTQQGCTC LIAVRCCTAC AACTRGACAA GTDFHAATGC AACCFHTTTG CTILMGAGGA AAAGGAGTLI GCNDEGAGGR CTTCTGAGFR CGGCAAFHRL CCATTQFCTG ACAGDYGHTW TTTACYWACT TGCFTRNQKI CCTCTGAWHY ...

Homology (or Comparative) Modelling that is

How to build a molecular model by homology

The tertiary structure of *homologous* proteins is better preserved than their primary structure (sequence)

Homology Modelling

Let's define the protein core



"Core": the structural conserved region Peripheral regions (indels) or "structural divergent regions" (SDR)

(Even extremely divergent) relatives in superfamilies contain well-conserved structural cores

CATHEDRAL: PLoS Comp. Biol. (2007), FLORA: PLoS Comp Biol. (2007), GRATH: J. Mol. Biol (2005), SSAP: J. Mol Biol (1989)

Root Mean Square Deviation (RMSD)



$$\mathsf{RMSD}_{\mathsf{N-1}} = \sqrt{\frac{\sum_{i=1}^{N-1} (d_i)^2}{N-1}}$$

It is a similarity measure for protein 3D structures



A (non biunivocal) relationship exists between the sequence identity and the structural similarity



Sequence identity

- Seq. Id. < 20%: "core" takes ~50% of the structure & r.m.s.d. of the backbone around 1.8 Å
- Seq. Id. > 50%: "core" takes ~ 90% of the structure & r.m.s.d. of the backbone around 1.0 Å

Chothia & Lesk, *EMBO Journal* (1986) vol <u>5</u>, pag. 823-826 *derived from 32 pairs of homologous proteins*







Step 2. Transfer of the *core* backbone coordinates from the template to the model

- Substitution of the side chains



Remember that the sequence alignment is our hypothesis of evolutionary correspondence between the target and template proteins!









The <u>sequence identity</u> with the template correlates with the extension of the protein core and affects the alignment quality (further improvable through MSAs and the structure analysis)

When Homology Modelling can trusthworthly be applied ?



Although... alignments of HMM profiles have challenged this "rule of thumb"

When Homology Modelling can trusthworthly be applied ?



Although... alignments of HMM profiles have challenged this "rule of thumb"

Rate of model confidence



Sequence identity

Model optimization:

- Serious errors cannot be corrected
 - We can eliminate sterically unfavorable interactions
 - We can optimize the backbone geometry in particular regions (e.g., around loops)





serious error in the protein packing



=> better packing of side chains

Energy optimization



Molecular mechanics (MM) approximates atoms as rigid spheres subjected to forces

MM force fields











Typical form:

$$V(\mathbf{r}^{N}) = \sum_{bonds} \frac{k_{i}}{2} (l_{i} - l_{i,0})^{2} + \sum_{angles} \frac{k_{i}'}{2} (\theta_{i} - \theta_{i,0})^{2}$$
$$+ \sum_{torsions} \frac{V_{i}}{2} [1 + \cos(i\omega - \gamma)] +$$
$$+ \sum_{i,j>i} \left\{ 4\varepsilon_{i} \left[\left(\frac{\sigma_{ij}}{r_{ij}} \right)^{6} - \left(\frac{\sigma_{ij}}{r_{ij}} \right)^{12} \right] + \frac{q_{i}q_{j}}{4\pi\varepsilon r_{ij}} \right\}$$

MM force fields





Example of Coulomb's interaction

Evaluating the model quality

"How much I can trust?"

In recognizing 'misfolded' models:

EM (energy minimization) & **MD** (molecular dynamics) are **inefficient**:

Normality indices can be more efficient:

e.g. regularity on the geometry, atomic contacts, distribution of polar and apolar residues, solvation potentials, etc., derived from the observation od deposited (experimental) 3D structures

Model quality assessment

Estimating the quality of protein structure models is a vital step

We may have a set of alternative models (e.g. from different modeling servers or based on alternative template structures and alignments) from which the best candidate shall be selected, or a singe model has been built from which we need to predict the quality in order to have an idea about its suitability for subsequent experiments

Traditionally, geometrical checks were performed on the models (see Ramachandran plots), however "statistical potentials" (also named "knowledge-based potentials") such as the packing quality or the inside/outside profile, based on the <u>comparison with experimental structures</u>, have been shown to be much more effective

Ramachandran plot



A traditional (and obsolete) model quality check

Model quality assessment "statistical (or "knowledge-based) potentials"

Statistical potentials or knowledge-based potentials are scoring functions derived from the analysis of known protein structures in the Protein Data Bank (PDB)

The most frequent something is, the most stable it is

$$u_{ij}(r) = -k_{\rm B}T\ln[g_{ij}(r)]$$

energy score =
$$\sum_{\text{atom pair}} u_{ij}(r)$$

K_B is the Boltzmann constant
T is the temperature *i* and *j* are two atoms
g_{ii}(r) is the pair distribution function (derived from a set of structure-known proteins)

Nowadays, one of the the most used/cited protein model quality check estimators is QMEAN (Qualitative Model Energy ANalysis, https://swissmodel.expasy.org/qmean/)

The QMEAN server provides access to three scoring functions for the quality estimation of protein structure models, which allow to:

i) rank a set of models and ii) identify potentially unreliable regions in them

Both single models and set of models can be analysed.

QMEAN provides both <u>global</u> (i.e. for the entire structure) and <u>local</u> (i.e. per residue) <u>quality estimates</u> on the basis of one single model

It is a linear combination of <u>statistical potential terms</u>: a new kind of torsion angle potential over three consecutive amino acids - to assess local geometry; a secondary structure-specific distance-dependent pairwise residue-level potential - to assess long-range interactions; a solvation potential term - to describe burial status of the residues.

QMEAN6 additionally uses two agreement terms evaluating the <u>agreement of predicted</u> (based on sequence) <u>and</u> <u>calculated secondary structure and solvent accessibility</u> Both global scores are originally in a range [0,1] with one being good.



The three scoring functions are:

QMEANDisCo uses the single terms of QMEAN as a basis.

In addition, it has a term predicting local per-residue quality estimates by assessing the <u>agreement of pairwise residue-</u> <u>residue distances with ensembles of distance constraints</u> (DisCo) extracted <u>from structures homologous</u> to the assessed model. If no homologues are found, the DisCo scores are not used.

All terms are combined using <u>neural networks</u> trained to predict per-residue scores in range [0,1]. The <u>global score</u> is the average per-residue score and the provided error estimate is based on global QMEANDisCo scores estimated for a large set of models of similar size to the input

Model quality assessment: ProQ2

ProQ2 also uses machine learning, specifically <u>support vector</u> <u>machines</u> (SVMs), to predict local as well as global quality of protein models

Scalar features from each protein model based on properties that can be derived from its <u>sequence</u>, e.g. <u>conservation</u>, <u>predicted secondary structure</u>, and <u>exposure</u>, or 3D coordinates, e.g. <u>atom-atom contacts</u>, <u>residue-residue</u> <u>contacts</u>, and <u>secondary structure</u>, are calculated and used as features to predict model correctness

http://bioinfo.ifm.liu.se/proq2/

The CASP experiment *Critical assessment of techniques for protein structure prediction*

CASP is a community-wide, worldwide experiment for protein structure prediction taking place every two years since 1994; in it the <u>accuracy</u> of different prediction methods is <u>blindly assessed</u> on common targets

The primary goal of CASP is to help advance the methods for predicting protein three-dimensional structure from its amino acid sequence

Like in a "world championship" in this field of protein structure prediction, over 100 research groups from all over the world participate in CASP on a regular basis

A model quality assessment section is included since CASP-7



The CASP experiment *Critical assessment of techniques for protein structure prediction*

Traditionally predictions in CASP are divided in:

<u>template-based</u> (including comparative modeling) & <u>template free</u> modeling

A revolution in CASP happened since 2020...

Target 89_2 Group 126 22 %





 $< 2.0 \text{ \AA}$ < 4.0 Å < 8.0 Å > 8.0 Å

Sample results from CASP 5, held in 2002

Target 128_2 Group 12 60 % $< 2.0 \text{ \AA} < 4.0 \text{ \AA} < 8.0 \text{ \AA} > 8.0 \text{ \AA}$ Target 137 Group 427 43 %



Sample results from CASP 5, held in 2002

Performance of predictors in CASP14 (2020)



GDT_TS: percentage of corresponding α-carbons within a 4 Å distance



Lesson 11. Content

- 1. The protein folding problem. Knowing a protein sequence is not enough to predict its folding by physico-chemical methods.
- 2. Comparative (Homology) modelling. Choice of template and alignment are the crucial steps of the method. The higher the target-template sequence identity, the better the expected quality of the model.
- 3. Comparative(Homology) modelling. Modelling of the protein core and loops. Serious errors in the core can be solved by correcting the target-template alignment. The loops (non-core) modelling is the method weakness.
- 4. Comparative(Homology) modelling. Optimization and quality check. The model quality can be assessed based on statistical potentials and evolutionary information.