



Structural Bioinformatics

GENOME 541

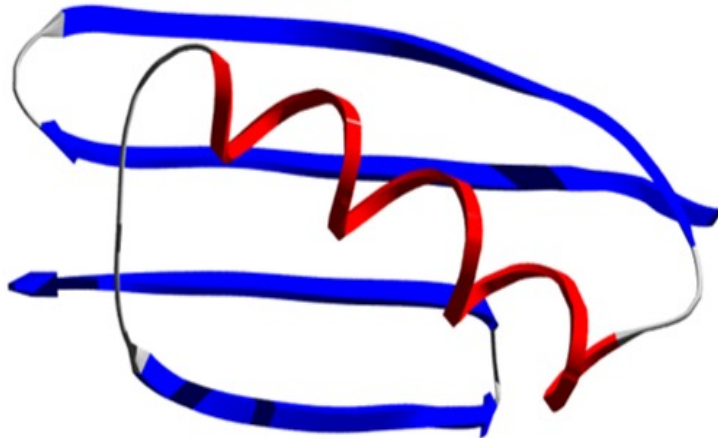
Spring 2023

Lecture 3: Protein Structure Prediction

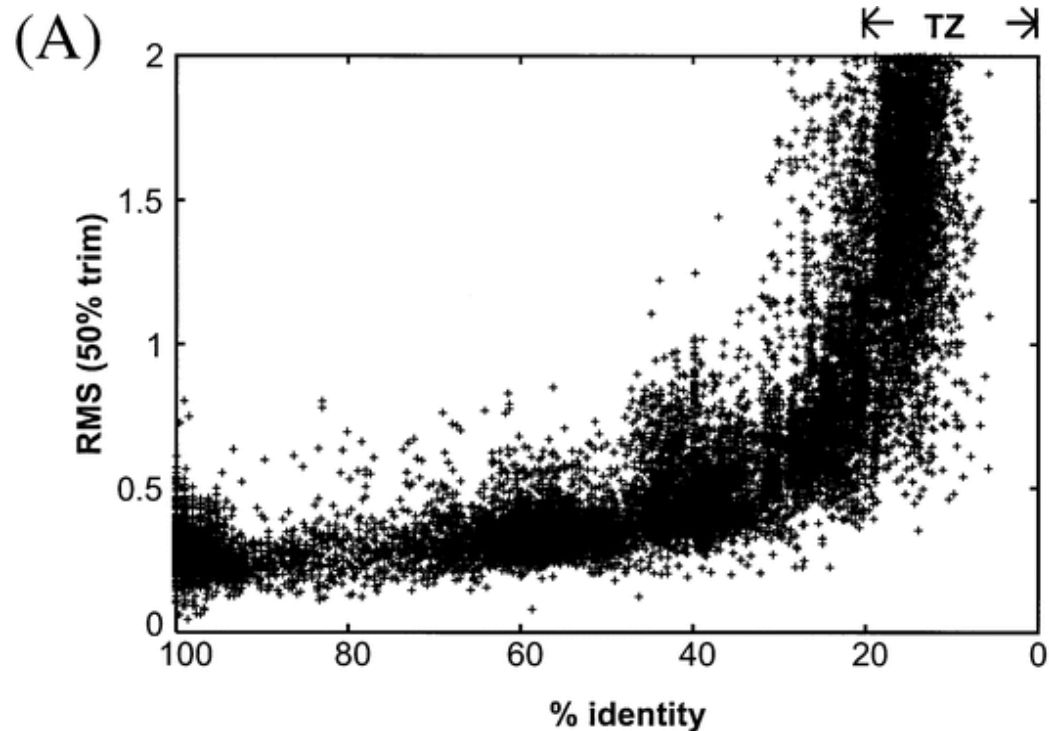
Frank DiMaio (dimaio@uw.edu)

Structure Prediction

DEIVKMSPIIRFYSSGNAGLRITYIGDHK
SCVMCTYWQNLLEYESGILLPQRSRTSR



Prediction Strategies



Wilson, Kreychman, Gerstein (2000)

Homology Modeling

- Proteins that share similar sequences share similar folds.
- Use known structures as the starting point for model building.

De Novo Structure Prediction

- Do not rely on global similarity with proteins of known structure
- Folds the protein from the unfolded state.

BLAST (Basic Local Alignment Search Tool)

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.

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W
C	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
T	-1	1	4	1	-1	1	0	1	0	0	0	-1	0	-1	-2	-2	-2	-2	-2	-3
P	-3	-1	1	7	-1	-2	-1	-1	-1	-1	-2	-2	-1	-2	-3	-3	-2	-4	-3	-4
A	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
N	-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
E	-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
H	-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
K	-3	0	0	-1	-1	-2	0	-1	1	1	-1	2	5	-1	-3	-2	-3	-3	-2	-3
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5	1	2	-2	0	-1	-1
I	-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

- Position-Specific Iterated BLAST
- 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

```

1 1bpi ..RPDFCLEPPYTGPCKARIIRYFYNA
2 1bpi ..RPDFCLEPPYTGPCKARIIRYFYNA
3 1bzxI ..R9DFCLEPPYTGPCKARIIRYFYNA
4 1fakI ..APDFCLEPPYDGPCKRALHLRYFYNA
5 1bunB ..RHFDCDKPPDTKICQTVVRAFYYKP
6 1bf0 ..PPWYCKEPPVRIGSCKKQFSSFYFKW

```

```

1 1bpi F C L E P P Y T G
2 1bpi F C L E P P Y T G
3 1bzxI F C L E P P Y T G
4 1fakI F C L E P P Y D G
5 1bunB D C D K P P D T K
6 1bf0 Y C K E P V R I G

```

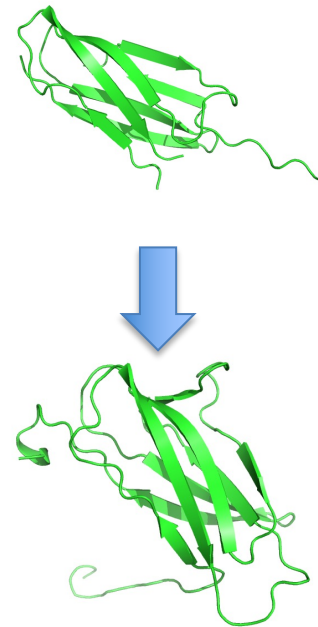
Number of A	0	0	0	0	0	0	0	0	0
Number of C	0	6	0	0	0	0	0	0	0
Number of D	1	0	1	0	0	0	1	1	0
Number of E	0	0	0	5	0	0	0	0	0
Number of F	4	0	0	0	0	0	0	0	0
Number of G	0	0	0	0	0	0	0	0	5
Number of H	0	0	0	0	0	0	0	0	0
Number of I	0	0	0	0	0	0	0	1	0
Number of K	0	0	1	1	0	0	0	0	1
Number of L	0	0	4	0	0	0	0	0	0
Number of M	0	0	0	0	0	0	0	0	0
Number of N	0	0	0	0	0	0	0	0	0
Number of P	0	0	0	0	6	5	0	0	0
Number of Q	0	0	0	0	0	0	0	0	0
Number of R	0	0	0	0	0	0	1	0	0
Number of S	0	0	0	0	0	0	0	0	0
Number of T	0	0	0	0	0	0	0	4	0
Number of V	0	0	0	0	0	0	0	0	0
Number of W	0	0	0	0	0	0	0	0	0
Number of Y	1	0	0	0	0	1	4	0	0
Number of .	0	0	0	0	0	0	0	0	0

- 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)

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:

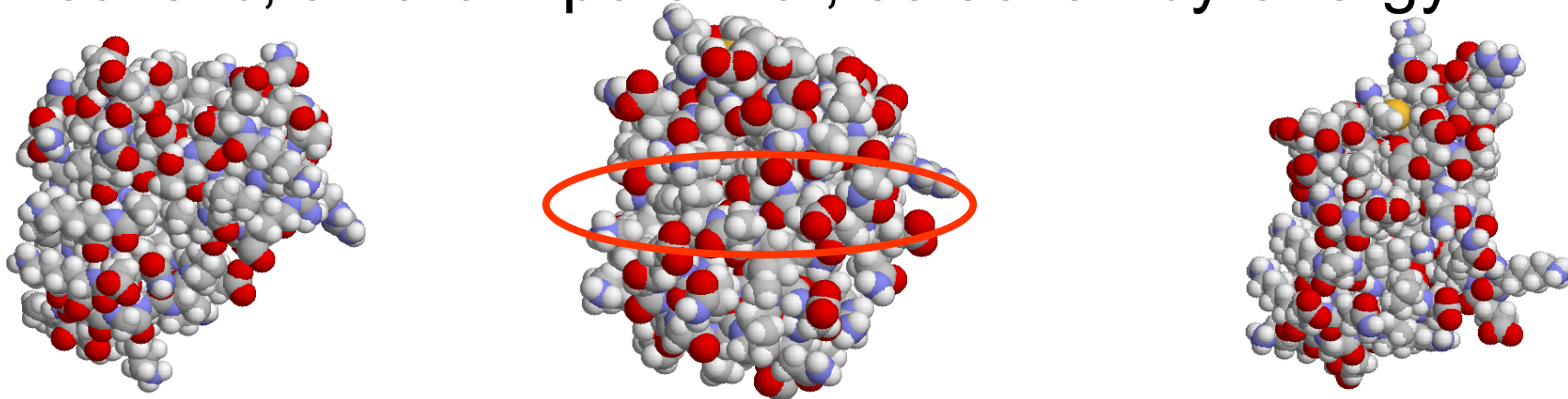
The native state is the lowest-energy conformation.

Structure Prediction Protocol

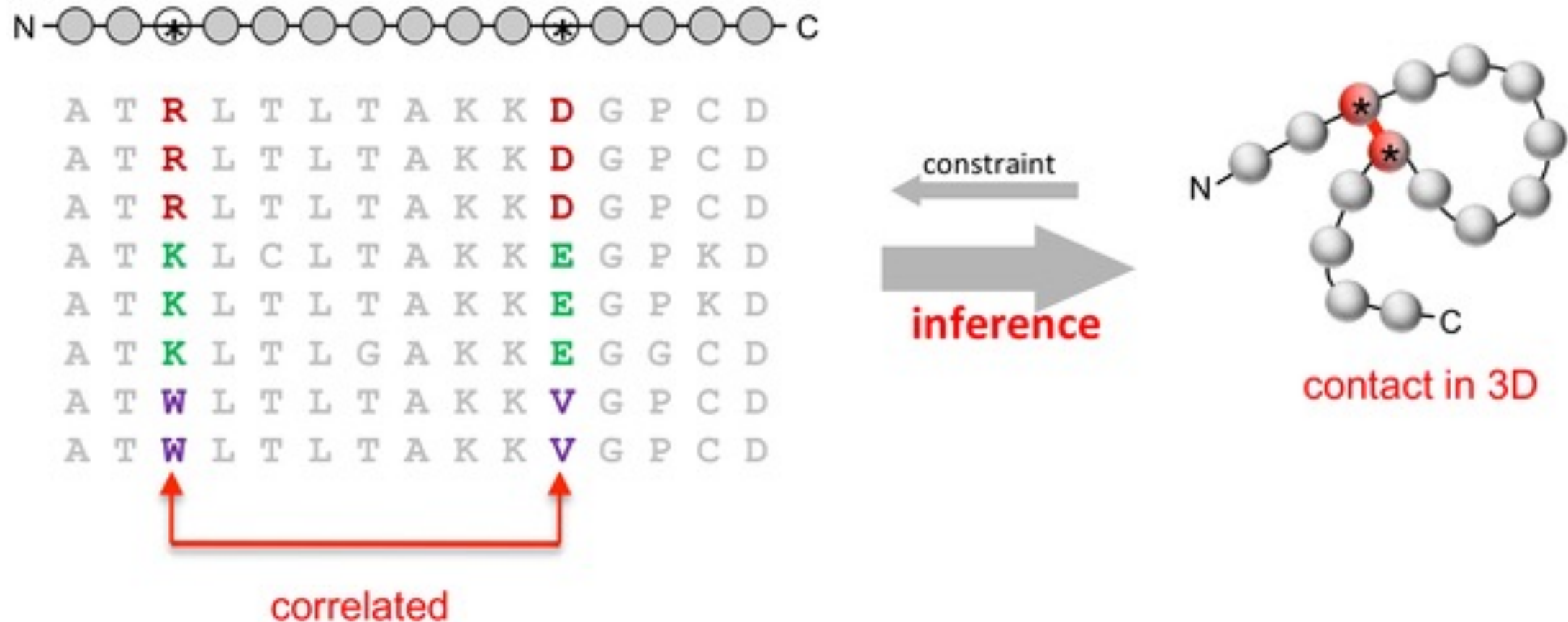
- Large-scale search of conformational space using a low-resolution potential



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



Correlated mutations carry information about distance relationships in protein structure.



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,

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

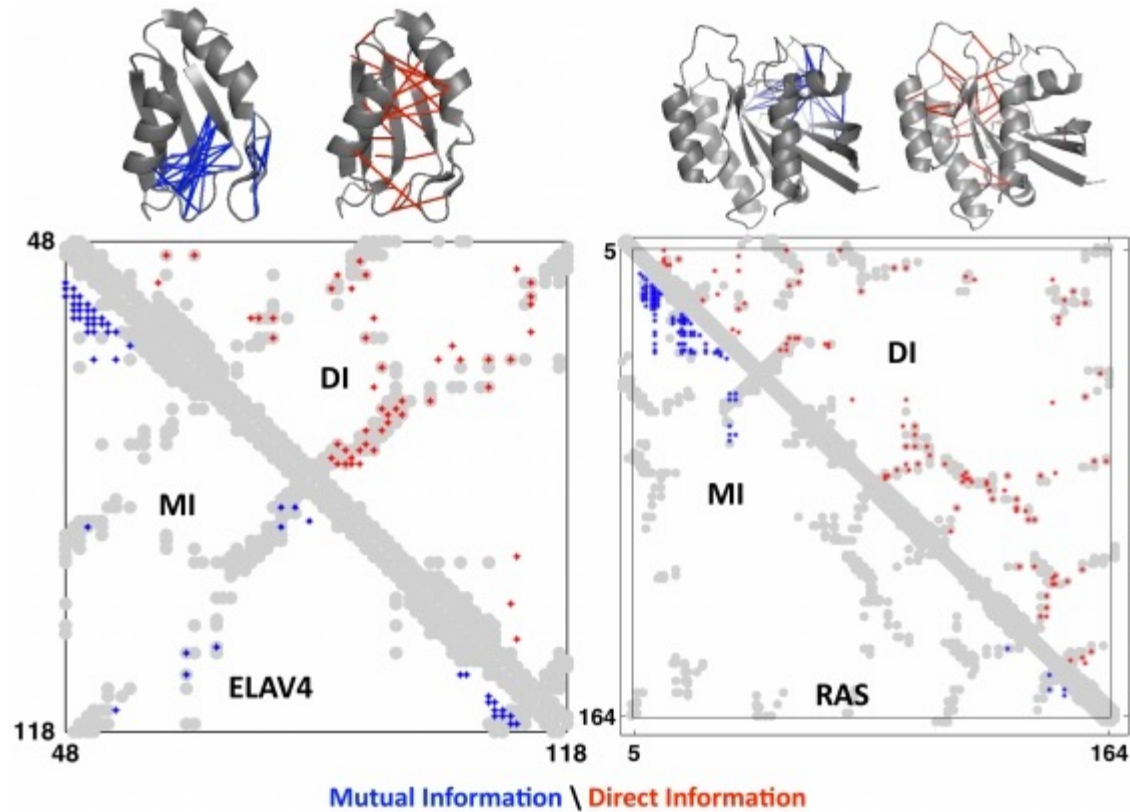
and to use the interaction parameters J_{ij} as predictions of spatial proximity among amino-acid pairs in the protein structure.

Problem: Z cannot be tractably computed

Solutions:

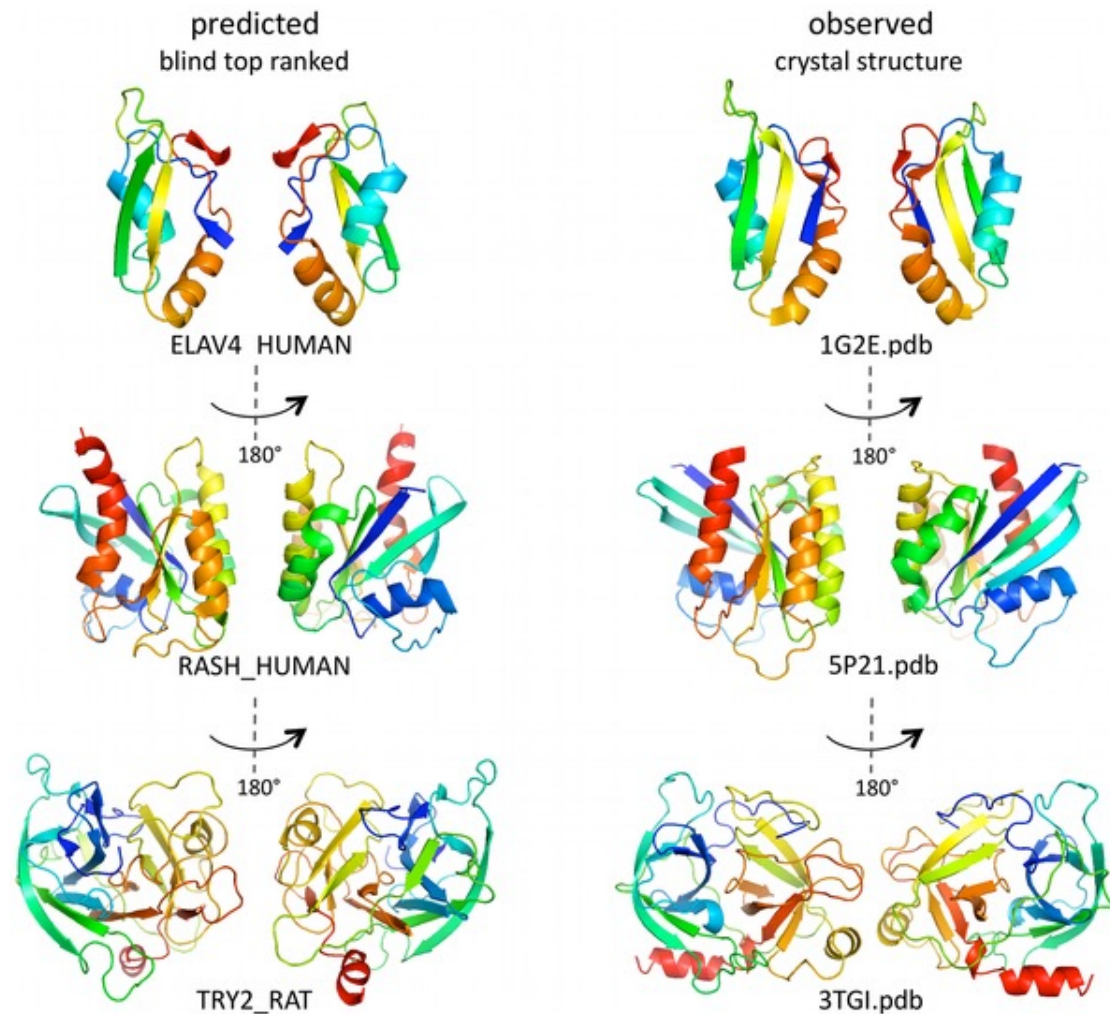
- Mean-field approach (mfDCA)
(<https://www.pnas.org/content/108/49/E1293>)
- Pseudo-likelihood (plmDCA)
(<https://journals.aps.org/pre/abstract/10.1103/PhysRevE.87.012707>)

Correlated mutations carry information about distance relationships in protein structure.



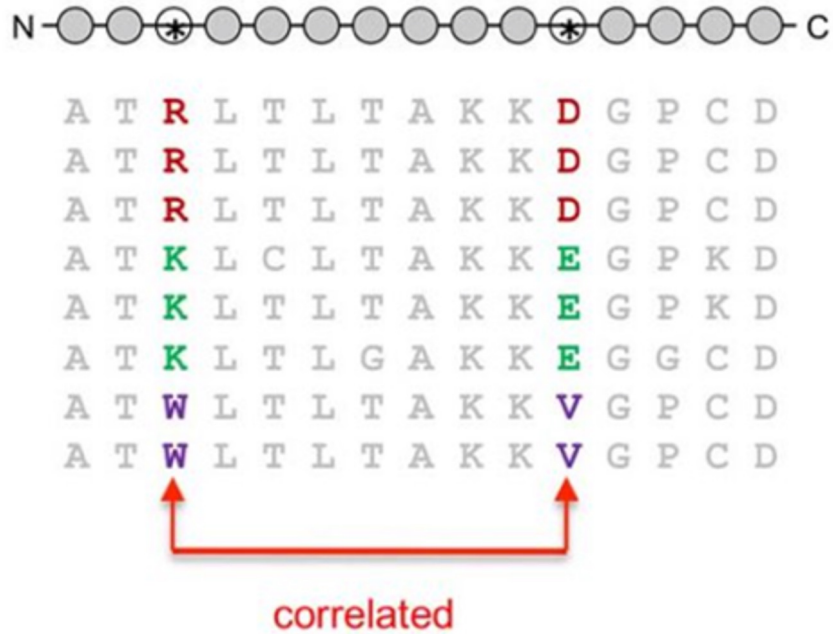
$$P(\mathbf{X}=\mathbf{x}) = \frac{1}{Z} \exp \left(\sum_{i=1}^L \left[\mathbf{v}_i(x_i) + \sum_{j>i}^L \mathbf{w}_{ij}(x_i, x_j) \right] \right).$$

Predicted 3D structures for three representative proteins.

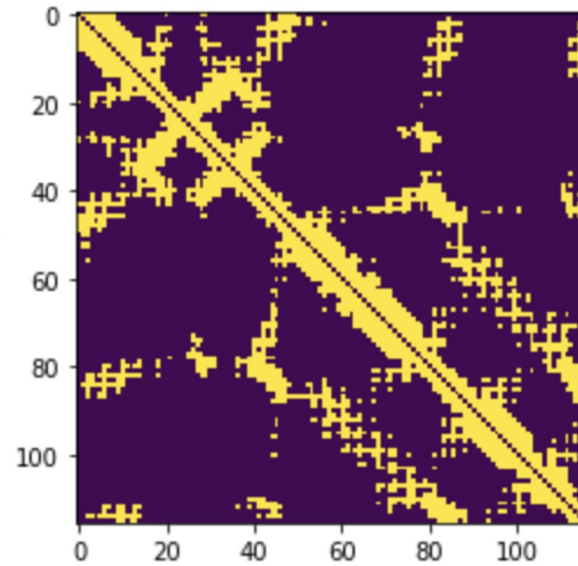


Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, et al. (2011) Protein 3D Structure Computed from Evolutionary Sequence Variation. PLOS ONE 6(12): e28766. <https://doi.org/10.1371/journal.pone.0028766> <http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0028766>

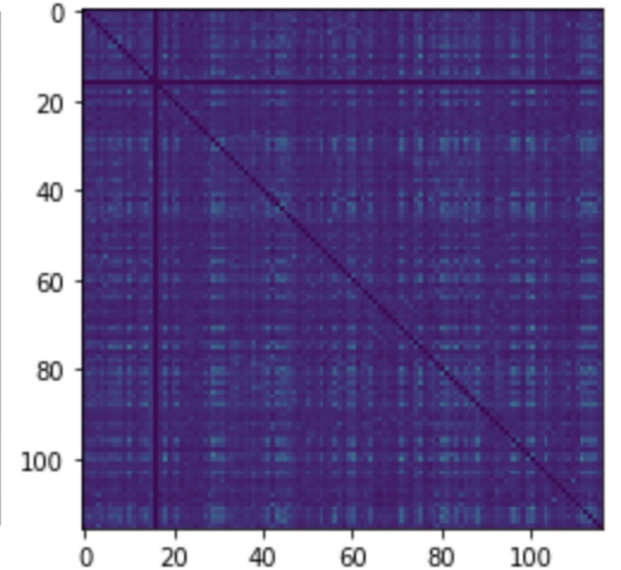
Coevolution guided modeling



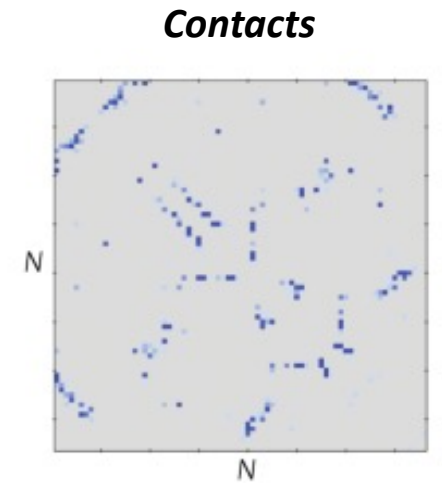
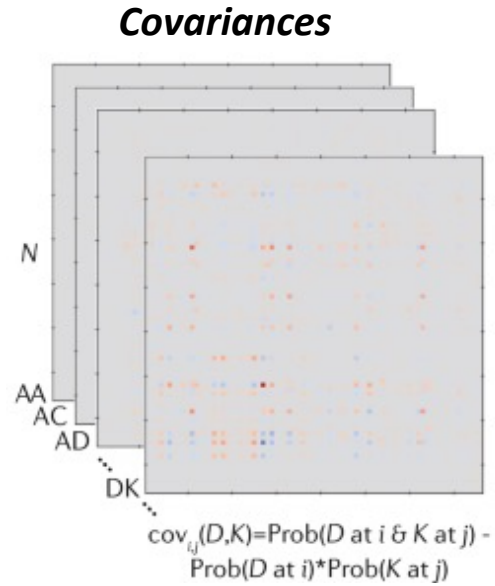
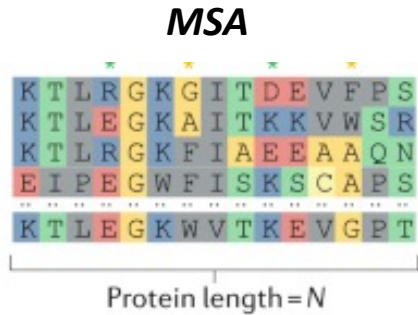
Native contact map



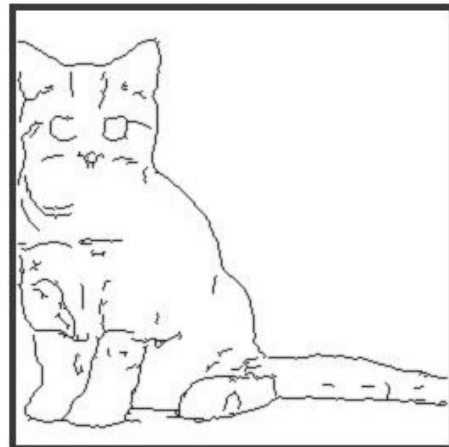
GREMLIN predictions on shallow MSAs
(Nseq=36, Nf=2.3)



Contact maps = Computer Images?



INPUT



OUTPUT



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.

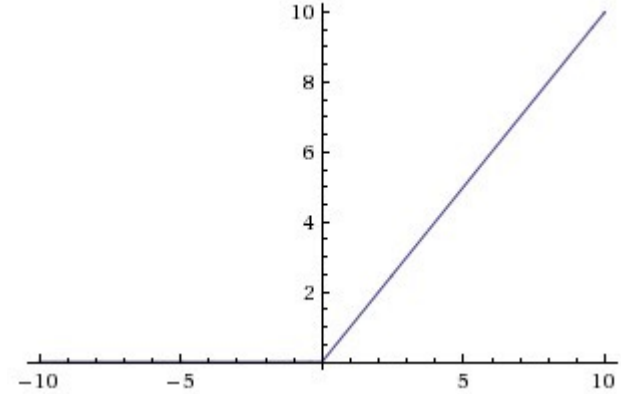
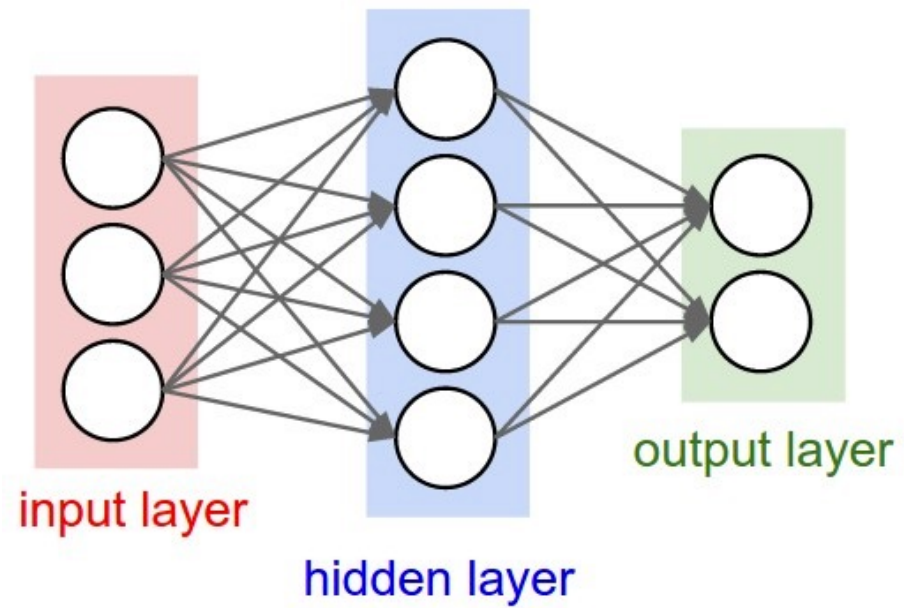
* jinboxu@gmail.com

Abstract

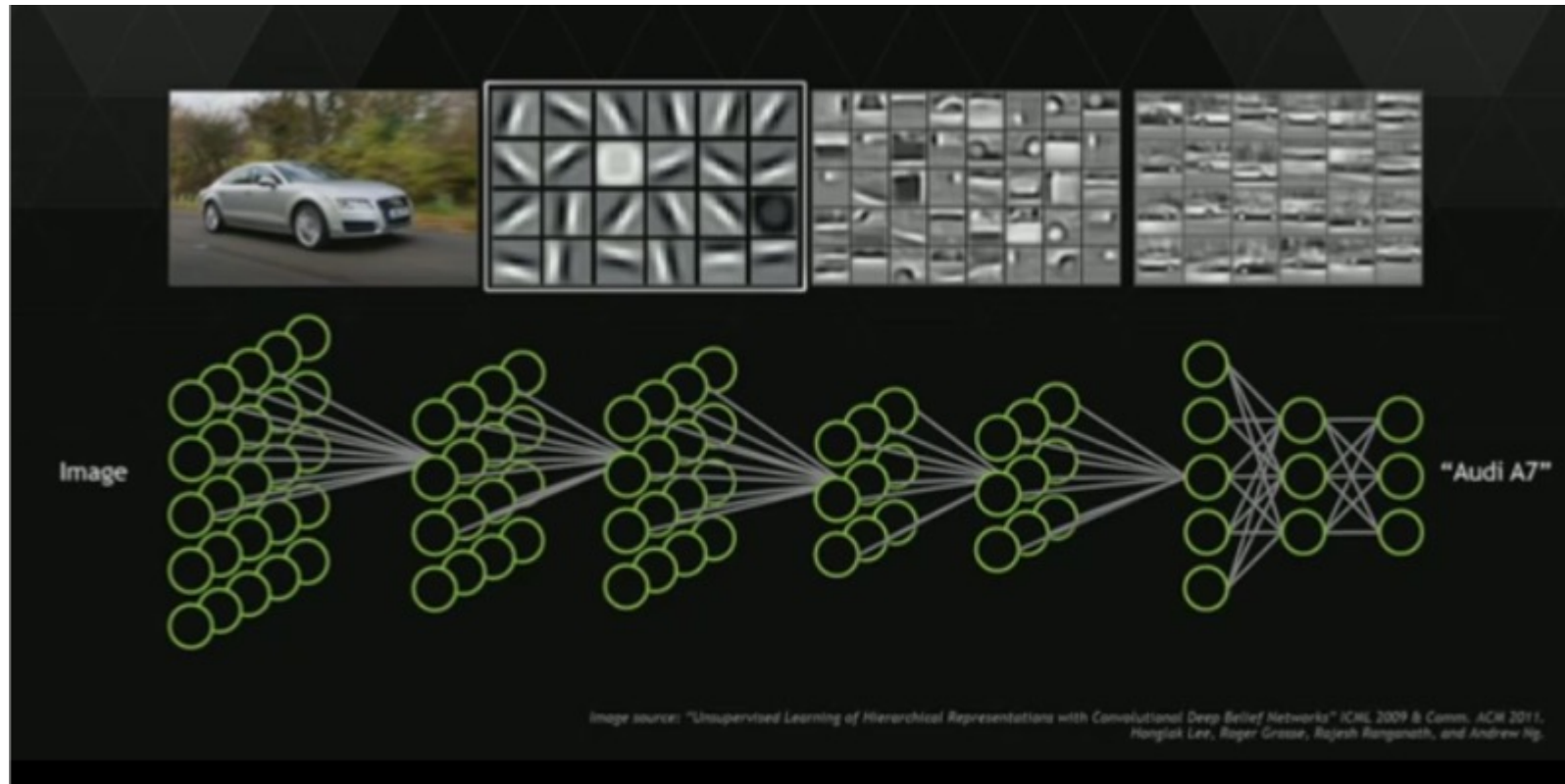
Motivation

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.

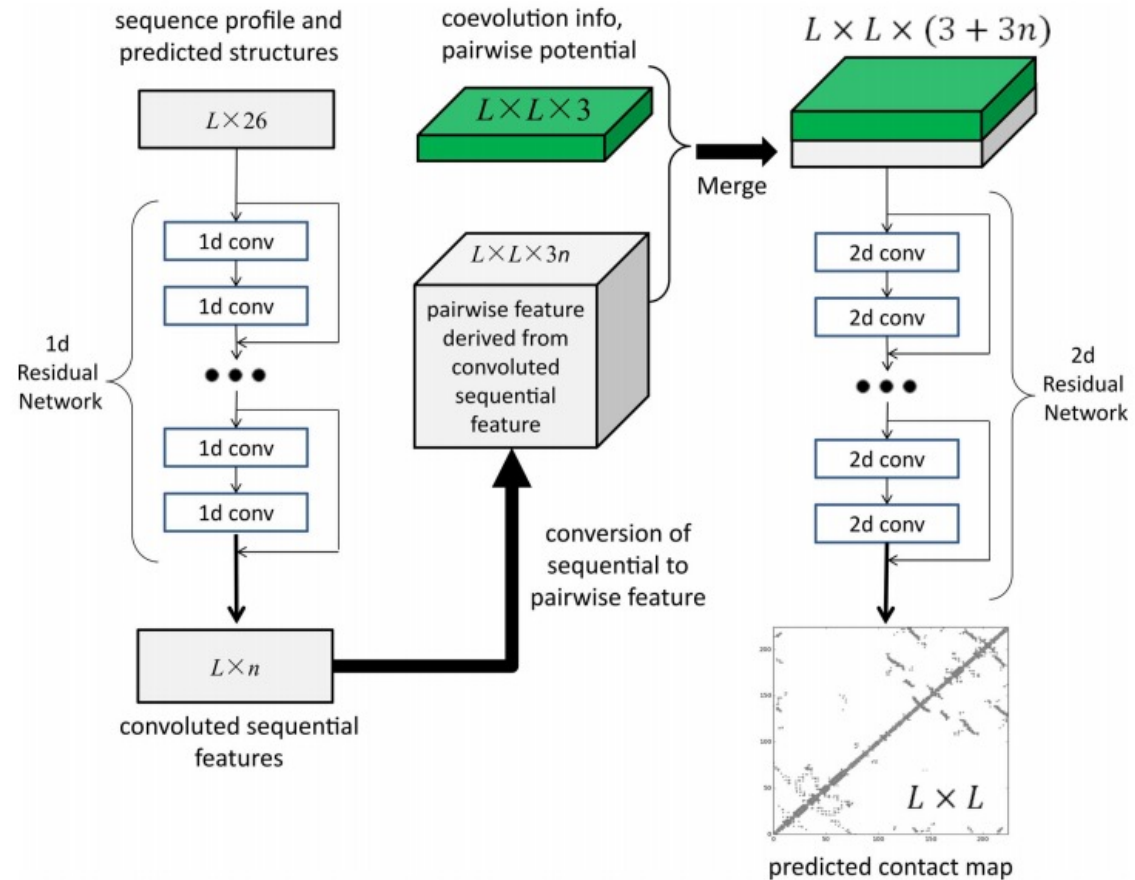
Neural networks



Convolutional neural networks



Learning a contact map from co-evolving residues



Inferring better contact maps (I)

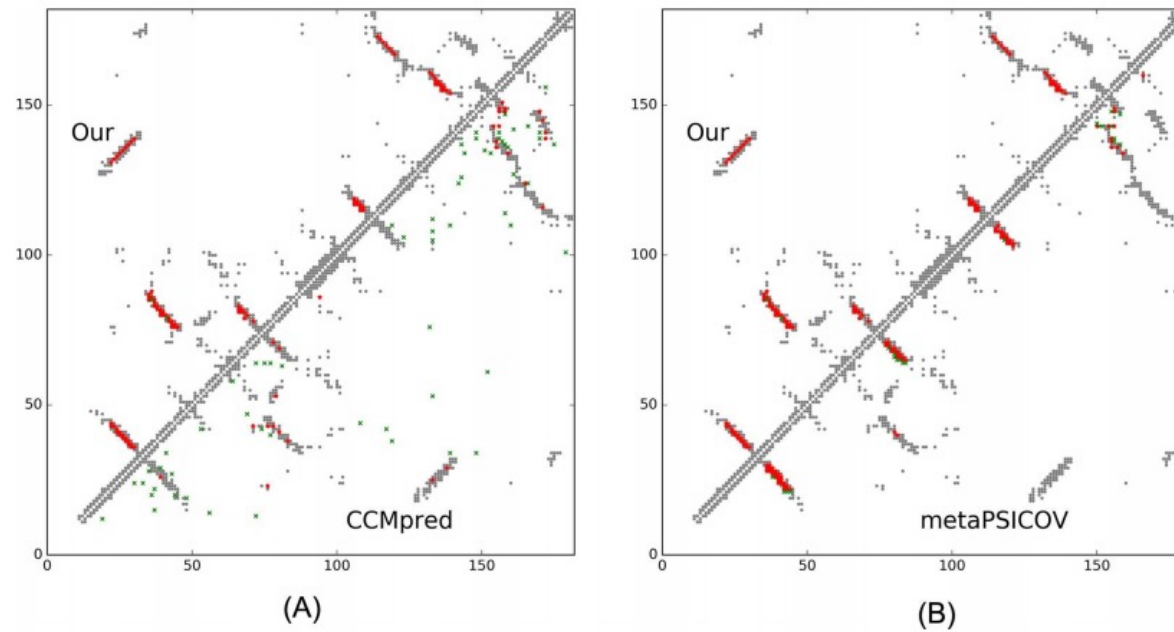
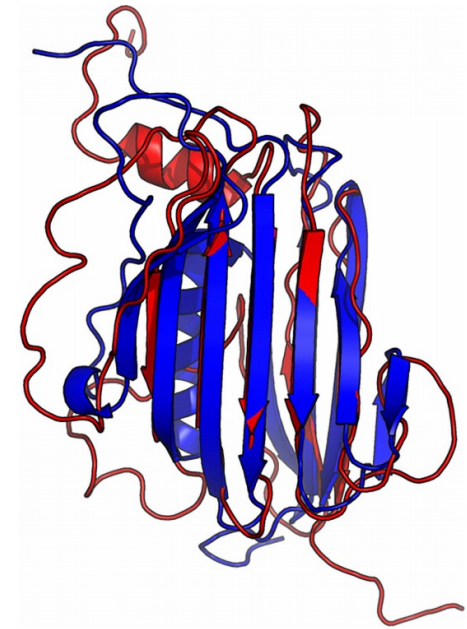


Fig 6. Overlap between top L/2 predicted contacts (in red or green) and the native contact map (in grey) for CAMEO target 2nc8A. Red (green) dots indicate correct (incorrect) prediction. (A) The comparison between our prediction (in upper-left triangle) and CCMpred (in lower-right triangle). (B) The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



Inferring better contact maps (II)

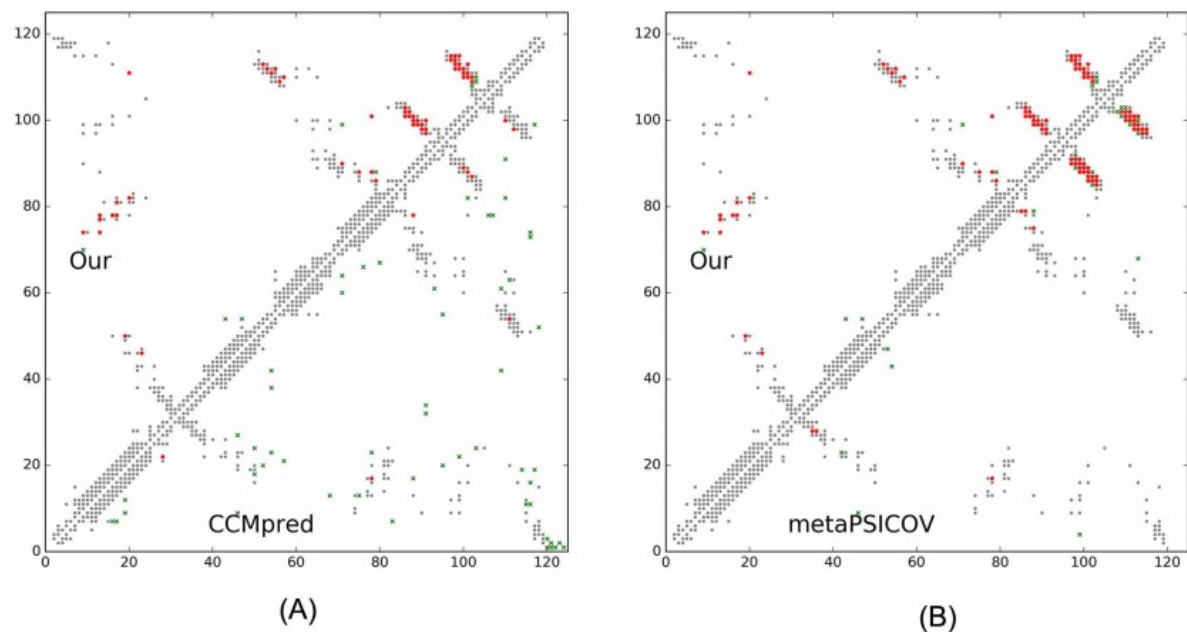
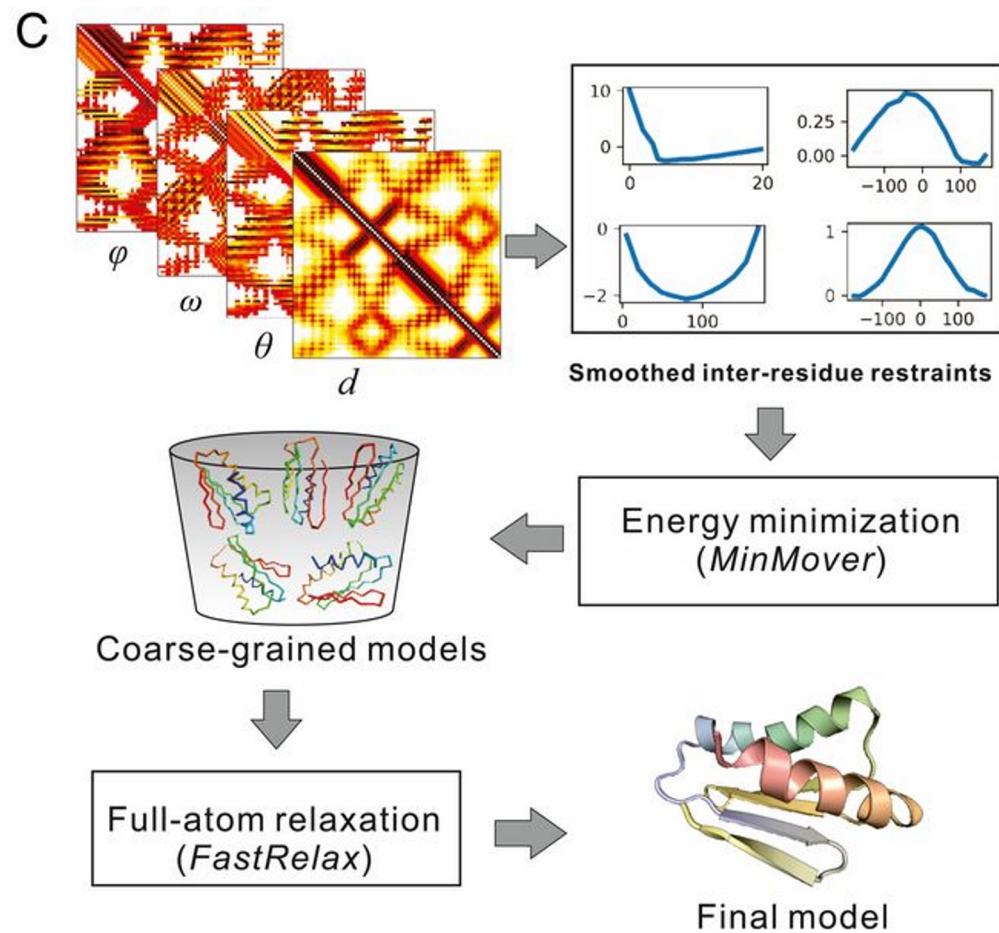
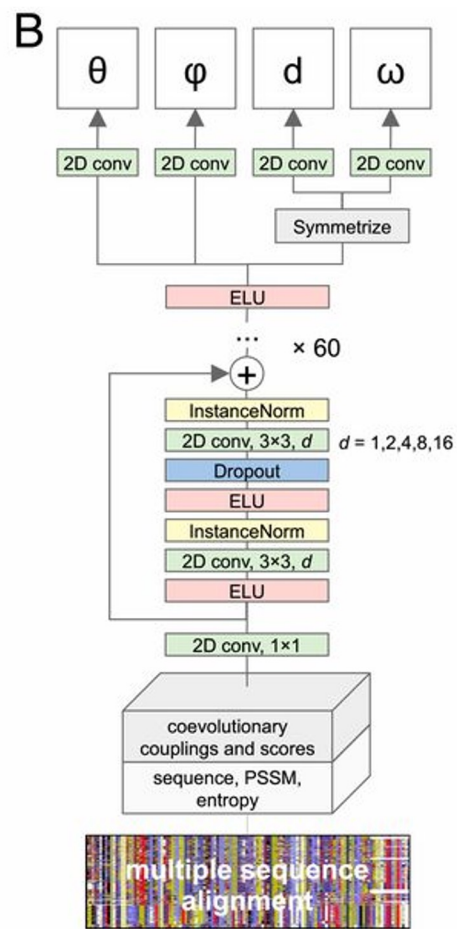


Fig 9. Overlap between top L/2 predicted contacts (in red or green) and the native contact map (in grey) for CAMEO target 5dcjA. Red (green) dots indicate correct (incorrect) prediction. (A) The comparison between our prediction (in upper-left triangle) and CCMpred (in lower-right triangle). (B) The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



trRosetta



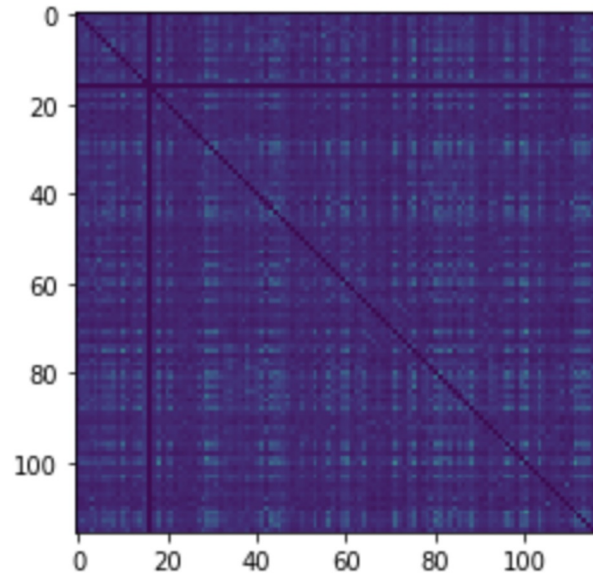
Improved protein structure prediction using predicted interresidue orientations

Jianyi Yang, Ivan Anishchenko, Hahnbeom Park, Zhenling Peng, Sergey Ovchinnikov, and David Baker

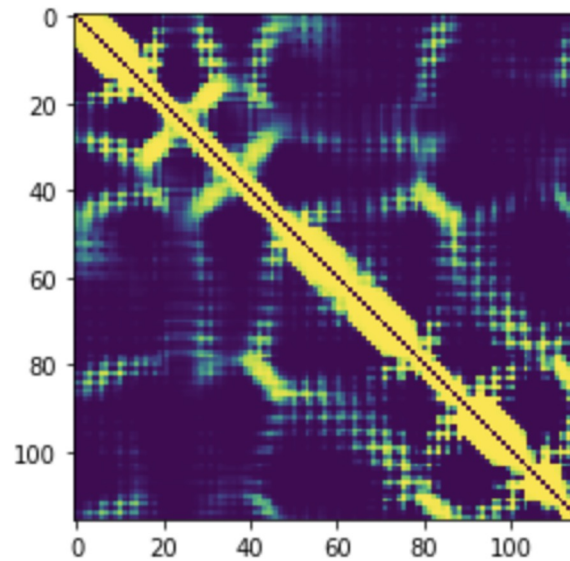
PNAS January 21, 2020 117 (3) 1496-1503; first published January 2, 2020 <https://doi.org/10.1073/pnas.1914677117>

Discovering hidden patterns with a learned model

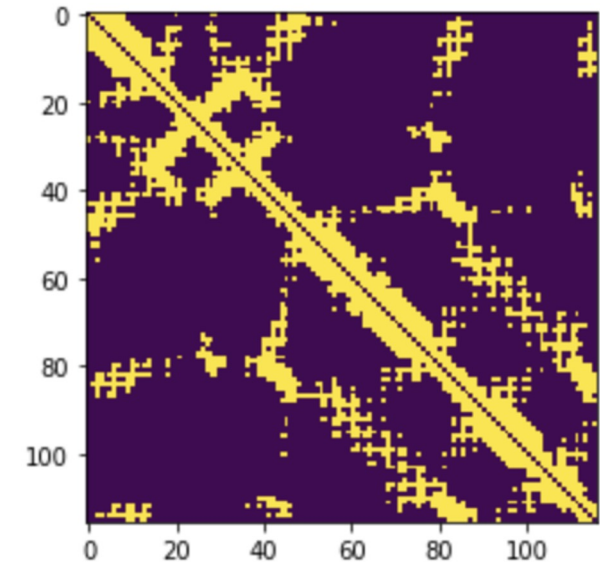
**Gremlin predictions
on shallow MSAs**
(Nseq=36, Nf=2.3)



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

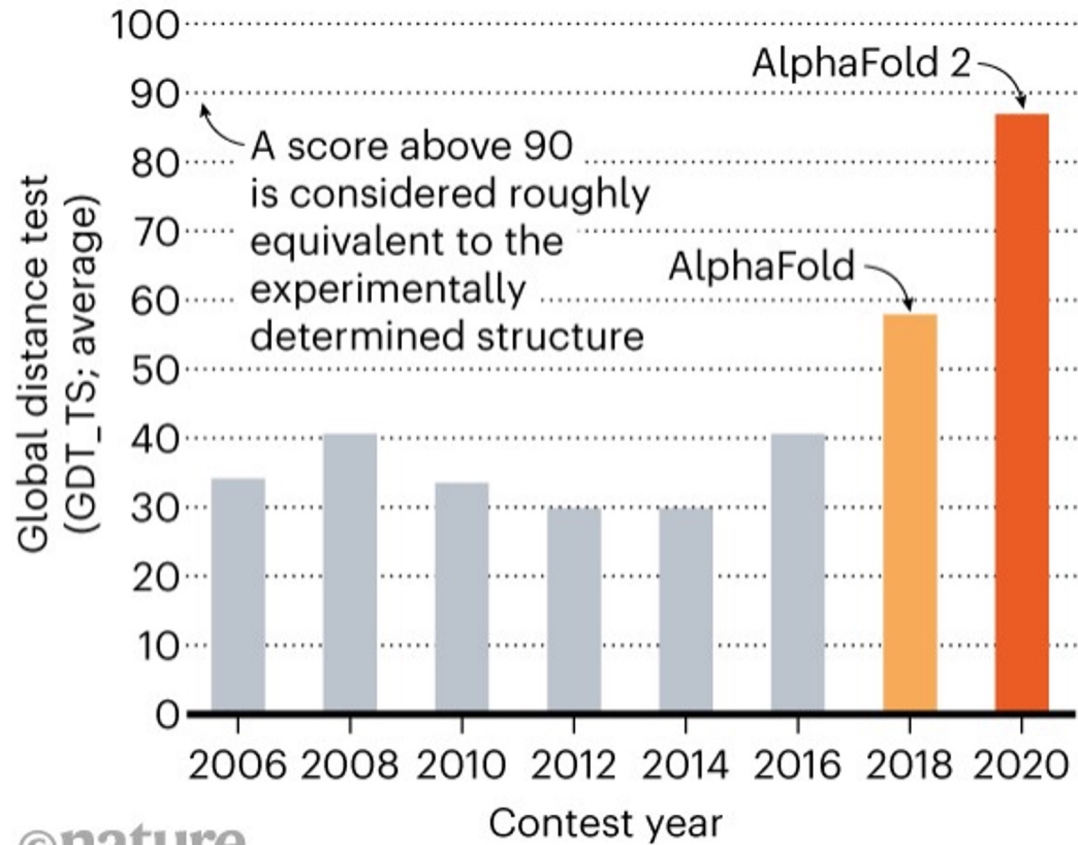


Native contact map



Improving protein structure prediction

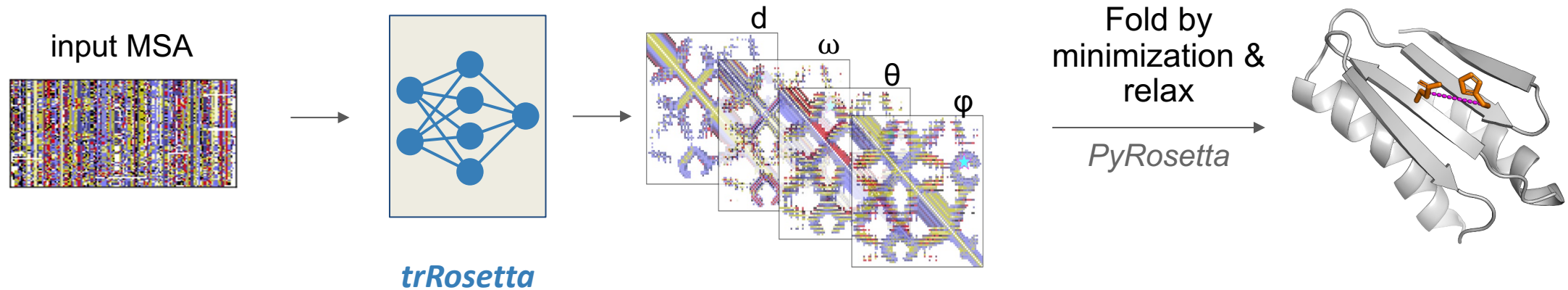
Free modeling accuracy in CASP



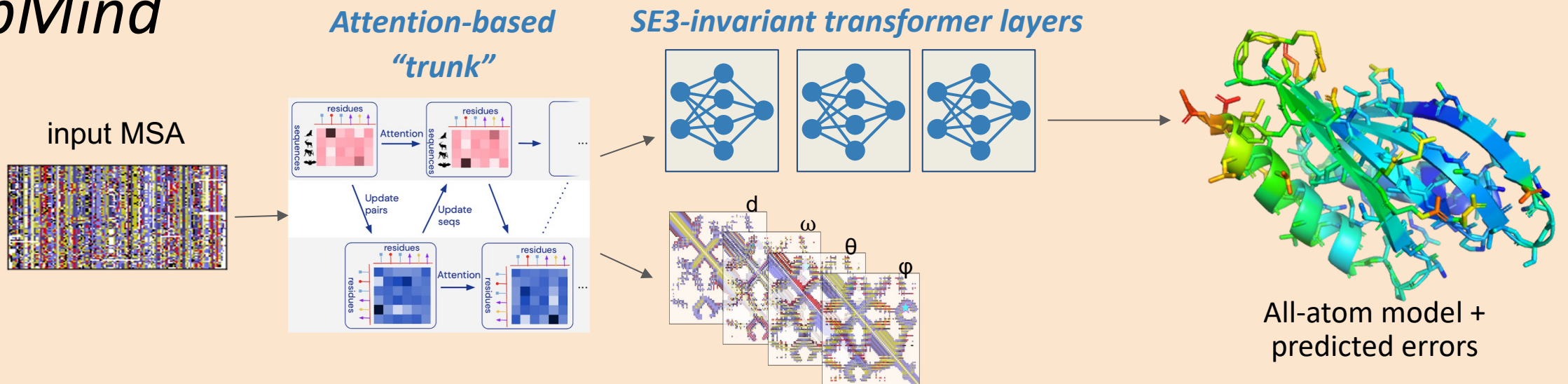
©nature

A differentiable end-to-end structure predictor

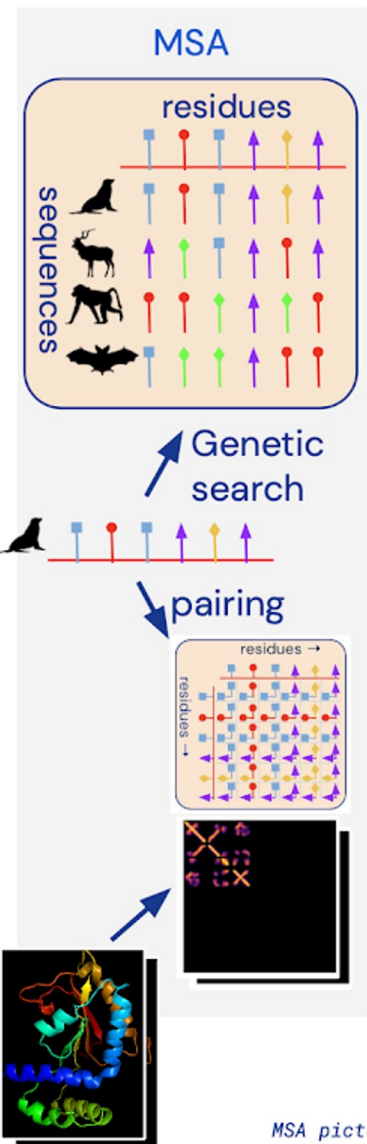
trRosetta



DeepMind



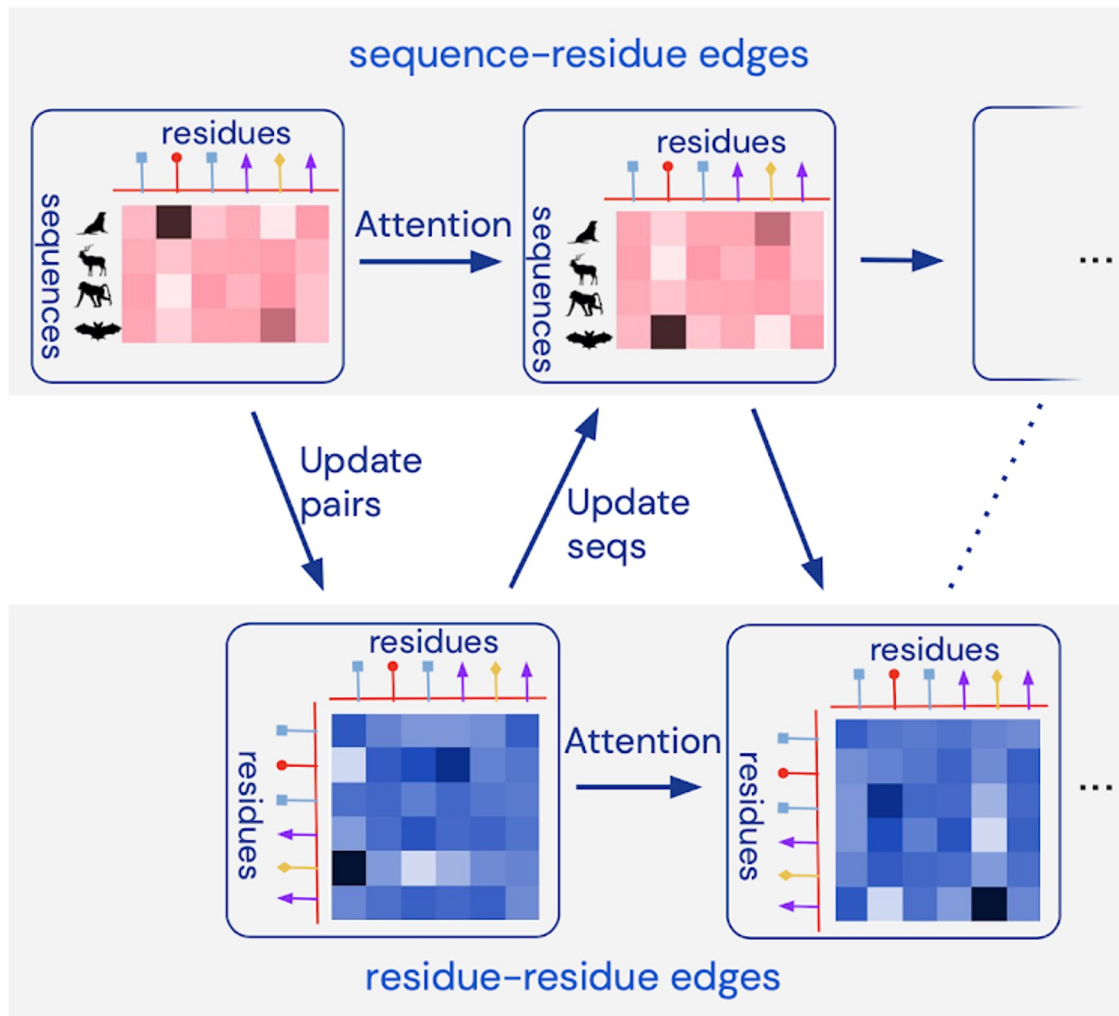
Embedding



templates

MSA picture inspired by: Riesselman, A.J., Ingraham, J.B. & Marks, D.S., Nature Methods (2018) doi:10.1038/s41592-018-0138-4

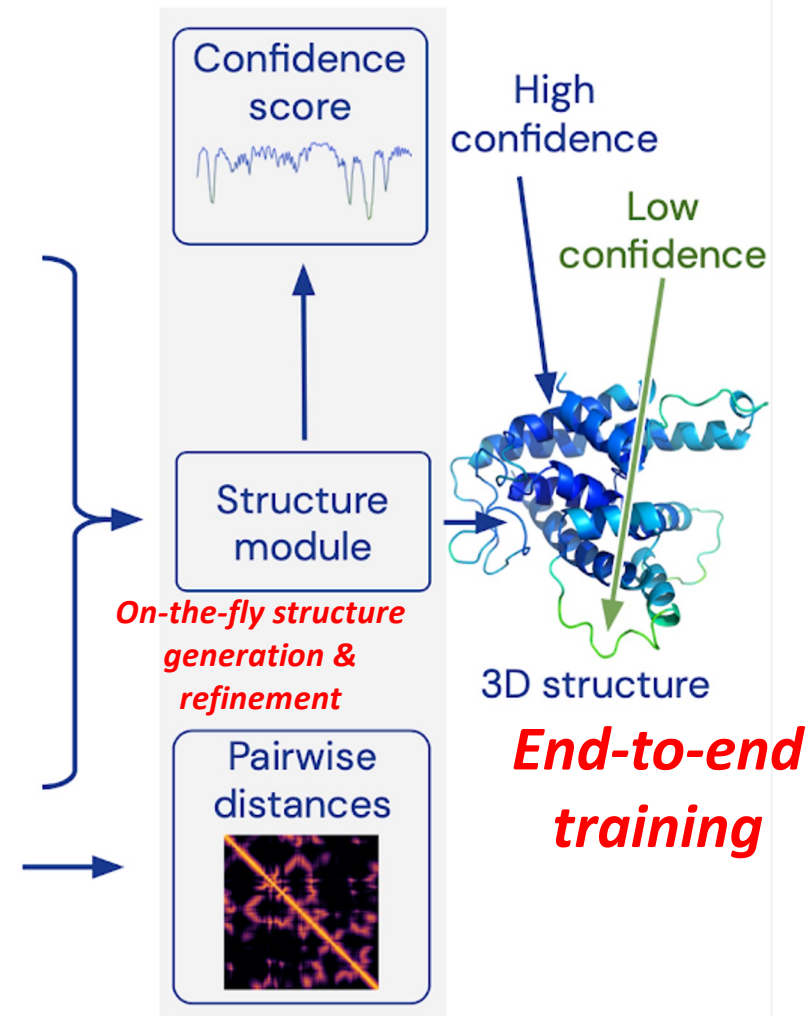
Trunk



**Iterative feature extraction
(attention instead of convolution)**

Heads

© 2020 DeepMind Technologies Limited

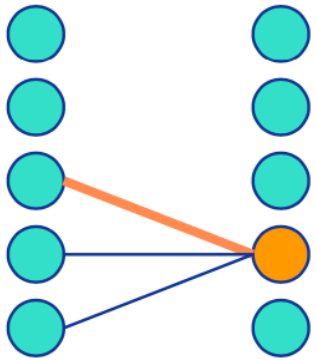


**On-the-fly structure
generation &
refinement**

**End-to-end
training**

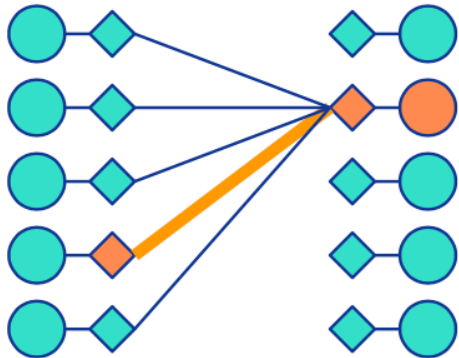


What would be a proper inductive bias for protein structure prediction?



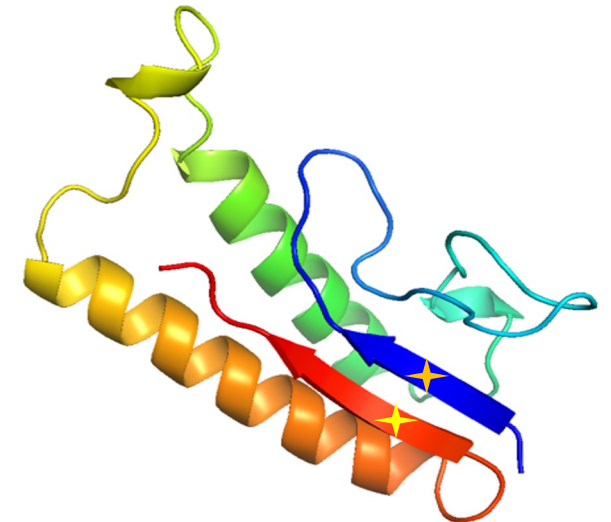
Convolutional Networks (e.g. computer vision)

- data in regular grid
- information flow to local neighbours

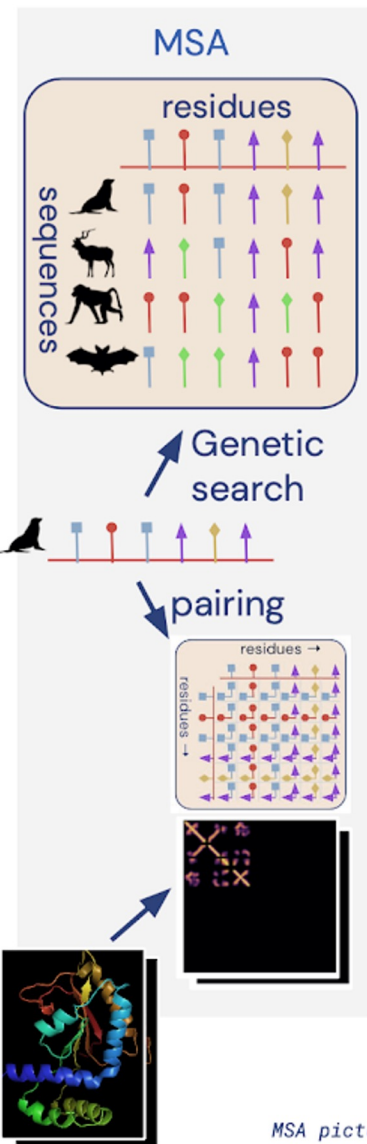


Attention Module (e.g. language)

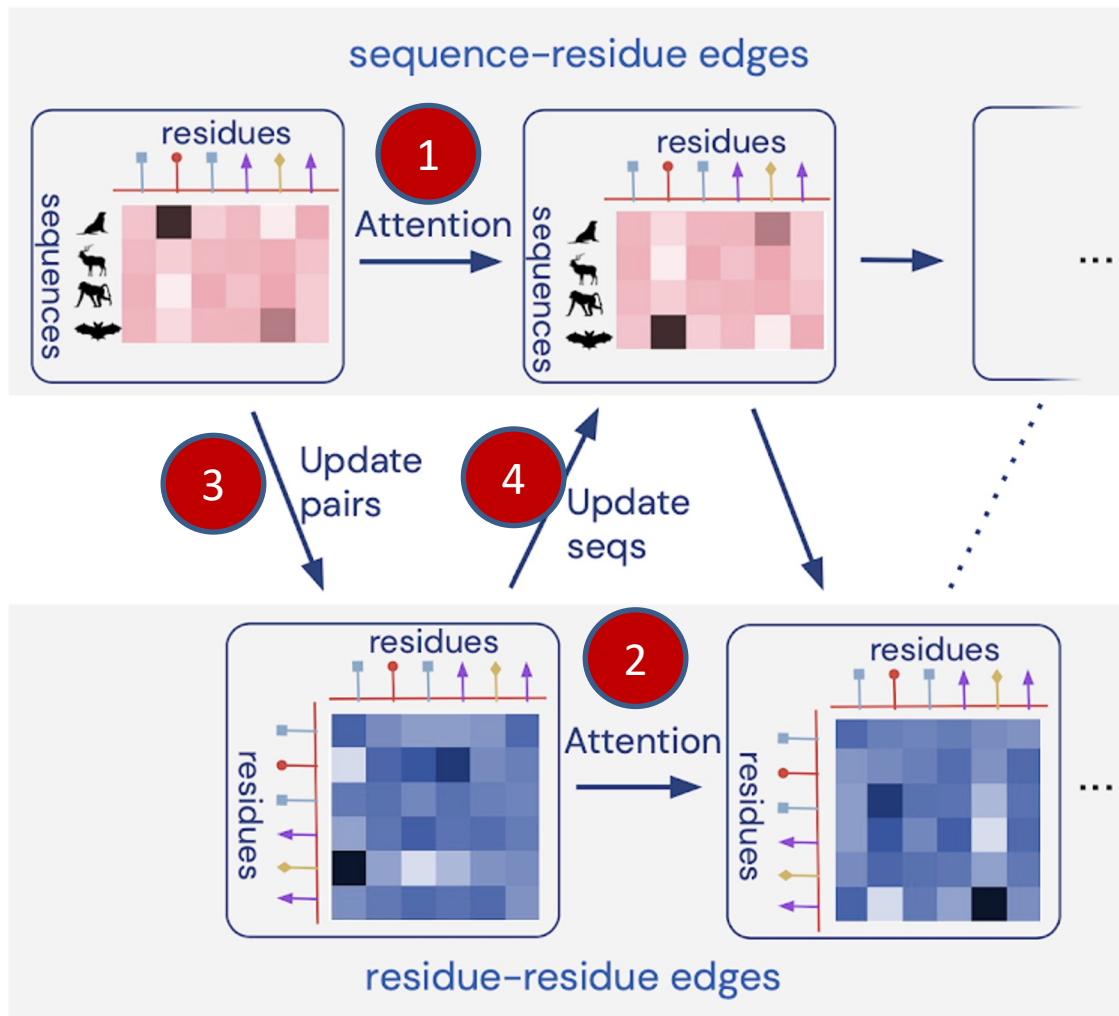
- data in unordered set
- information flow dynamically controlled by the network (via keys and queries)



Embedding

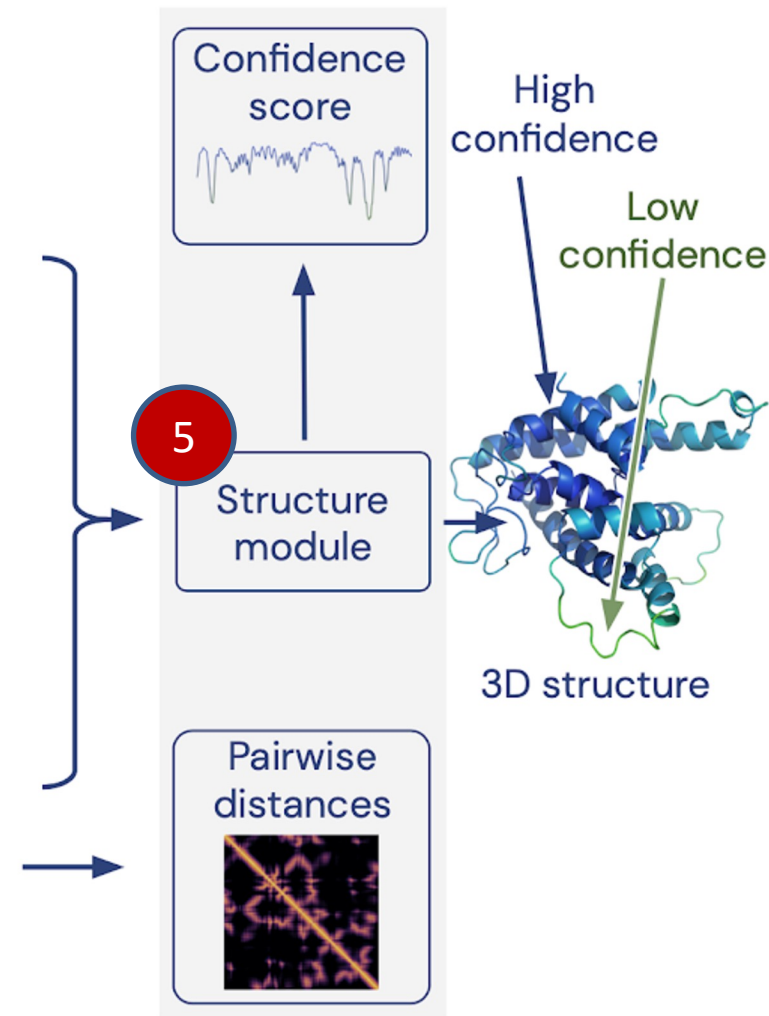


Trunk



Heads

© 2020 DeepMind Technologies Limited

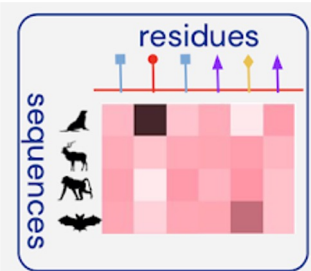
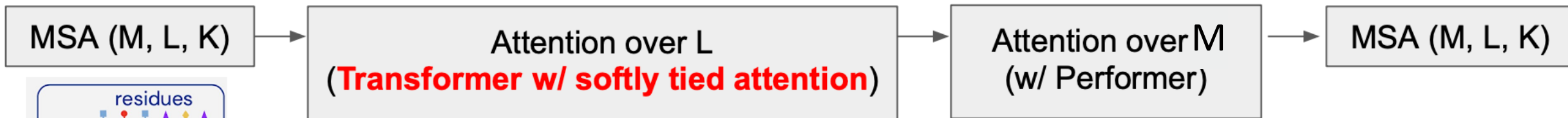


MSA picture inspired by: Riesselman, A.J., Ingraham, J.B. & Marks, D.S., Nature Methods (2018) doi:10.1038/s41592-018-0138-4

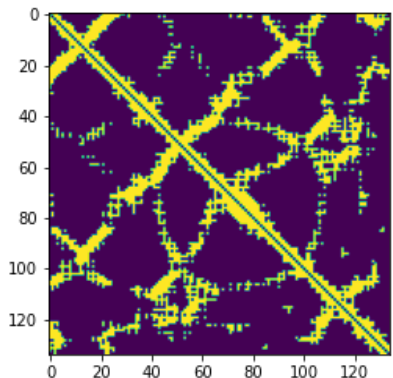
templates



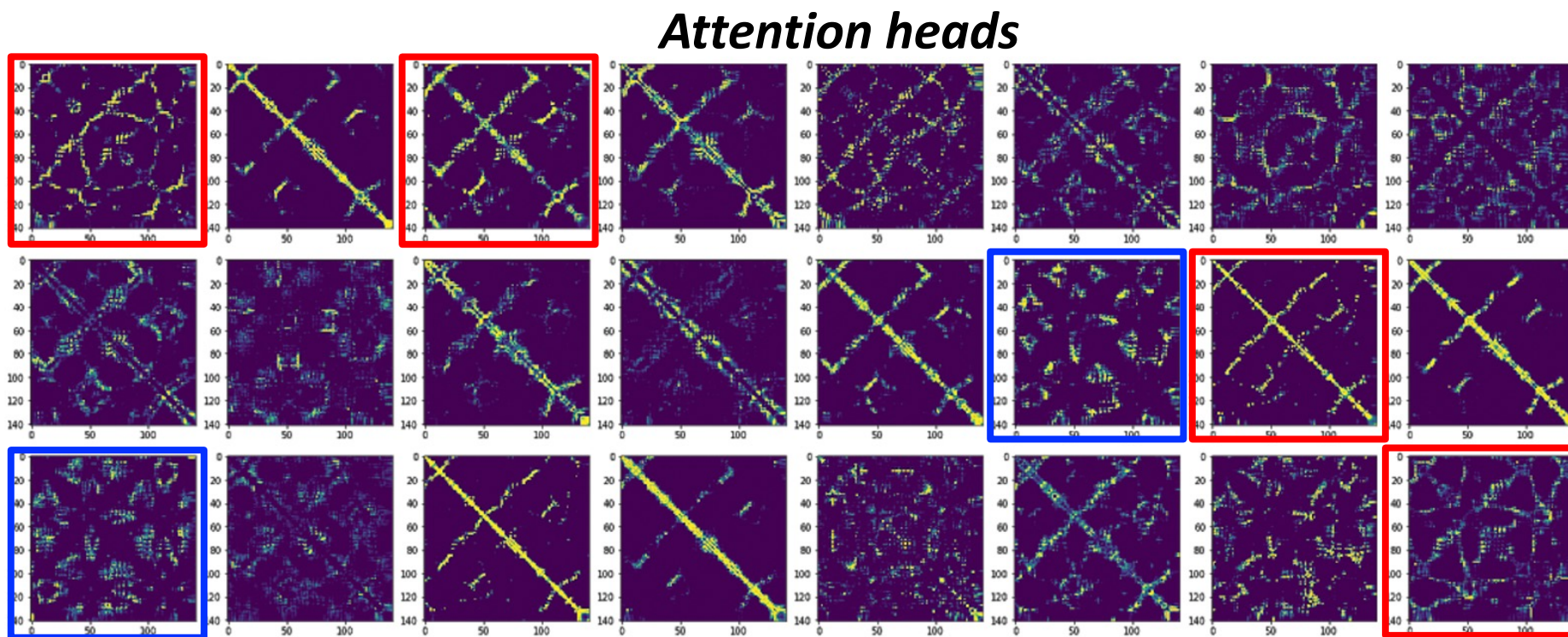
Component 1: MSA updates via self-attention



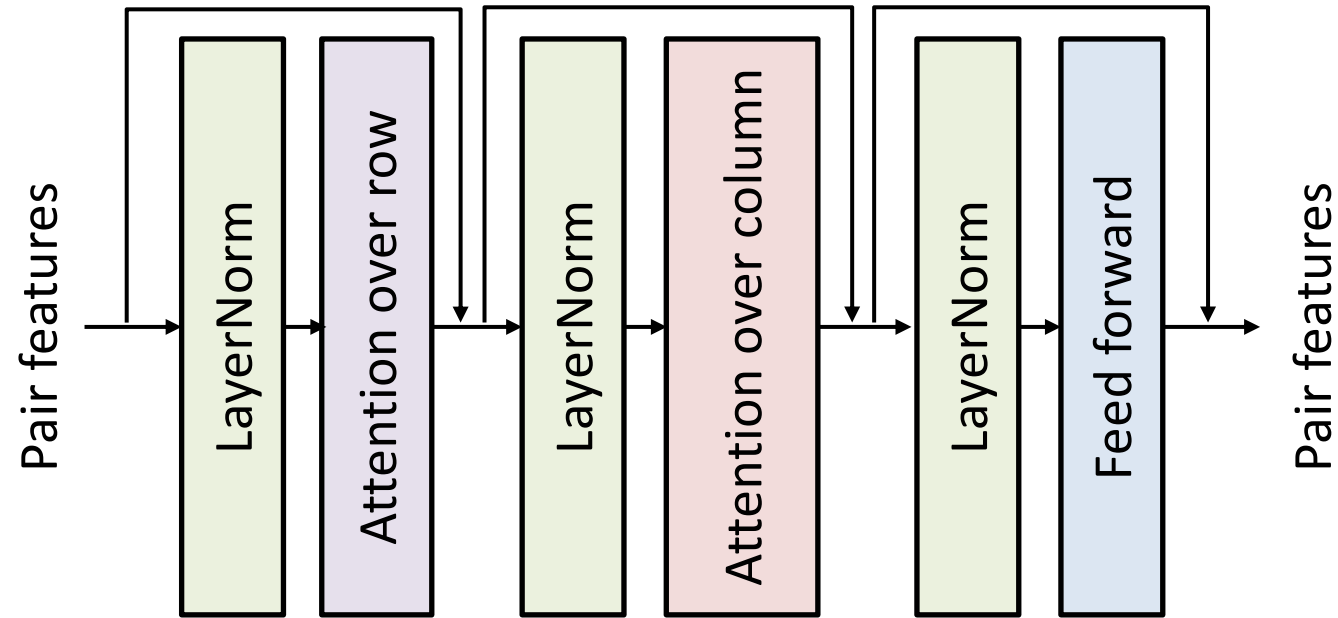
True contact map



Last 3 blocks

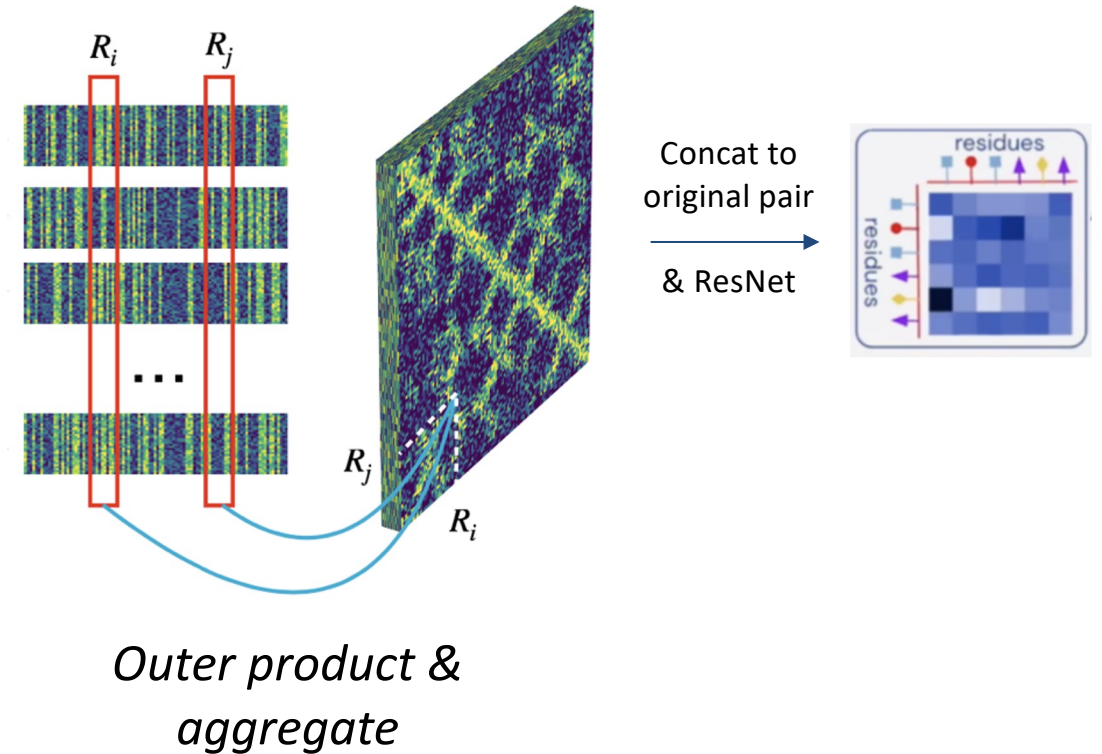
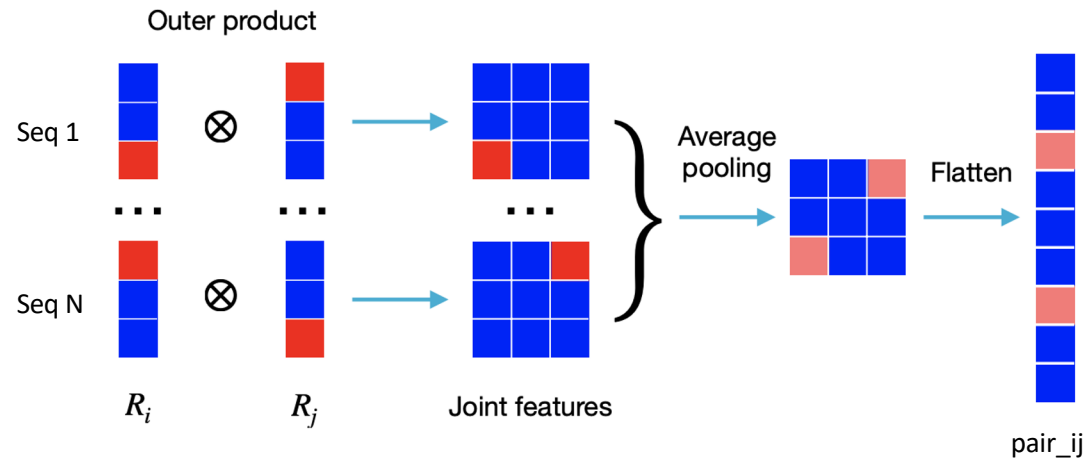


Component 2: Update pair features via self-attention



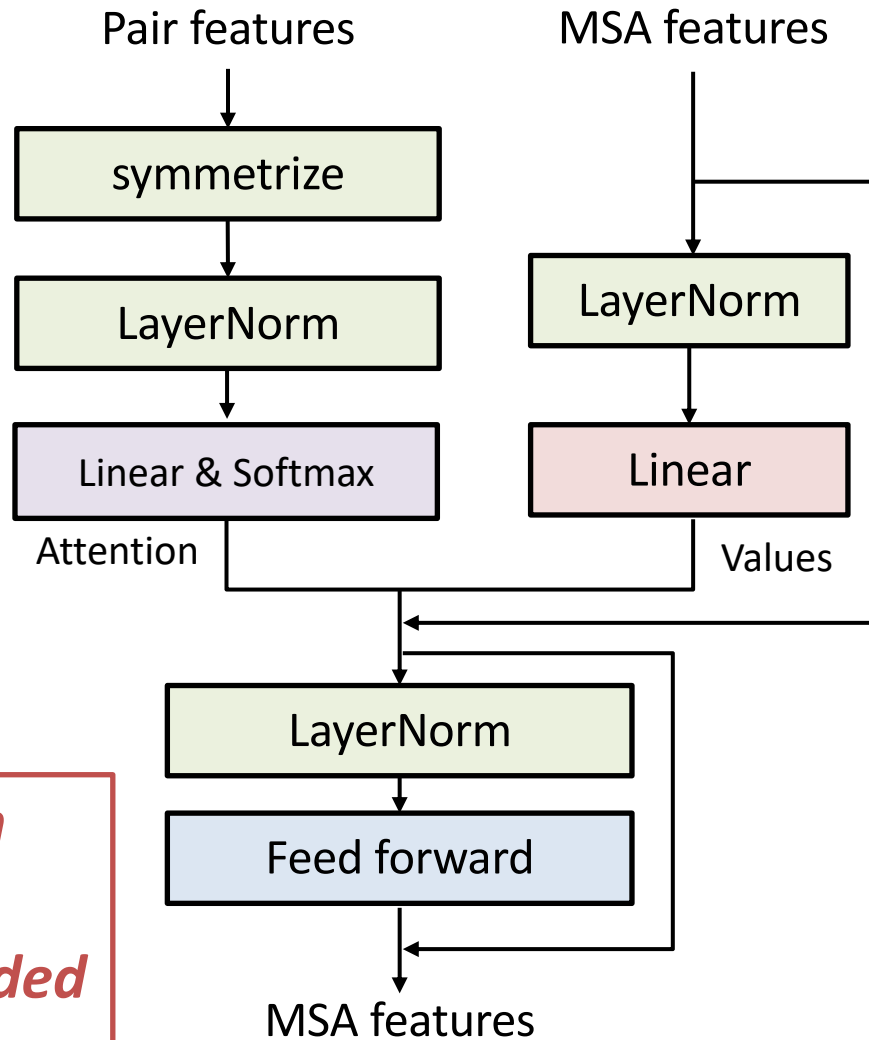
*Axial Attention (attention over rows then columns)
to reduce memory requirements & computation time*

Component 3: Extract pair features from MSA



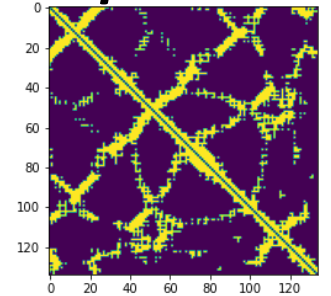
Non-interacting pairs \rightarrow Broader distribution
Interacting pairs (co-mutating) \rightarrow Sharper distribution

Component 4: Update MSA based on pair features

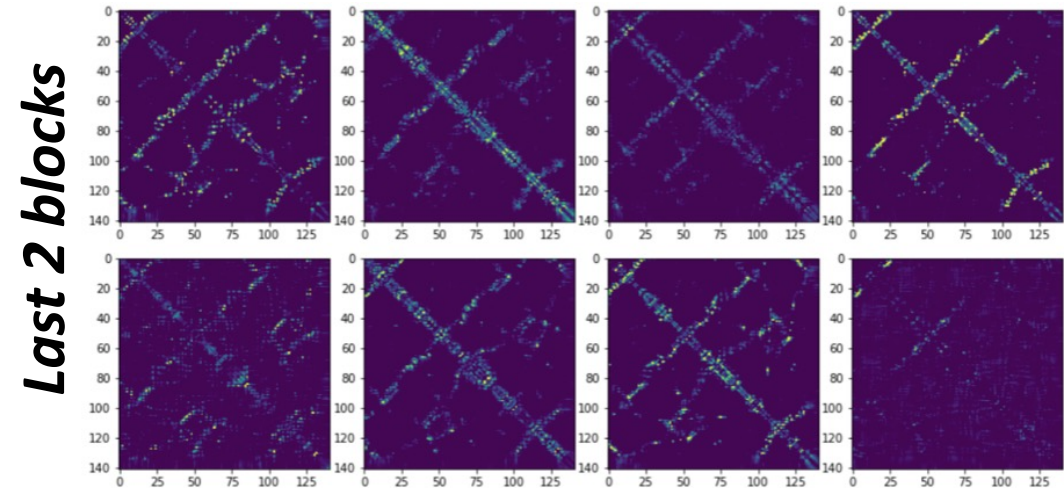


Attention from structure information encoded in pair features

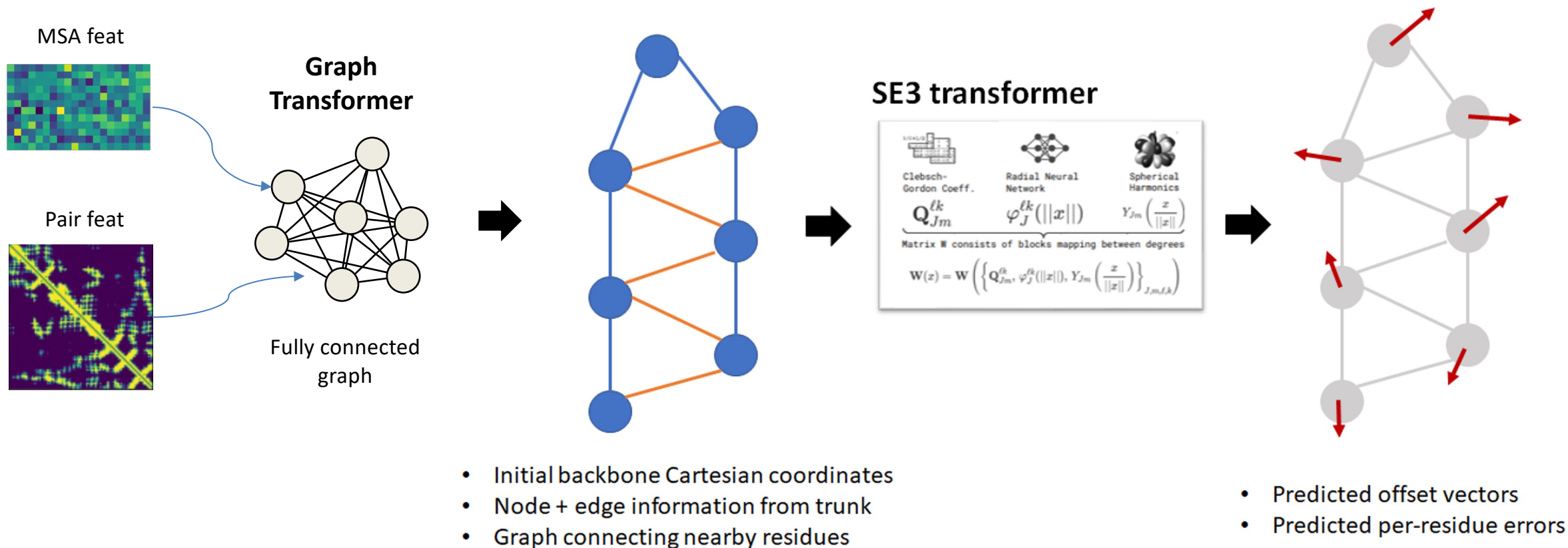
True contact map



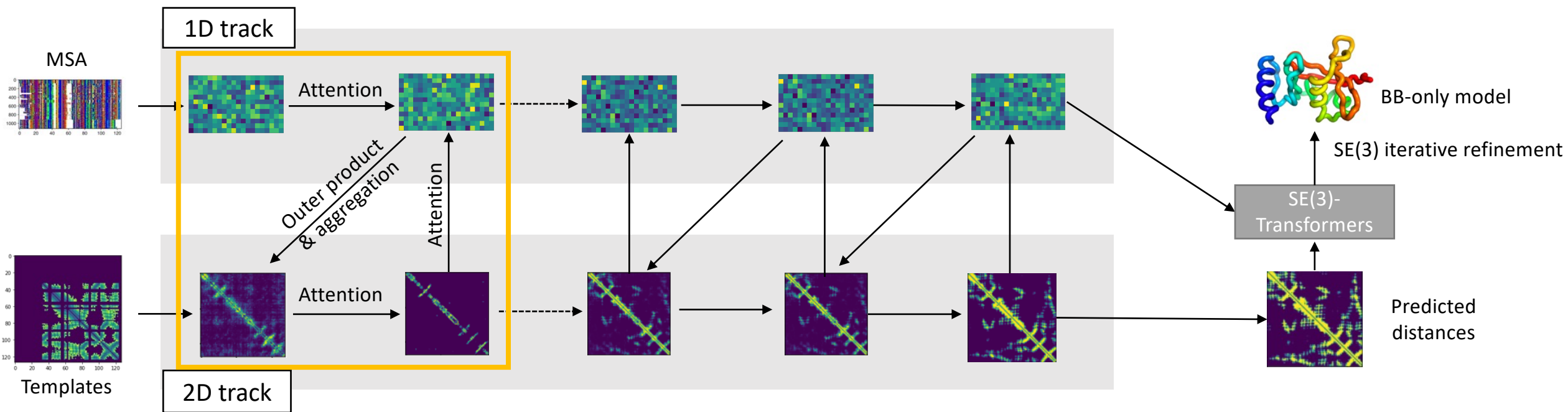
Attention heads



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



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?

