

Structural Bioinformatics GENOME 541 Spring 2022

Lecture 3: Protein Structure Prediction Frank DiMaio (dimaio@uw.edu)

Last lecture



Typically, proteins fold by progressive formation of nativelike structures.

Folding energy surface is highly connected with many different routes to final folded state.

Structure Prediction

DEIVKMSPIIRFYSSGNAGLRTYIGDHKSCVMCTYWQNLLTYESGILLPQRSRTSR



Prediction Strategies

Homology Modeling

• Proteins that share similar sequences share similar folds.

• Use known structures as the starting point for model building.

• Can not be used to predict structure of new folds.

De Novo Structure Prediction

• Do not rely on global similarity with proteins of known structure

•Folds the protein from the unfolded state.

• Very difficult problem, search space is gigantic

Similar Sequences Share Similar Structures



Wilson, Kreychman, Gerstein (2000)

BLAST (<u>Basic Local Alignment Search Tool</u>)

BLAST is a fast sequence alignment algorithm that identifies high-scoring local alignments by finding short exact matches (seeds) and extending outward. BLAST uses the BLOSUM62 aa substitution matrix by default.

	С	S	Т	Р	A	G	Ν	D	Е	Q	Η	R	K	Μ	Ι	L	V	F	Y	W
С	9	-1	-1	-3	0	-3	-3	-3	-4	-3	-3	-3	-3	-1	-1	-1	-1	-2	-2	-2
S	-1	4	1	-1	1	0	1	0	0	0	-1	-1	0	-1	-2	-2	-2	-2	-2	-3
Т	-1	1	4	1	-1	1	0	1	0	0	0	-1	0	-1	-2	-2	-2	-2	-2	-3
Р	-3	-1	1	7	-1	-2	-1	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	-4
Α	0	1	-1	-1	4	0	-1	-2	-1	-1	-2	-1	-1	-1	-1	-1	-2	-2	-2	-3
G	-3	0	1	-2	0	6	-2	-1	-2	-2	-2	-2	-2	-3	-4	-4	0	-3	-3	-2
Ν	-3	1	0	-2	-2	0	6	1	0	0	-1	0	0	-2	-3	-3	-3	-3	-2	-4
D	-3	0	1	-1	-2	-1	1	6	2	0	-1	-2	-1	-3	-3	-4	-3	-3	-3	-4
Ε	-4	0	0	-1	-1	-2	0	2	5	2	0	0	1	-2	-3	-3	-3	-3	-2	-3
Q	-3	0	0	-1	-1	-2	0	0	2	5	0	1	1	0	-3	-2	-2	-3	-1	-2
Η	-3	-1	0	-2	-2	-2	1	1	0	0	8	0	-1	-2	-3	-3	-2	-1	2	-2
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5	2	-1	-3	-2	-3	-3	-2	-3
Κ	-3	0	0	-1	-1	-2	0	-1	1	1	-1	2	5	-1	-3	-2	-3	-3	-2	-3
Μ	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5	1	2	-2	0	-1	-1
Ι	-1	-2	-2	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4	2	1	0	-1	-3
L	-1	-2	-2	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4	3	0	-1	-2
V	-1	-2	-2	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4	-1	-1	-3
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6	3	1
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7	2
W	-2	-3	-3	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11

PSI-BLAST

- <u>Position-Specific Iterated BLAST</u>
- Allows more distantly related sequences to be identified
- Steps
 - 1. Use BLAST to identify related sequences
- 2. Create a profile from related sequences
- 3. Search for related sequences using this profile

Sequence Profile

- For each column in a MSA count how often each amino acid occurs
- Combine with prior information about substitution frequencies (ie. BLOSUM62)
- Convert counts to log odds scores. End product is a Position-Specific Scoring Matrix (PSSM)

1 lbpi 2 lbpi 3 lbzxI 4 lfakI 5 lbunB 6 lbf0	· · · · · · · · · · · · · · · · · · ·		FCI FCI FCI FCI DCD	EPI EPI EPI EPI KPI	YT YT YD DT VRI		KAI KAI RAI RAI	RII RII LHL VVR	RYF RYF RYF RYF AFY SFY	YNJ YNJ YNJ YNJ YKJ FKJ	
1 1bpi		F	C	г	E	P	2	Y	T	G	
2 1bpi		F	C	г	E	P	P	Y	T	G	
3 1bzx	1bzxI				E	P	P	Y	т	G	
4 1fak	1fakI					P	P	Y	D	G	
5 1bun	в	D	C	D	K	P	P	D	т	K	
6 1bf0		Y	C	ĸ	E	P	v	R	I	G	
Number of Number of	ACDEFCHHKIMN ORSEVWY .	00104000000000000000000	060000000000000000000000000000000000000	001000001400000000000000000000000000000	000000000000000000000000000000000000000	000000000000000000000000000000000000000	000000000000000000000000000000000000000	001000000000001000040	001000010000000040000	000005001000000000000000000000000000000	

Homology Modeling

- Identify homologous
 protein sequences
- Build model by
 - 1. "Threading" residues in corresponding positions of homologous structure
 - 2. Sampling conformations of unaligned residues
 - 3. All-atom refinement

MNDD--VDIQ---QSYP-FSI... LTDSQLAQVAAFVNNYPNVEL...



De novo protein structure prediction

MQIFVKTLTGKTIT LEVEPSDTIENVKA KIQDKEGIPPDQQR LIFAGKQLEDGRTL SDYNIQKESTLHLV LRLRGG



Thermodynamic hypothesis:

The native state is the lowest-energy conformation.

Structure Prediction Protocol

• Large-scale search of conformational space using a low-resolution potential

• Refinement of candidate models in a physically realistic, all-atom potential; selection by energy







Insights from Folding Studies

- 1. Local (*sequence-specific*) interactions strongly bias conformational sampling.
- 2. Folding is guided by hydrophobic burial, assembly of secondary structure, excluded volume.
- Native interactions on average stronger / more consistent than non-native interactions => native minima broader than nonnative minima.

Fragment-based Methods (Rosetta)

- **Hypothesis:** the PDB database contains all the possible conformations that a short region of a protein chain might adopt
- How do we choose fragments that are most likely to correctly represent the query sequence?





Fragment Libraries

- A unique library of fragments is generated for each 9-residue window in the query sequence.
- Assume that the distributions of conformations in each window reflects conformations this segment would sample.
- Regions with very strong local preferences will not have a lot of diversity in the library. Regions with weak local preferences will have more diversity in the library.

Generating Structures from Fragments

- Low resolution energy function used in initial search through conformational space
- Side chains represented by single "centroid" pseudoatom
- Major contributions from
 - Hydrophobic burial
 - Beta-strand pairing
 - Steric overlap
 - Specific residue pair interactions





High-resolution model refinement



sidechain rotamers

- Rosetta "relax" structure refinement:
 - Discrete sidechain optimization via Simulated Annealing Monte Carlo
 - Gradient-based minimization of energy with respect to *torsion angles*
- Potential function: Rosetta all-atom energy
 - Lennard-Jones,
 - LK implicit solvation,
 - Coloumb electrostatics
 - orientation-dependent hydrogen bonding,
 - PDB derived torsional potential



Correlated mutations carry information about distance relationships in protein structure.





Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, et al. (2011) Protein 3D Structure Computed from Evolutionary Sequence Variation. PLOS ONE 6(12): e28766. https://doi.org/10.1371/journal.pone.0028766 http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0028766

Learning the DCA (direct coupling analysis) matrix

The essence of DCA is then to assume that the rows, i.e. our aligned homologous proteins, are independent events drawn from a Potts-model probability distribution,

$$P(\boldsymbol{\sigma}) = \frac{1}{Z} \exp\left(\sum_{i=1}^{N} h_i(\sigma_i) + \frac{1}{2} \sum_{i,j=1}^{N} J_{ij}(\sigma_i,\sigma_j)\right),\tag{1}$$

and to use the interaction parameters J_{ij} as predictions of spatial proximity among amino-acid pairs in the protein structure.

Problem: **Z** cannot be tractably computed

Solutions:

- Mean-field approach (mfDCA) (<u>https://www.pnas.org/content/108/49/E1293</u>)
- Pseudo-likelihood (plmDCA) (<u>https://journals.aps.org/pre/abstract/10.1103/PhysRevE.87.012707</u>)

Predicted 3D structures for three representative proteins.



Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, et al. (2011) Protein 3D Structure Computed from Evolutionary Sequence Variation. PLOS ONE 6(12): e28766. https://doi.org/10.1371/journal.pone.0028766 http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0028766 Correlated mutations carry information about distance relationships in protein structure.



Coevolution guided modeling

GREMLIN predictions on shallow MSAs (Nseq=36, Nf=2.3)



Native contact map



Contact maps = Computer Images?



Contacts



OUTPUT



Convolutional neural networks





RESEARCH ARTICLE

Accurate De Novo Prediction of Protein Contact Map by Ultra-Deep Learning Model

Sheng Wang[®], Siqi Sun[®], Zhen Li, Renyu Zhang, Jinbo Xu*

Toyota Technological Institute at Chicago, Chicago, Illinois, United States of America

These authors contributed equally to this work.
 * jinboxu@gmail.com

Abstract

Motivation

Protein contacts contain key information for the understanding of protein structure and function and thus, contact prediction from sequence is an important problem. Recently exciting progress has been made on this problem, but the predicted contacts for proteins without many sequence homologs is still of low quality and not very useful for de novo structure prediction.

Learning a contact map from co-evolving residues



Inferring better contact maps (I)





Fig 6. Overlap between top L/2 predicted contacts (in red or green) and the native contact map (in grey) for CAMEO target 2nc8A. Red (green) dots indicate correct (incorrect) prediction. (A) The comparison between our prediction (in upper-left triangle) and CCMpred (in lower-right triangle). (B) The comparison between our prediction (in upper-left triangle) and MetaPSICOV (in lower-right triangle).

Inferring better contact maps (II)



Fig 9. Overlap between top L/2 predicted contacts (in red or green) and the native contact map (in grey) for CAMEO target 5dcjA. Red (green) dots indicate correct (incorrect) prediction. (A) The comparison between our prediction (in upper-left triangle) and CCMpred (in lower-right triangle). (B) The comparison between our prediction (in upper-left triangle) and MetaPSICOV (in lower-right triangle).



trRosetta

В



Improved protein structure prediction using predicted interresidue orientations

^(D) Jianyi Yang, ^(D) Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and ^(D) David Baker

PNAS January 21, 2020 117 (3) 1496-1503; first published January 2, 2020 https://doiorg.offcampus.lib.washington.edu/10.1073/pnas.1914677117

Discovering hidden patterns with a learned model

Gremlin predictions on shallow MSAs

(Nseq=36, Nf=2.3)



trRosetta predictions on shallow MSAs (Nseq=36, Nf=2.3)



Native contact map



Improving protein structure prediction

Free modeling accuracy in CASP



A differentiable end-to-end structure predictor

trRosetta







What would be a proper inductive bias for protein structure prediction?



Convolutional Networks (e.g. computer vision)

- data in regular grid
- information flow to local neighbours





Attention Module (e.g. language)

- data in unordered set
- information flow dynamically controlled by the network (via keys and queries)





Component 1: MSA updates via self-attention



Component 2: Update pair features via self-attention



Axial Attention (attention over rows then columns) to reduce memory requirements & computation time

Component 3: Extract pair features from MSA



 Ri
 Rj

 Goncat to
 Concat to

 Viginal pair
 & ResNet

Outer product &

residues

Non-interacting pairs \rightarrow Broader distribution Interacting pairs (co-mutating) \rightarrow Sharper distribution

aggregate

Ju, Fusong, et al. "CopulaNet: Learning residue co-evolution directly from multiple sequence alignment for protein structure prediction." *bioRxiv* (2020).

Component 4: Update MSA based on pair features



Component 5: SE(3)-Transformer for structure refinement



Graph connecting nearby residues

.

Predicted per-residue errors

RosettaFold 2-track model: Reproduce Alphafold 2 based on underlying principles



12 two-track blocks (orange box) + SE(3)-Transformer at the end Trained on protein structures in PDB (clustered w/ seqID cutoff 30%)

What happens during iteration?



What questions still remain?

- Predicting proteins *without* MSA information
- Predicting conformational states
- Predicting effects of mutation
- Predicting complex structures (particularly *pathogen/host* interactions)
- Predicting protein/nucleic acid complexes and RNA structures (Thursday!)
- Predicting protein/ligand complexes