

Protein structure prediction: the past, the present, and the future

Dec. 11, 2022

CASP15 meeting

Minkyung Baek

minkbaek@snu.ac.kr

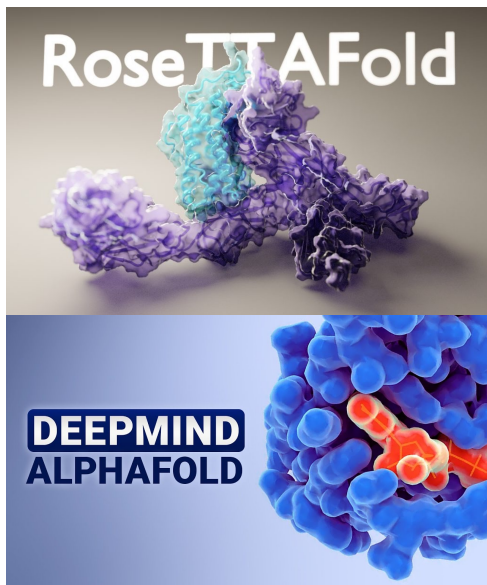


INSTITUTE FOR
Protein Design

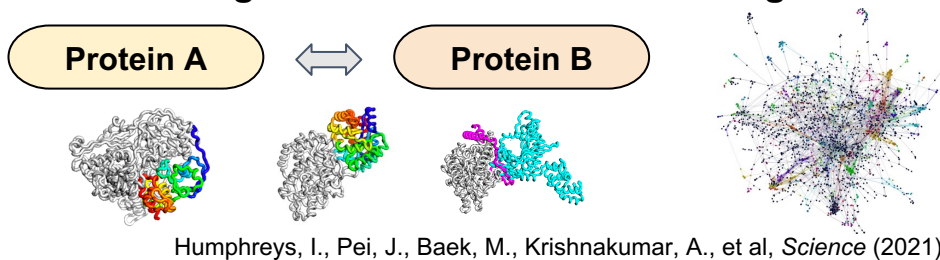
UNIVERSITY *of* WASHINGTON



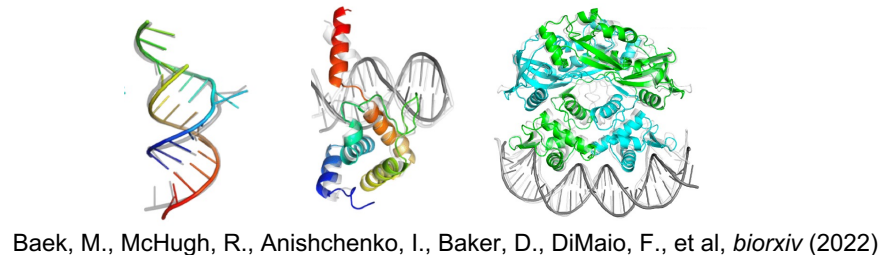
AI-based protein modeling



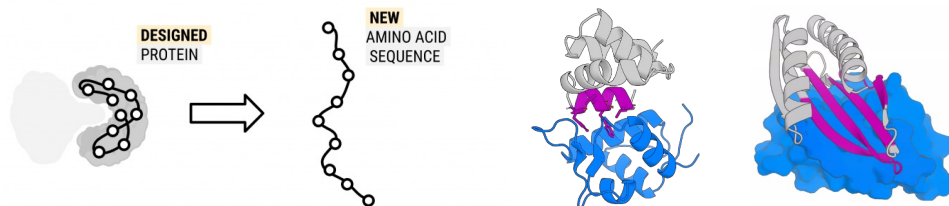
Large-scale *in silico* PPI screening



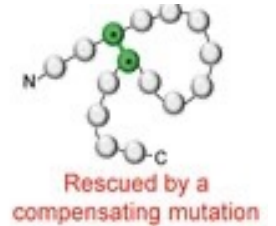
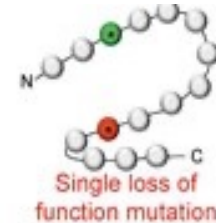
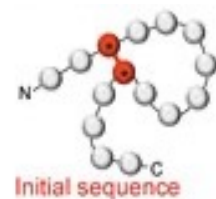
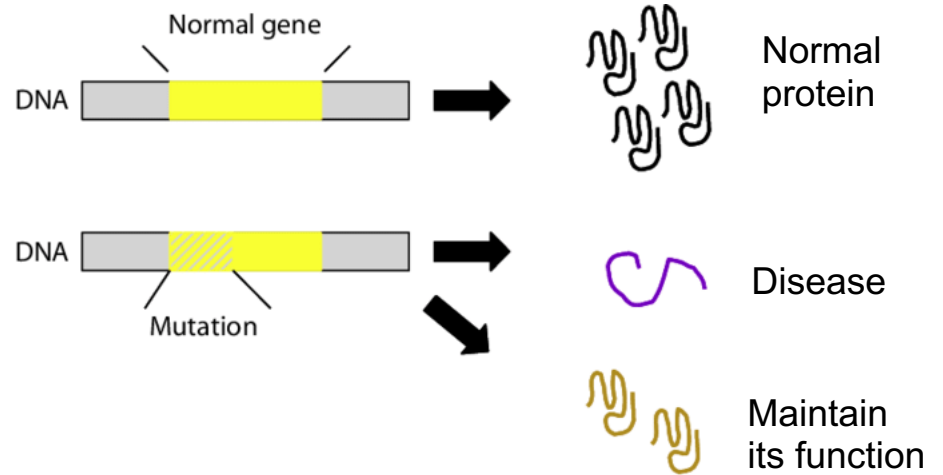
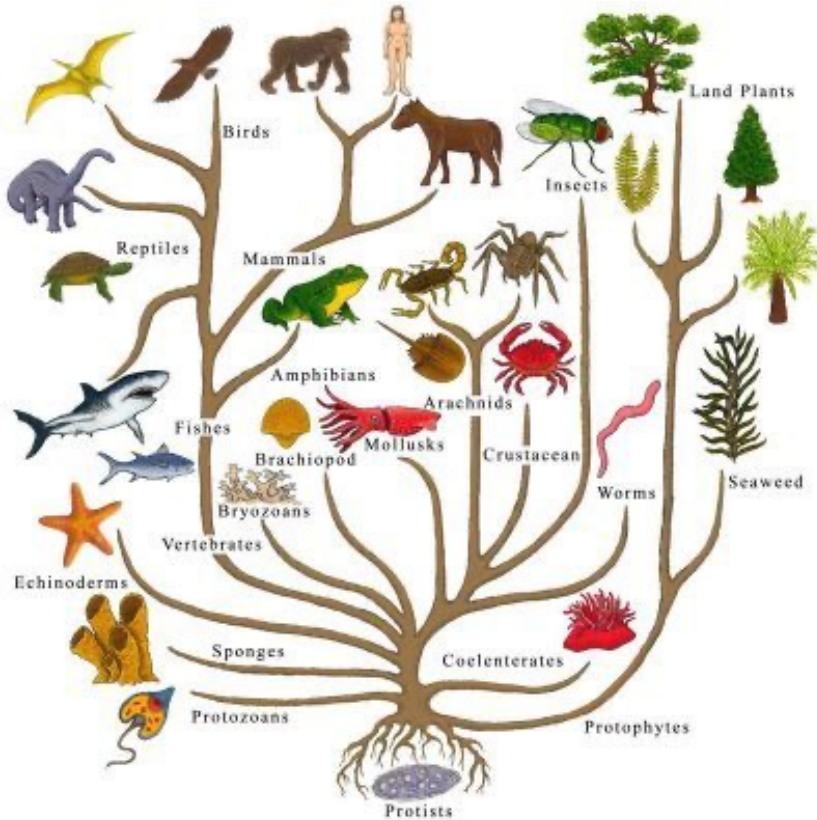
Nucleic acid structure & interaction prediction



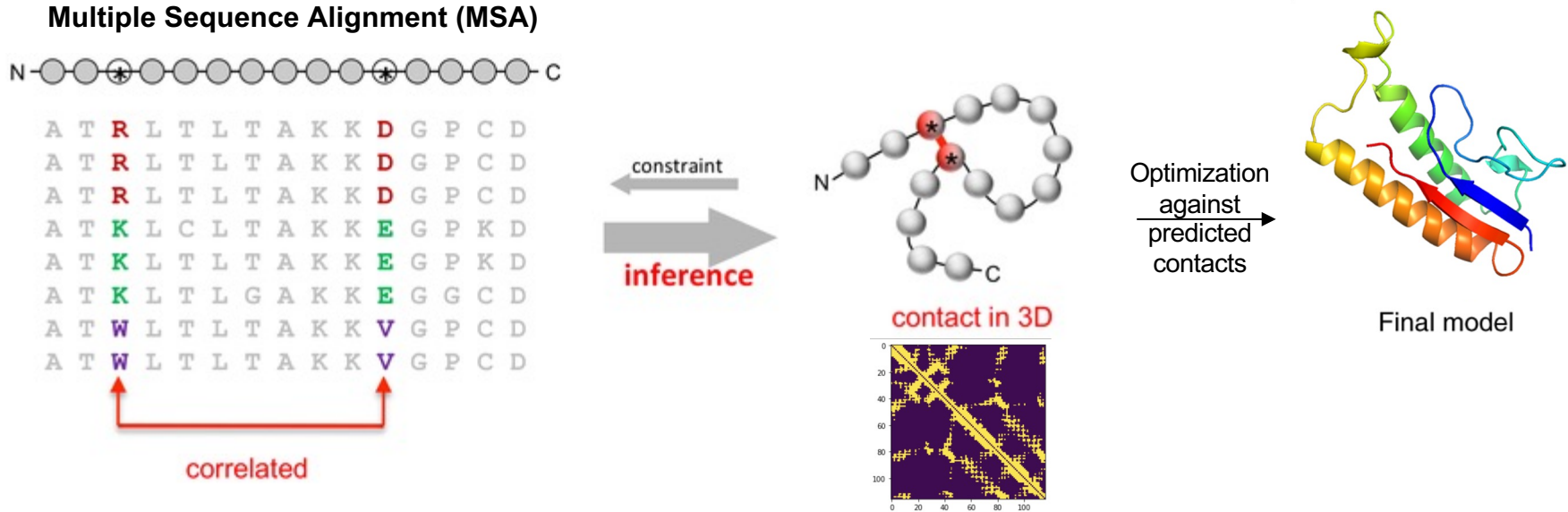
De novo functional protein design



Protein structure prediction using evolution history



Coevolution guided modeling



Q: How can we find the coevolution pattern from MSA?

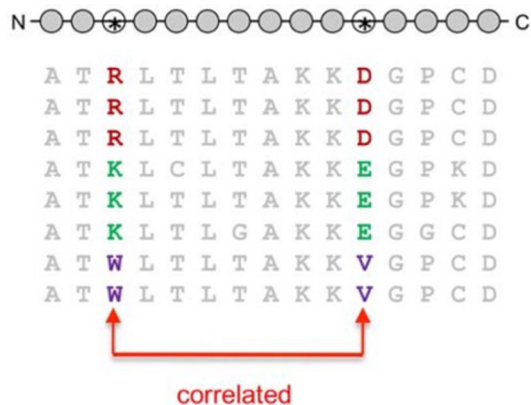
Applying deep learning to protein structure prediction



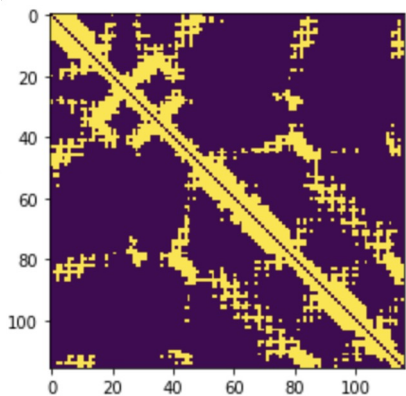
- Deep learning algorithms attempt to learn (multiple levels of) representation by using a **hierarchy of multiple layers**
- Exceptional effective at **learning patterns**
- If you provide the system **tons of training examples**, it begins to understand hidden patterns in data and respond in useful ways

From MSA to 3D structures with deep learning

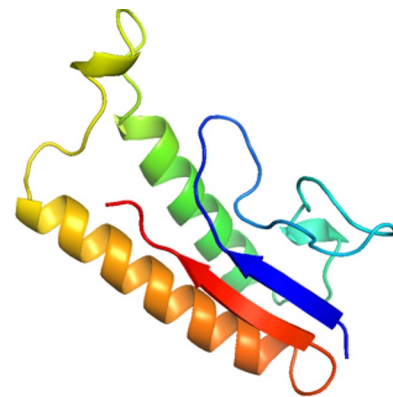
Multiple sequence alignments



Residue pairwise interaction

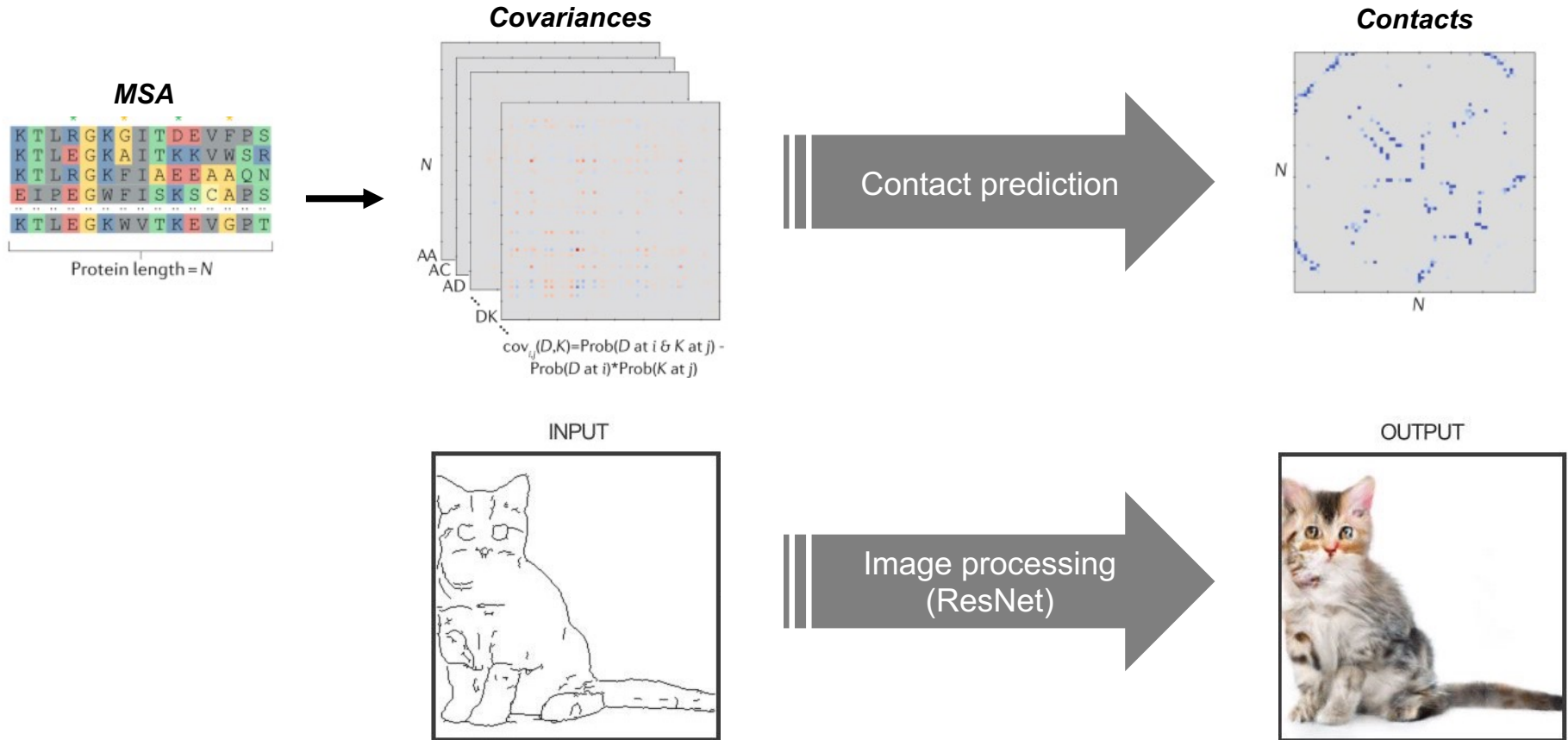


3D atomic coordinates

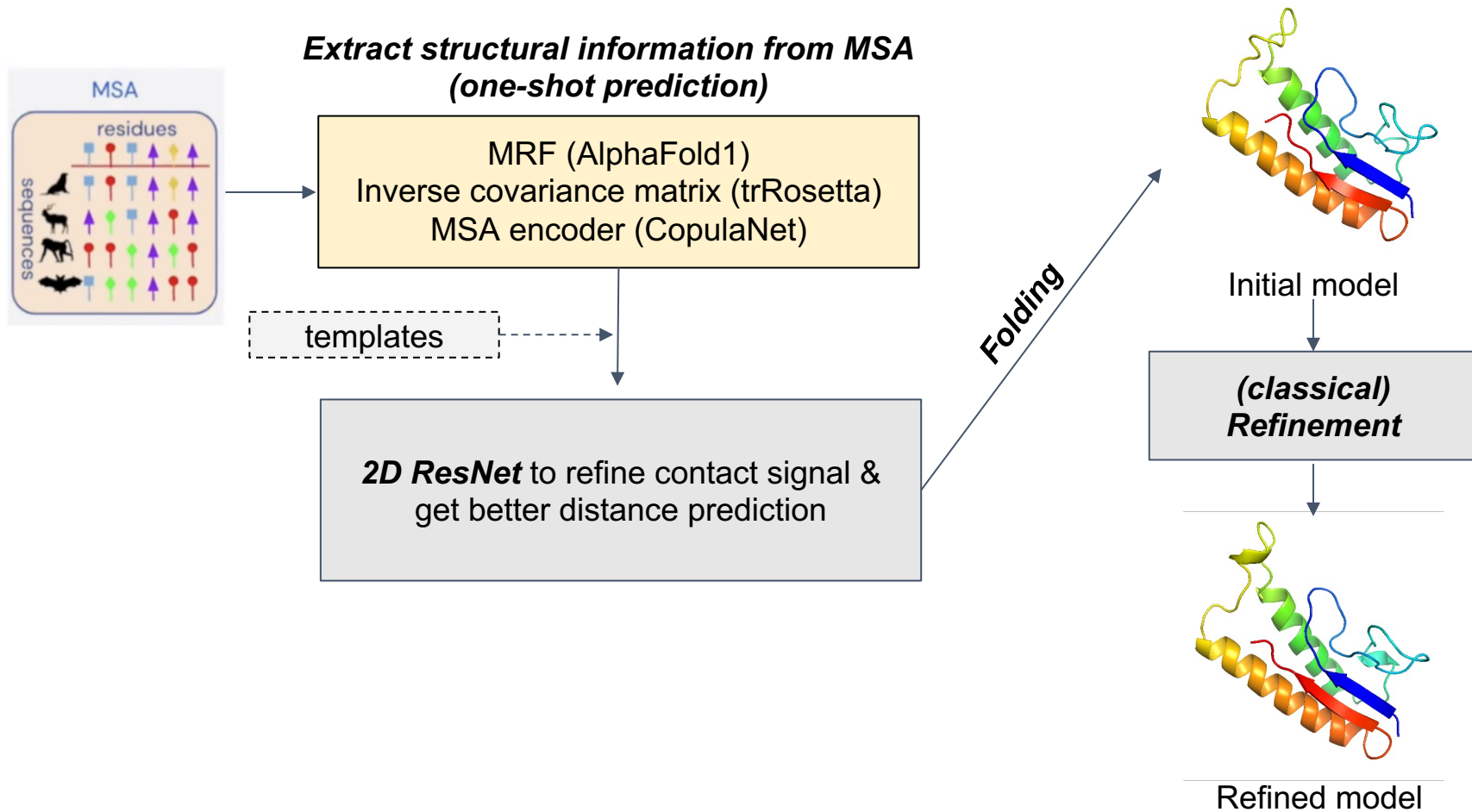


How to extract residue pairwise interactions & build 3D models based on that?

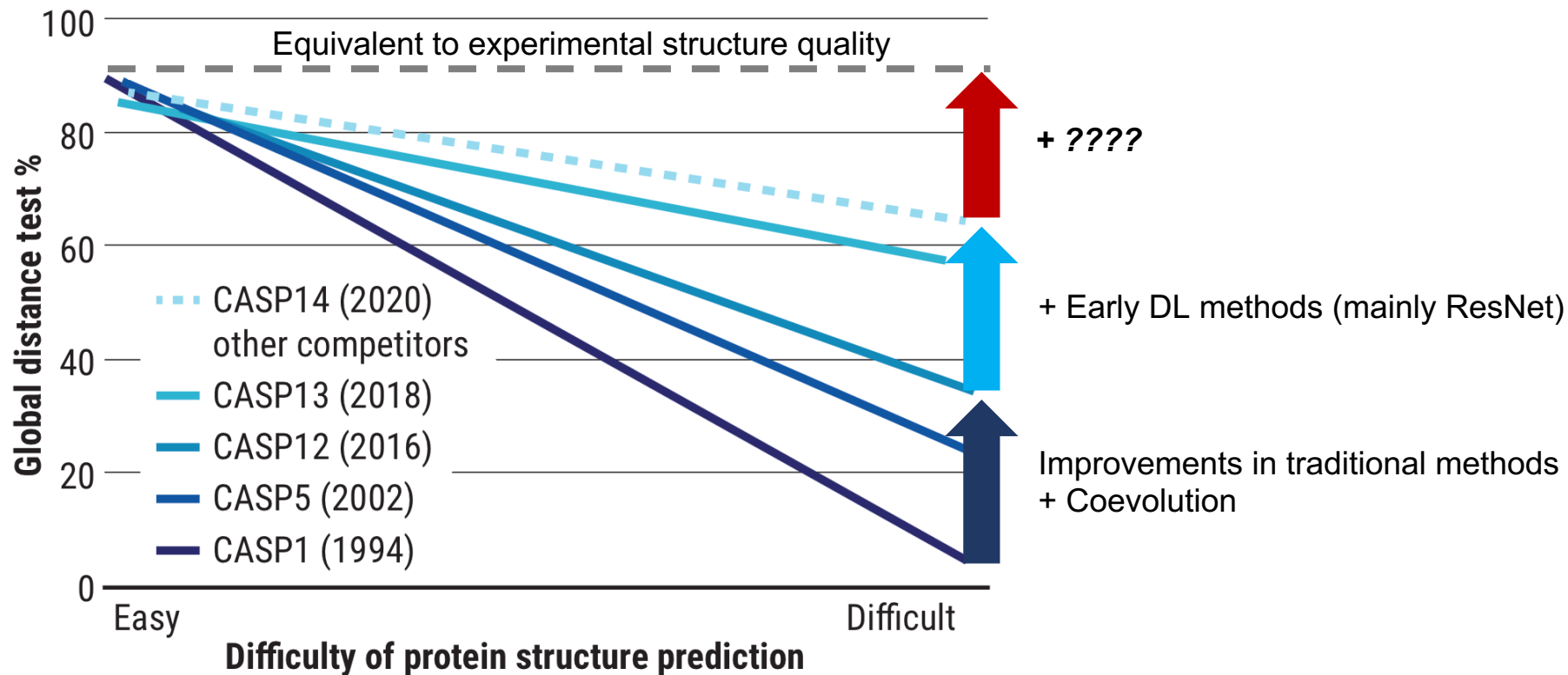
Early attempts w/ DL: ResNet-based approach



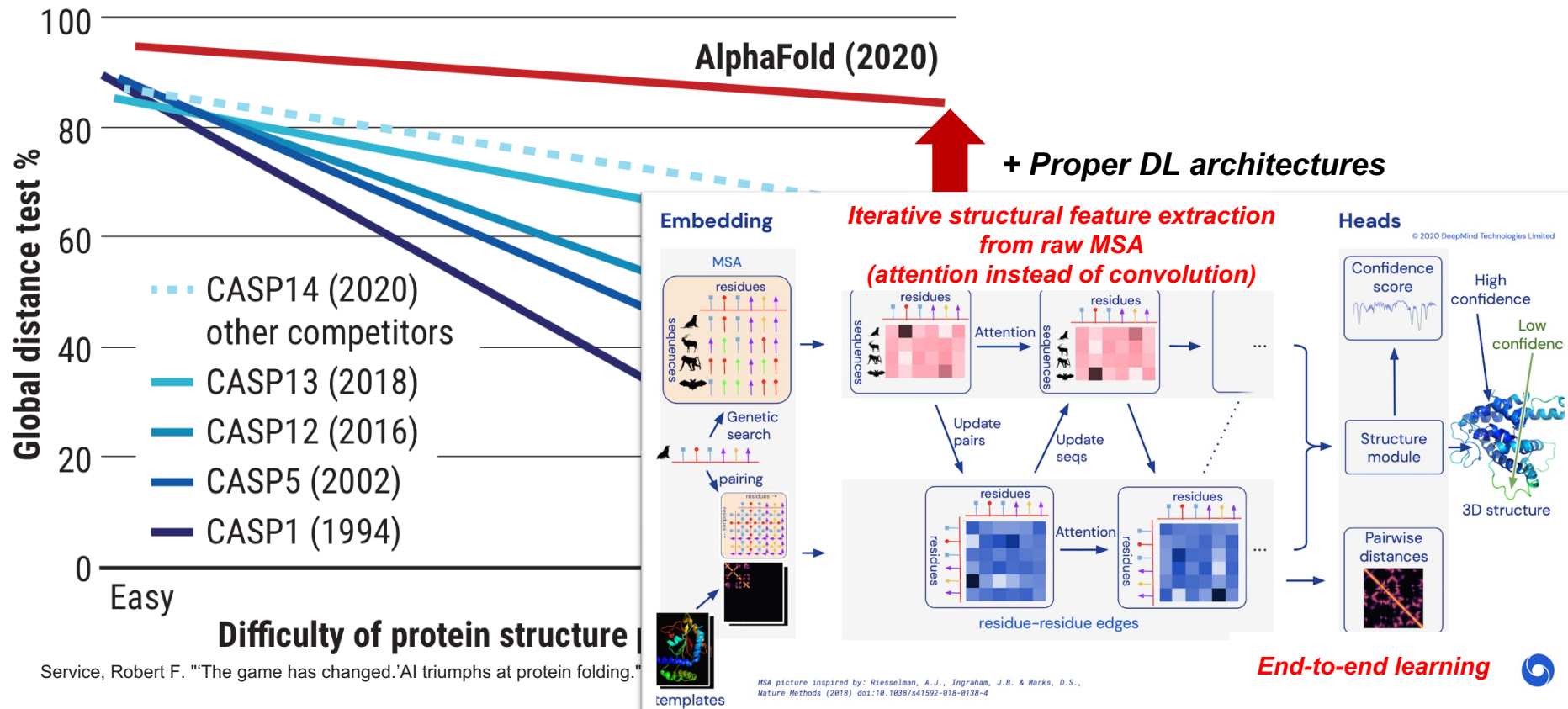
Early attempts w/ DL: ResNet-based approach



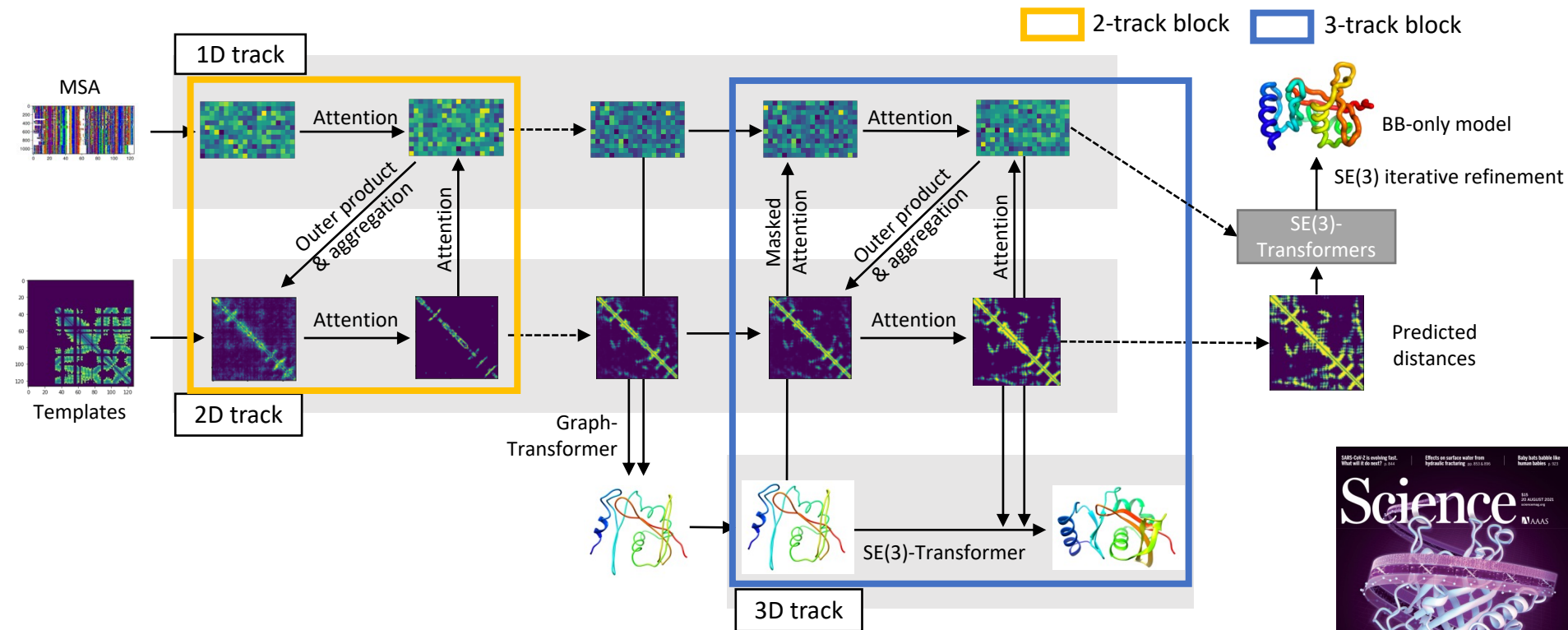
Progress in protein structure prediction



AlphaFold2: A game changer



RoseTTAFold: Can academia make something like AlphaFold?

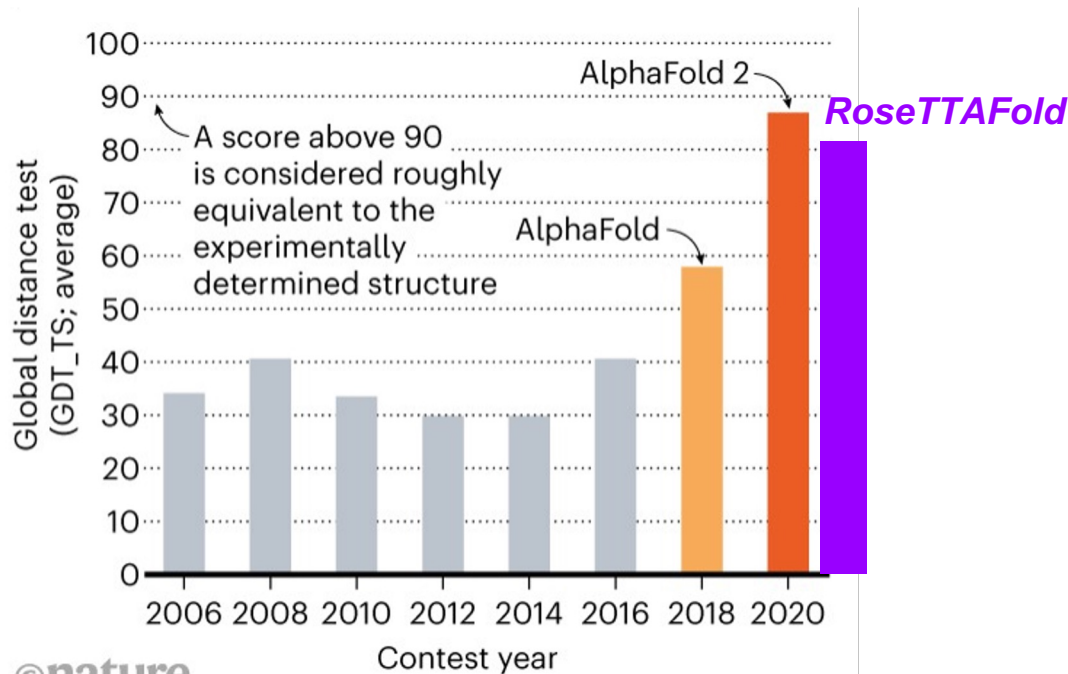


Baek, M., et al, *Science* (2021)

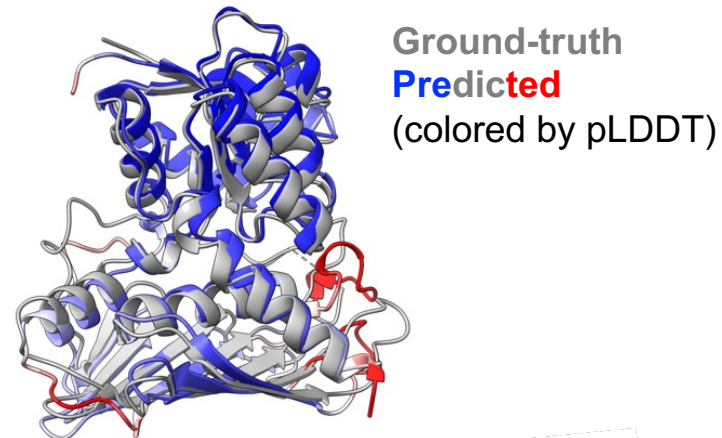


RoseTTAFold: Academia can do!

Free modeling accuracy in CASP14

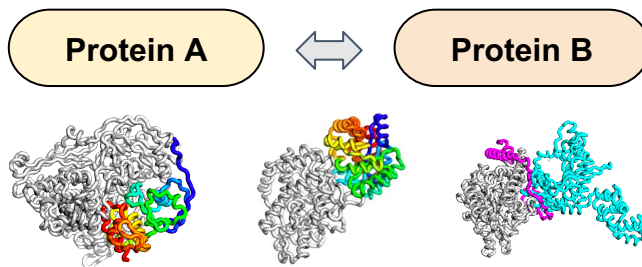


©nature



Beyond accurate modeling of protein structures

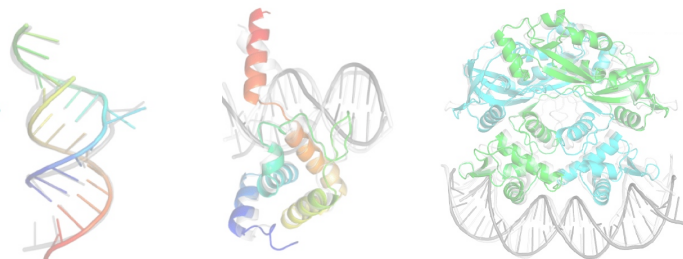
Large-scale *in silico* PPI screening



- 1) Do **A** and **B** interact?
- 2) What is the structure of **AB**?

Humphreys, I., Pei, J., Baek, M., Krishnakumar, A., et al, *Science* (2021)

Nucleic acid structure and interaction modeling



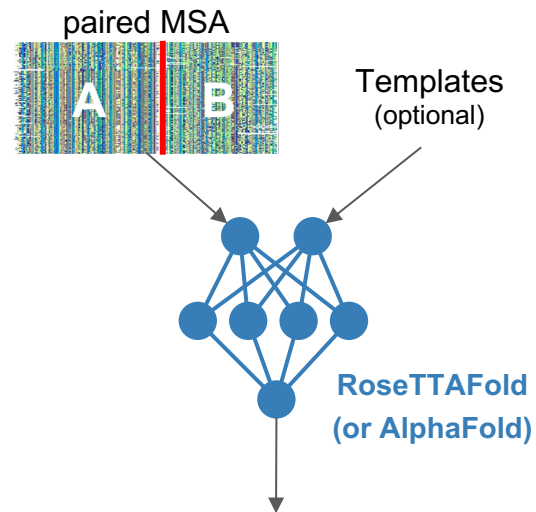
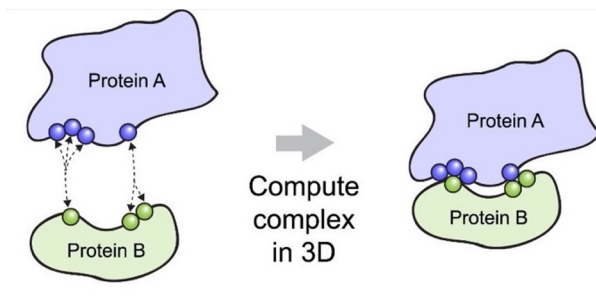
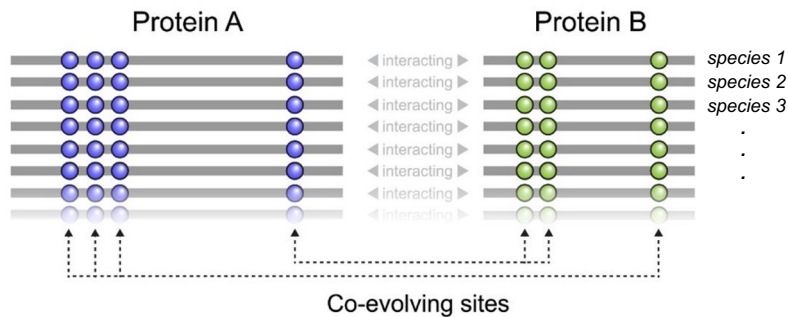
Baek, M., et al, *bioRxiv* (2022)

De novo functional protein design



Wang, J., et al., *Science* (2022)

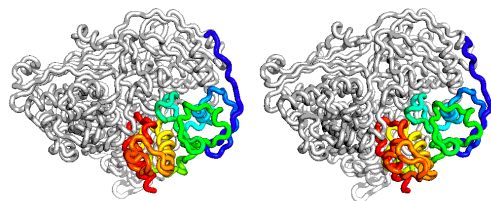
Protein-protein complex structure prediction



- 1) Do **A** and **B** interact?
- 2) What is the structure of **AB**?

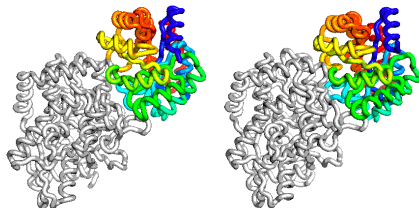
Protein-protein complex structure prediction

Aldehyde oxidoreductase



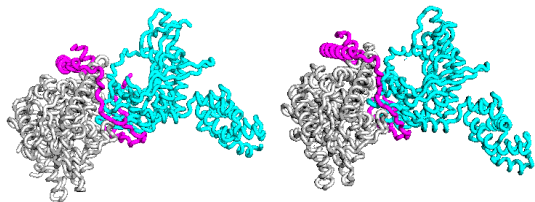
TM-score: 95

Tryptophan synthase



TM-score: 92

tRNA-dependent
amidotransferase



TM-score: 89



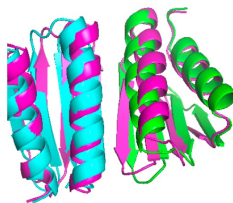
Minkyung Baek @minkbaek · Jul 20, 2021

Adding a big enough number for "residue_index" feature is enough to model hetero-complex using AlphaFold (green&cyan: crystal structure / magenta: predicted model w/ residue_index modification).
#AlphaFold #alphafold2

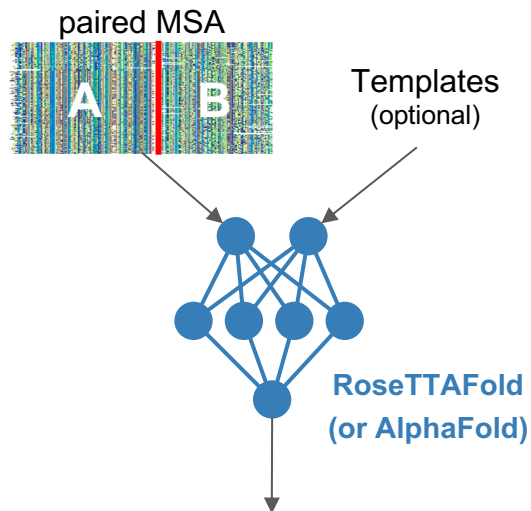
```
to residue index  
residue_index']  
  
in each chain
```

+= 200

```
dex'] = idx_res
```

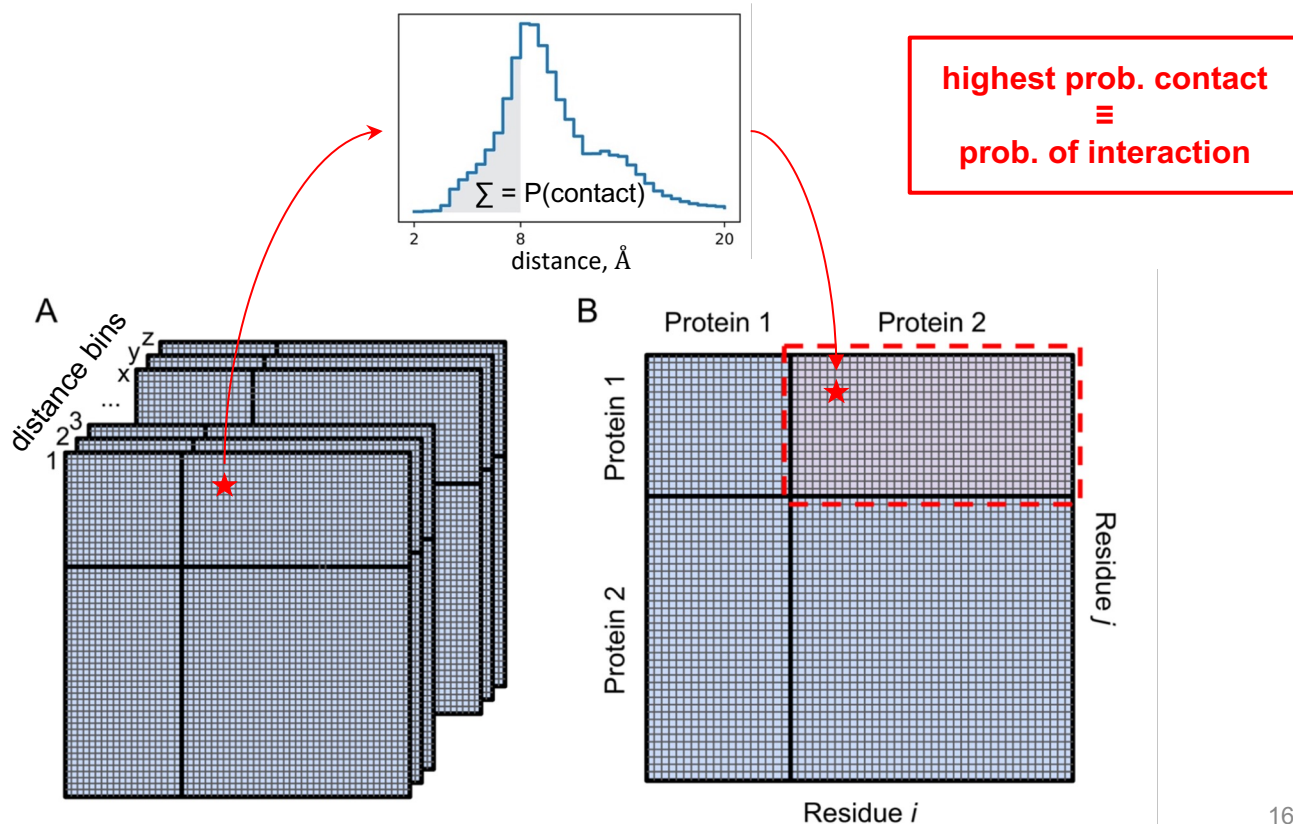
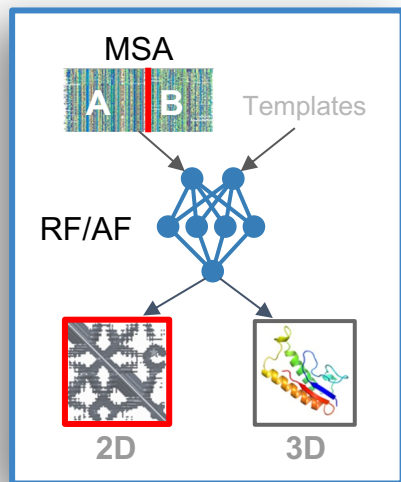


7 85 376



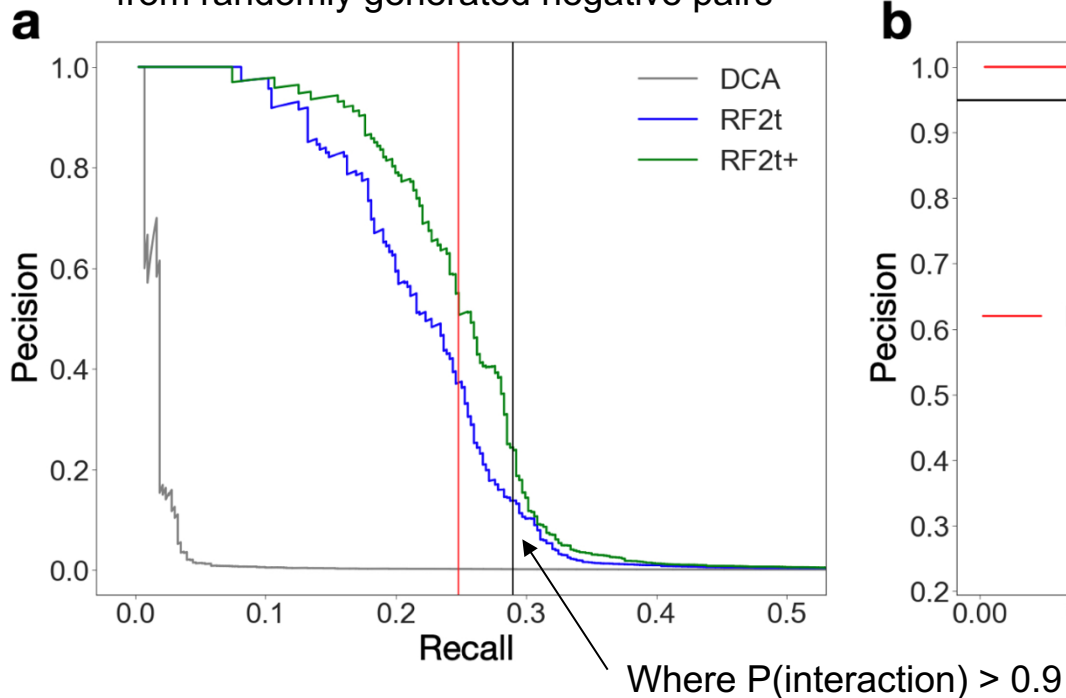
- 1) Do **A** and **B** interact?
- 2) What is the structure of **AB**?

Predicting protein-protein interactions

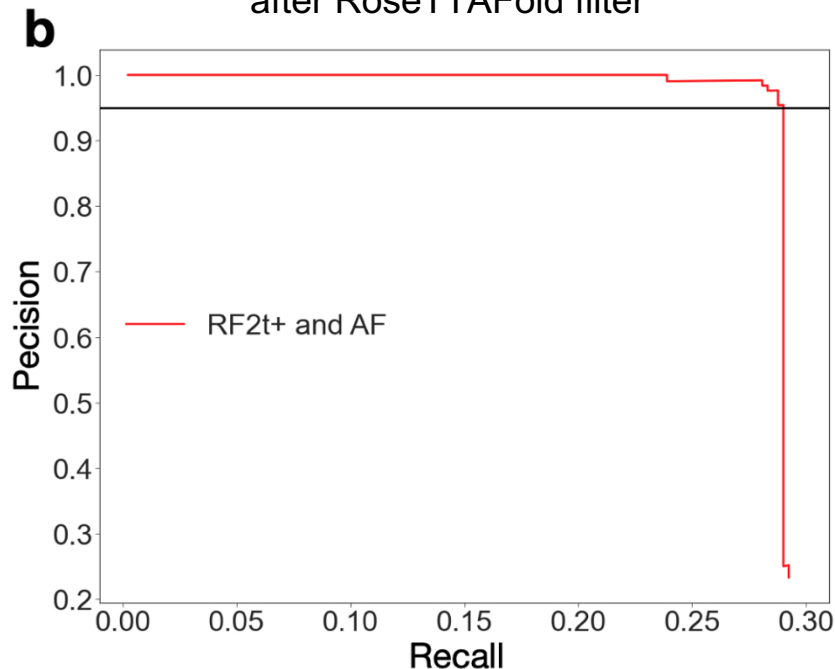


Predicting protein-protein interactions

Distinguish 768 gold-standard pairs from randomly generated negative pairs

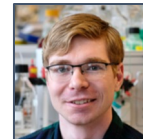


Distinguish 717 gold-standard pairs after RoseTTAFold filter



In silico PPI screening: Yeast interactome

- 4.3 million protein pairs



Ian



Aditya

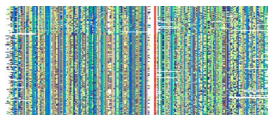


Qian



Jimin

pMSA



2-track RoseTTAFold
(10.7M parameters)
~ 10 sec per 1000aa complex

5,495
pairs

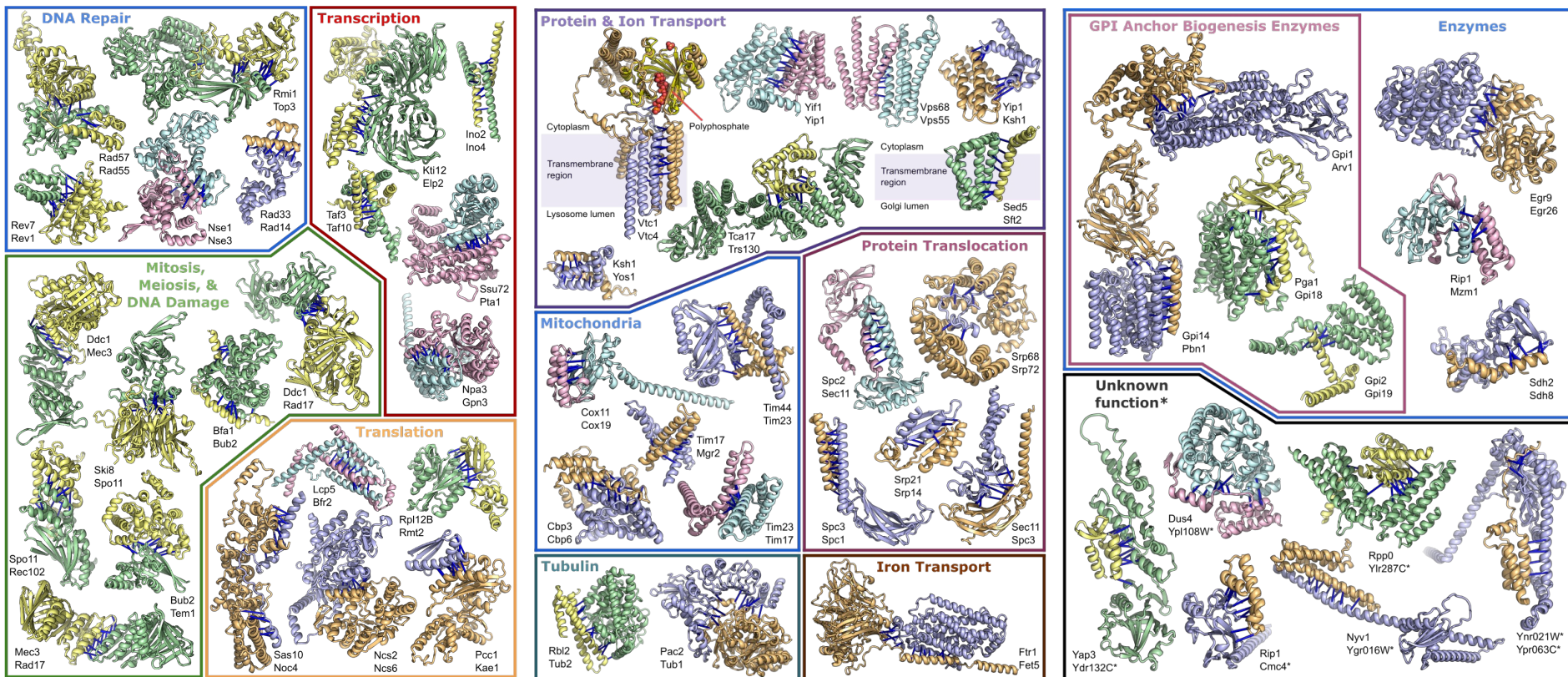


AlphaFold
~1800 sec per 1000aa complex

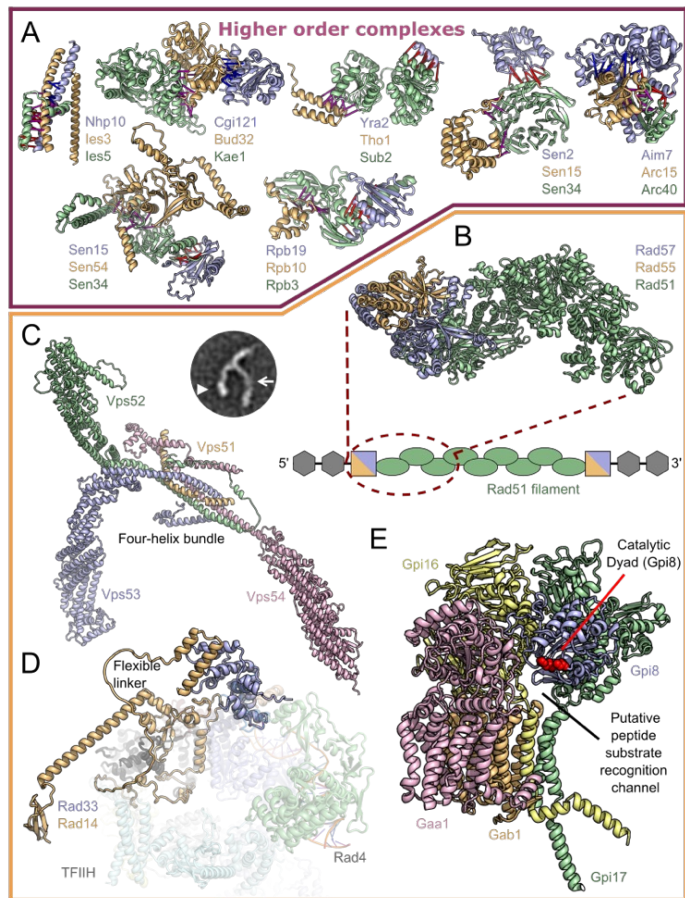


1,505 pairs
likely to interact
+
Manual study

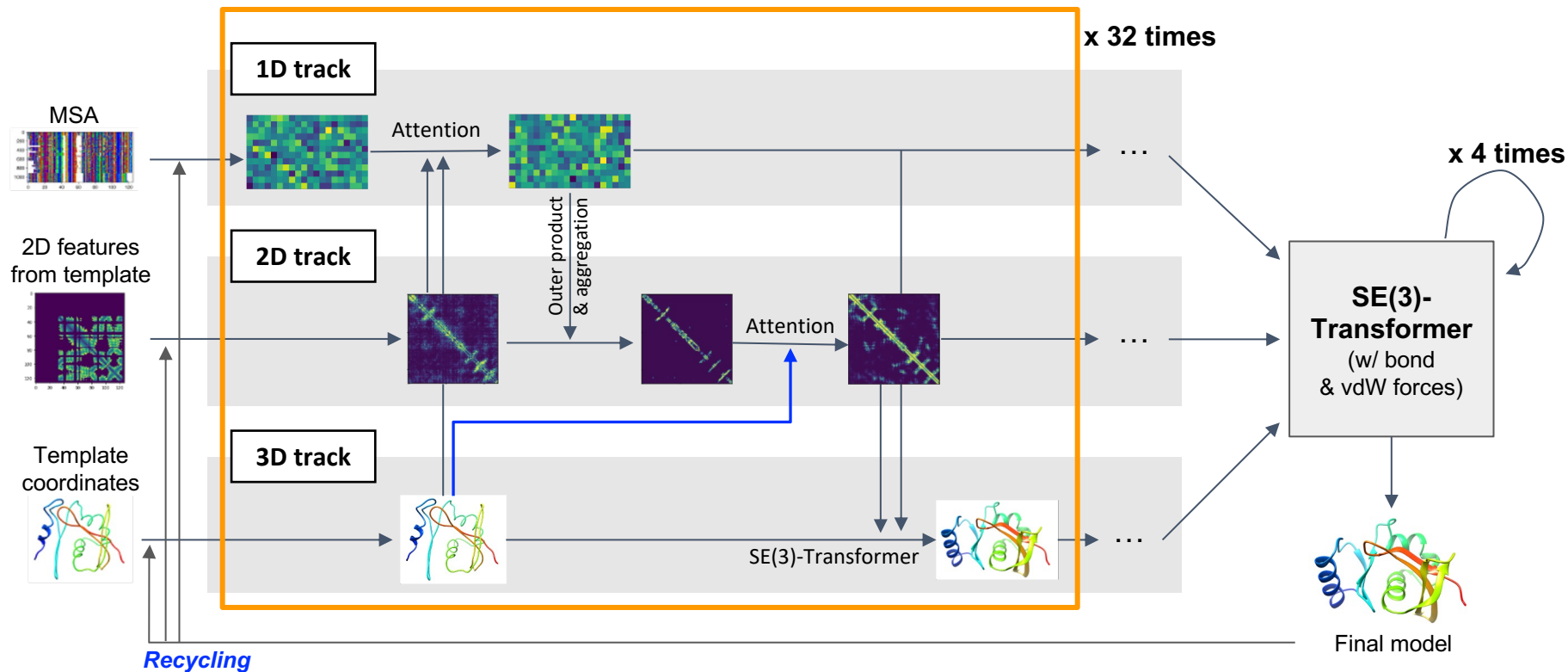
In silico PPI screening: Yeast interactome



In silico PPI screening: Yeast interactome

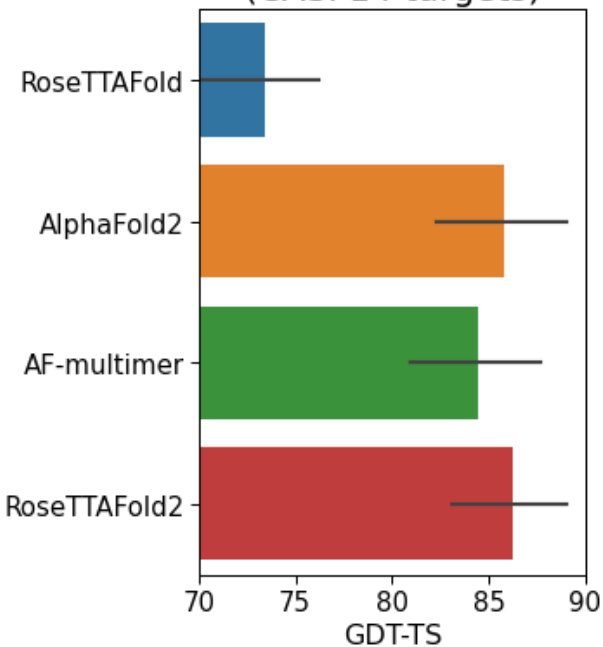


RoseTTAFold2: improving RF for better modeling & screening

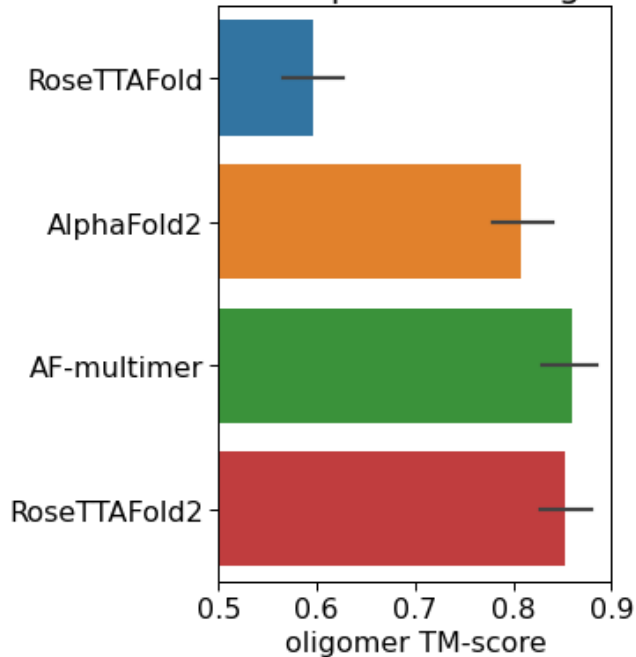


RoseTTAFold2: improving RF for better modeling & screening

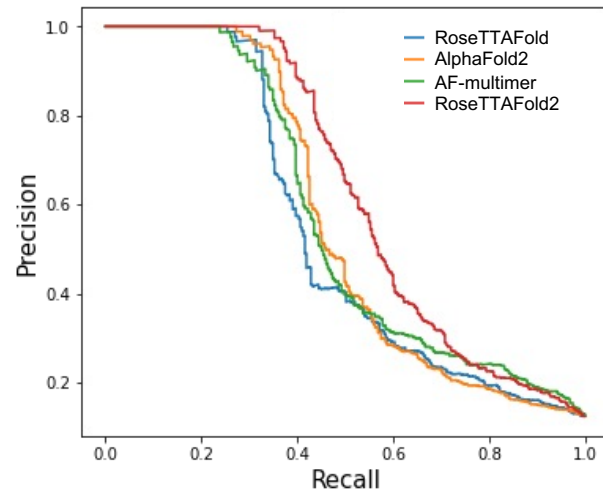
Tertiary structure modeling
(CASP14 targets)



Complex modeling

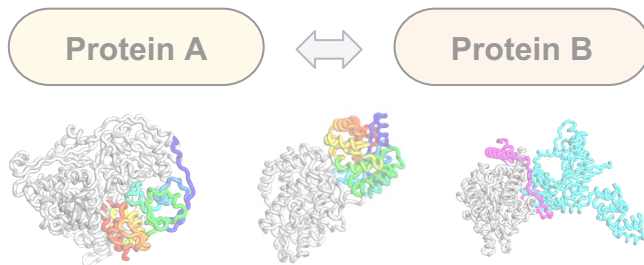


PPI screening benchmark



Beyond accurate modeling of protein structures

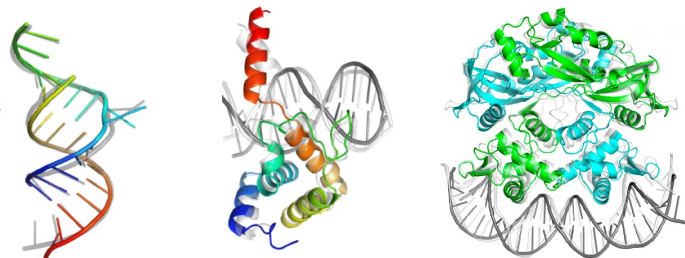
Large-scale *in silico* PPI screening



- 1) Do **A** and **B** interact?
- 2) What is the structure of **AB**?

Humphreys, I., Pei, J., Baek, M., Krishnakumar, A., et al, *Science* (2021)

Nucleic acid structure and interaction modeling



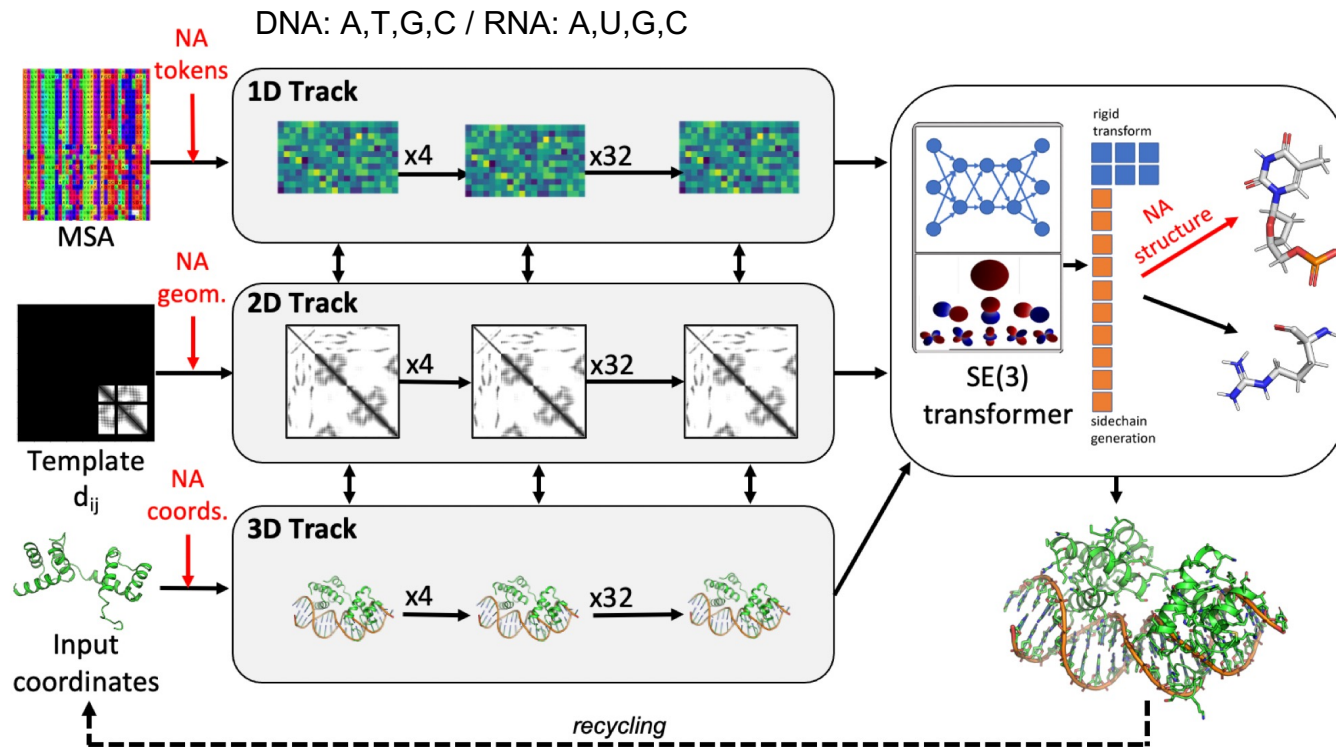
Baek, M., et al, *bioRxiv* (2022)

De novo functional protein design



Wang, J., et al., *Science* (2022)

Extending RoseTTAFold – nucleic acid prediction

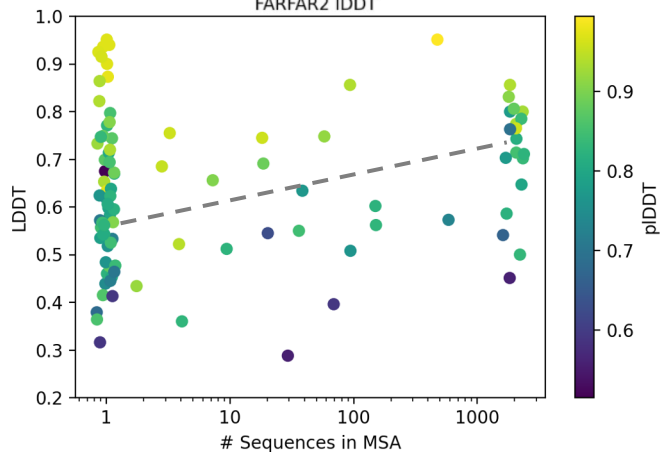
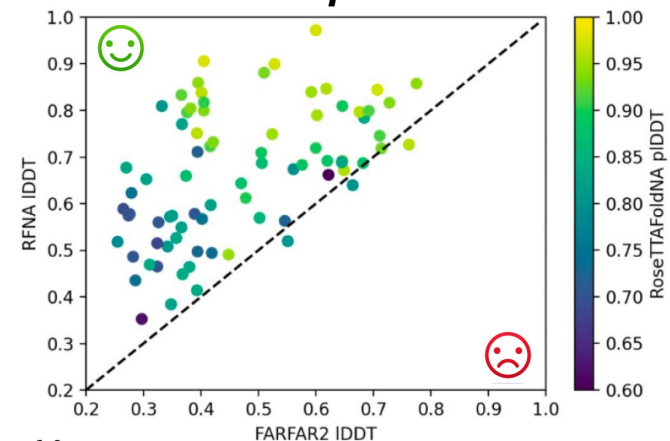


Frank DiMaio

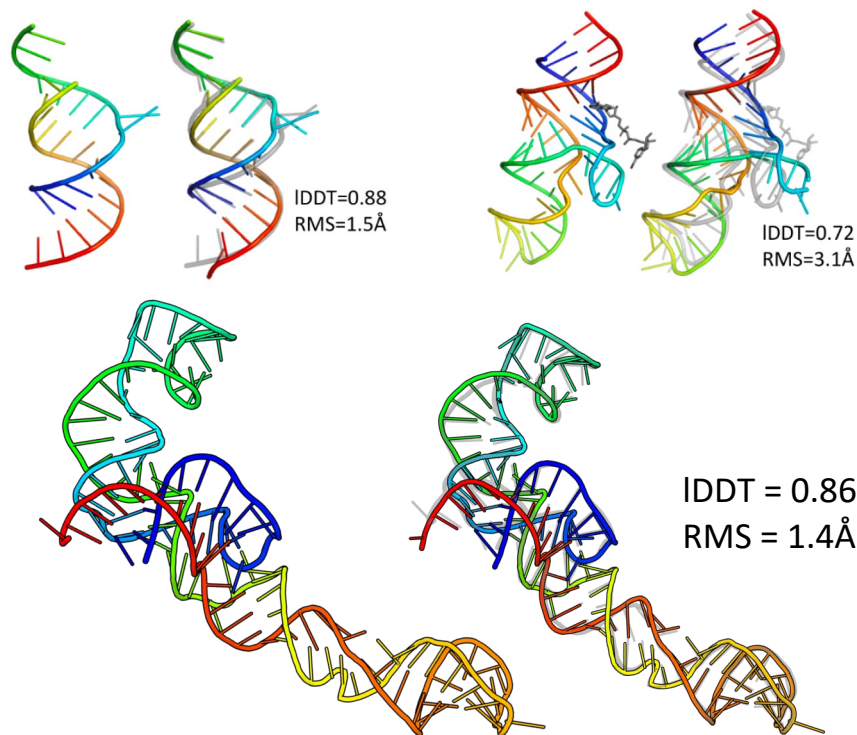
Trained on protein only set (>26k clusters) + **nucleic acid included structures (~3k clusters)**

RoseTTAFoldNA shows promising results

RNA structure prediction

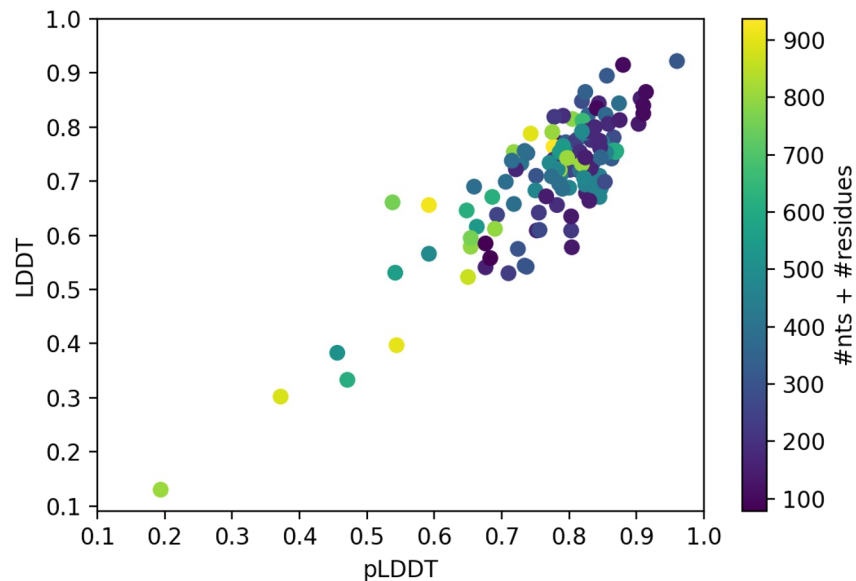


Ground-truth (left) vs Predicted (right)

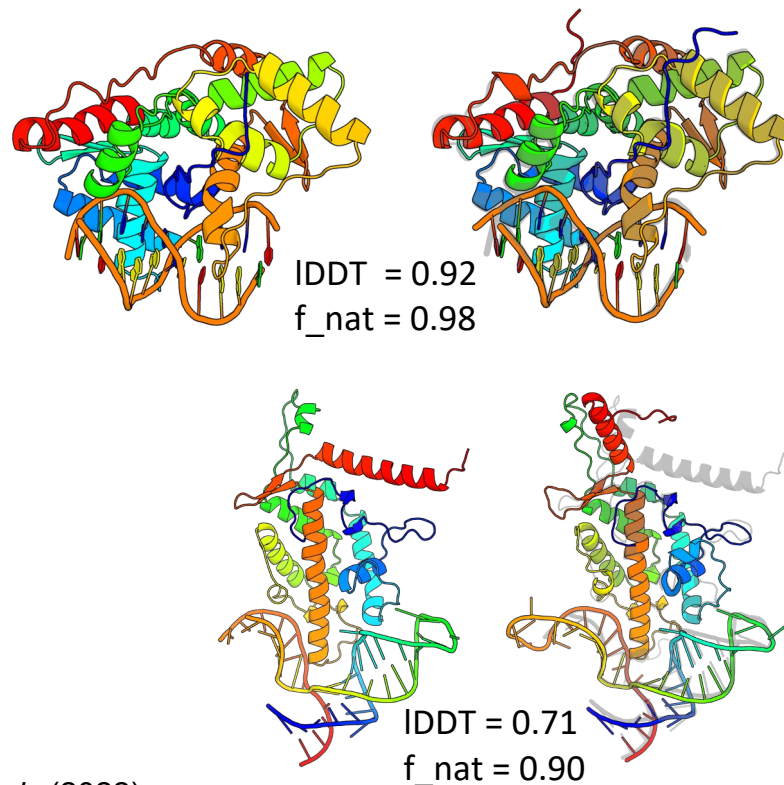


RoseTTAFoldNA shows promising results

Protein-nucleic acid complex prediction

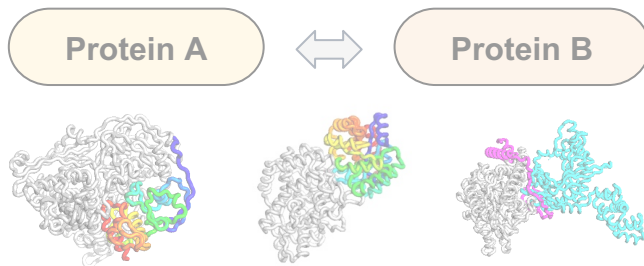


Ground-truth (left) vs Predicted (right)



Beyond accurate modeling of protein structures

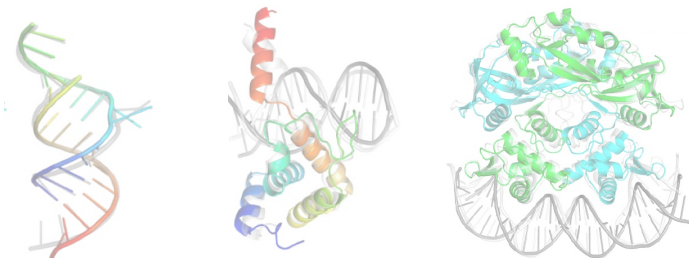
Large-scale *in silico* PPI screening



- 1) Do **A** and **B** interact?
- 2) What is the structure of **AB**?

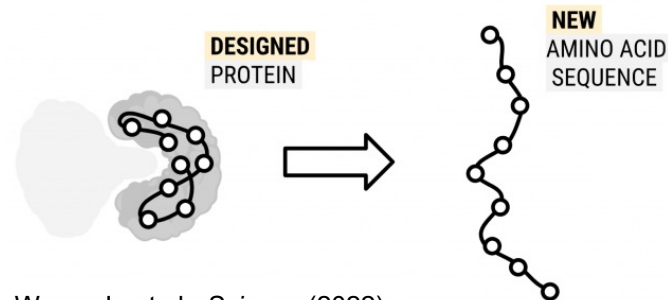
Humphreys, I., Pei, J., Baek, M., Krishnakumar, A., et al, *Science* (2021)

Nucleic acid structure and interaction modeling



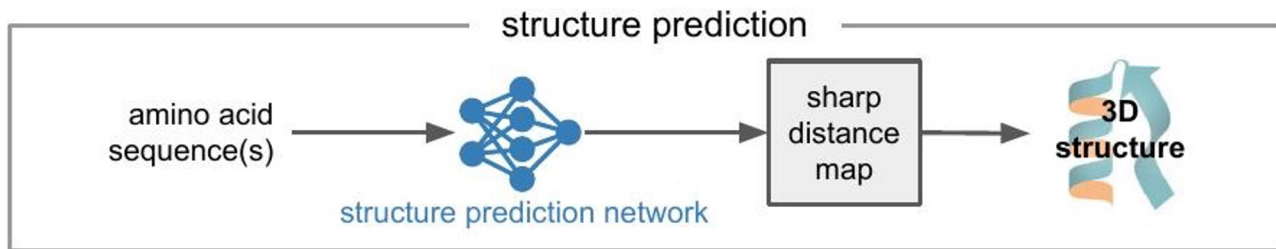
Baek, M., et al, *bioRxiv* (2022)

De novo functional protein design

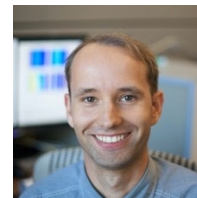


Wang, J., et al., *Science* (2022)

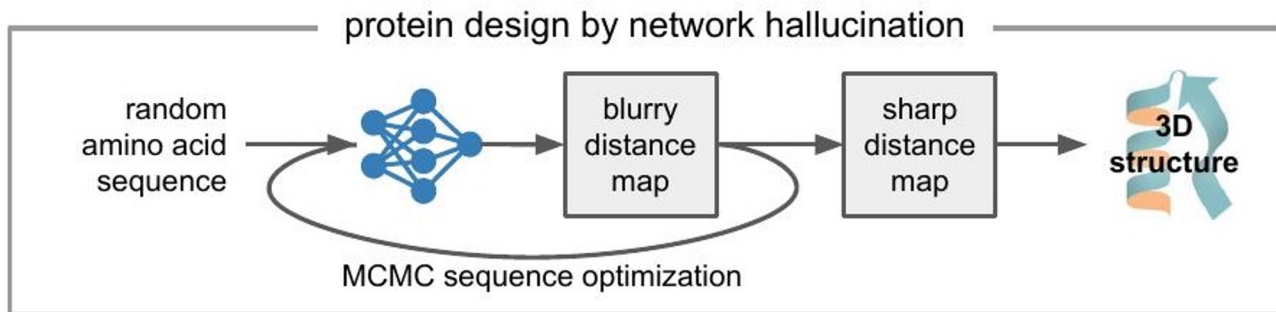
Hallucination: Optimize sequence to have a structure



method development



Ivan



experimental validation



Sam

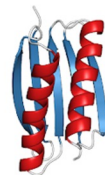
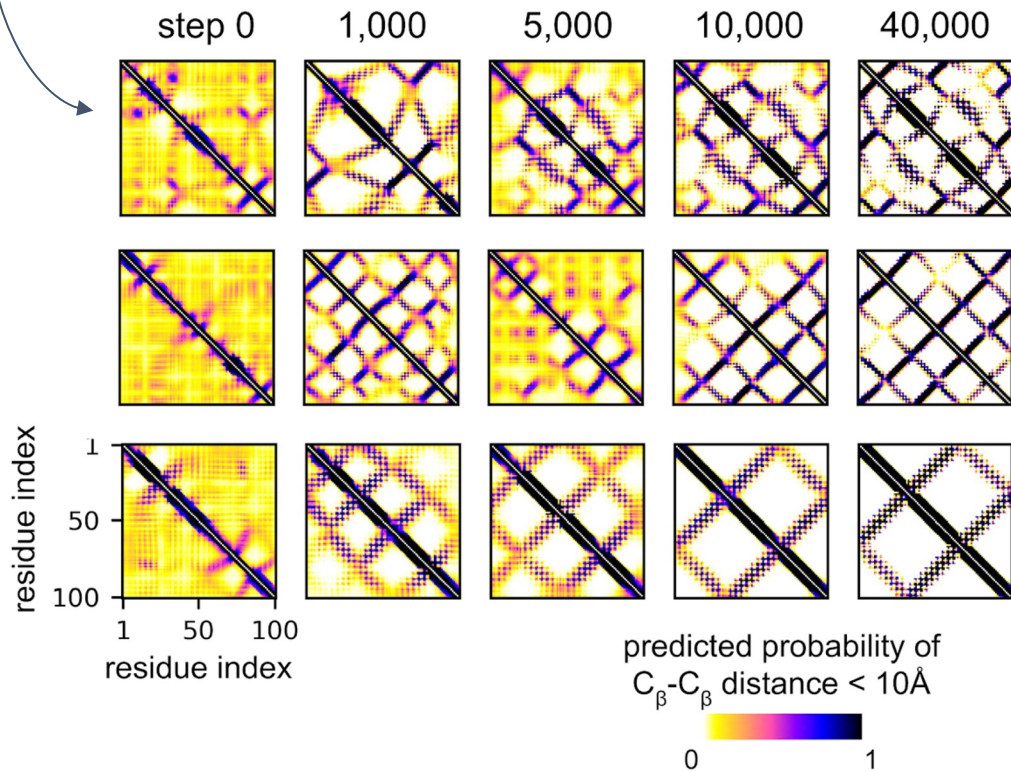


Tamuka

random
sequence

emergence of structure

optimized
sequence



mixed α and β



all- β



all- α

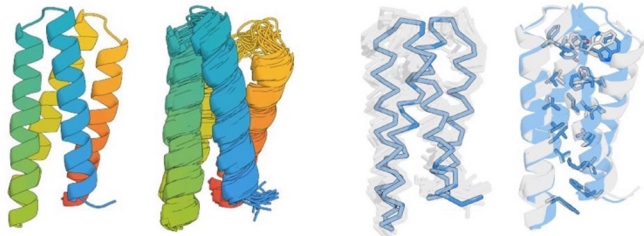
0515

Hallucination

NMR

Hallucination/NMR

1.82 Å bb RMSD over 100 aa



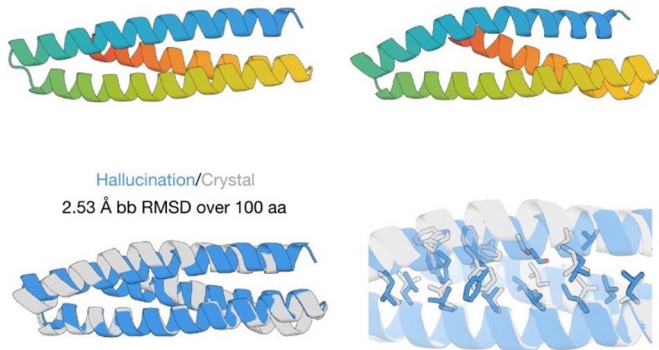
0217

Hallucination

Crystal

Hallucination/Crystal

2.53 Å bb RMSD over 100 aa



Structures of 3 hallucinations were confirmed experimentally

0738

Hallucination

Crystal

Hallucination/Crystal

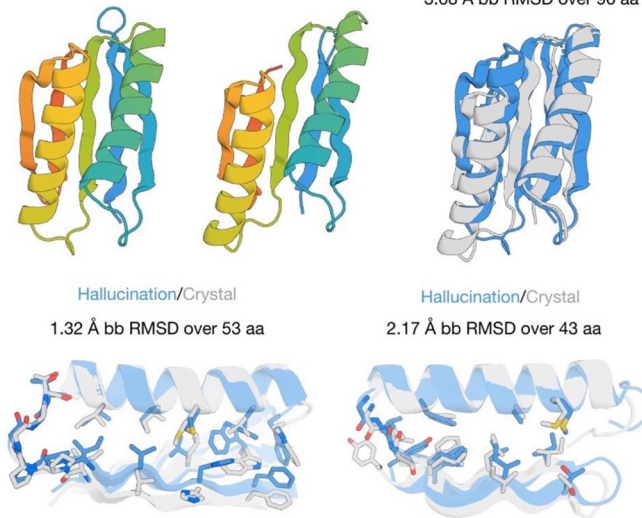
3.68 Å bb RMSD over 96 aa

Hallucination/Crystal

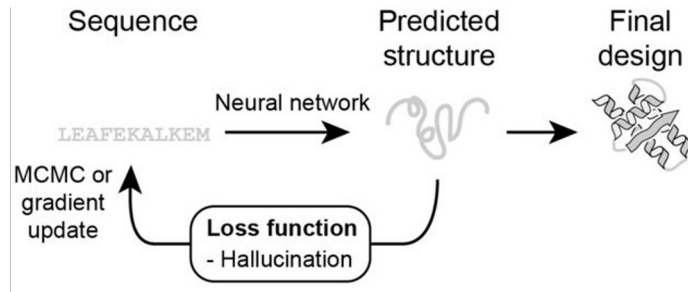
1.32 Å bb RMSD over 53 aa

Hallucination/Crystal

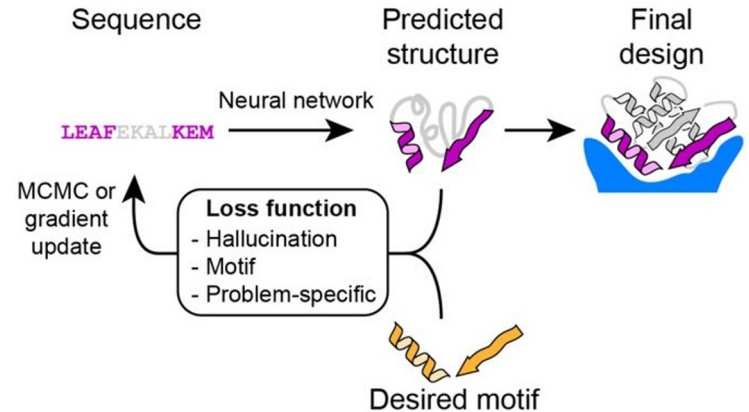
2.17 Å bb RMSD over 43 aa



Constrained hallucination: design a new protein having a given functional motif



Free hallucination:
generate novel protein folds



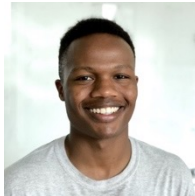
Constrained hallucination:
generate scaffolds harboring
pre-specified functional sites



Jue



Doug

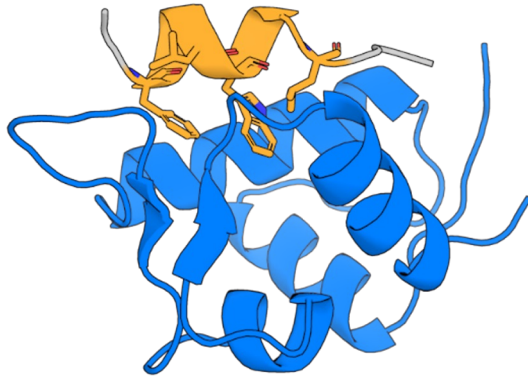


Sidney

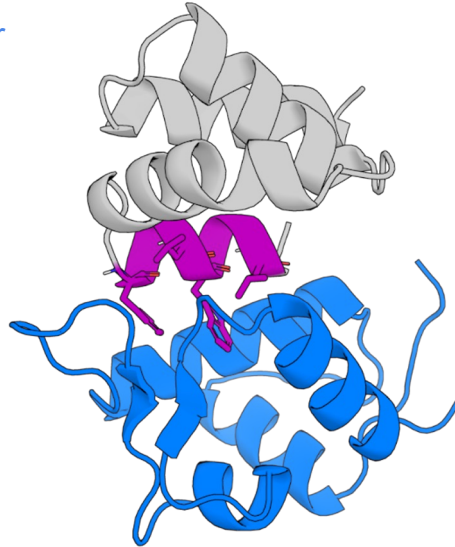
Applications of constrained hallucination

Scaffolding p53 helix to bind cancer-signaling protein mdm2

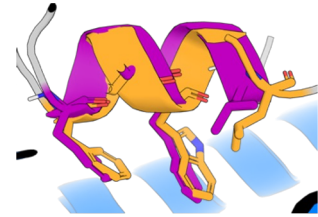
Native motif
Design motif
Binding partner



Native



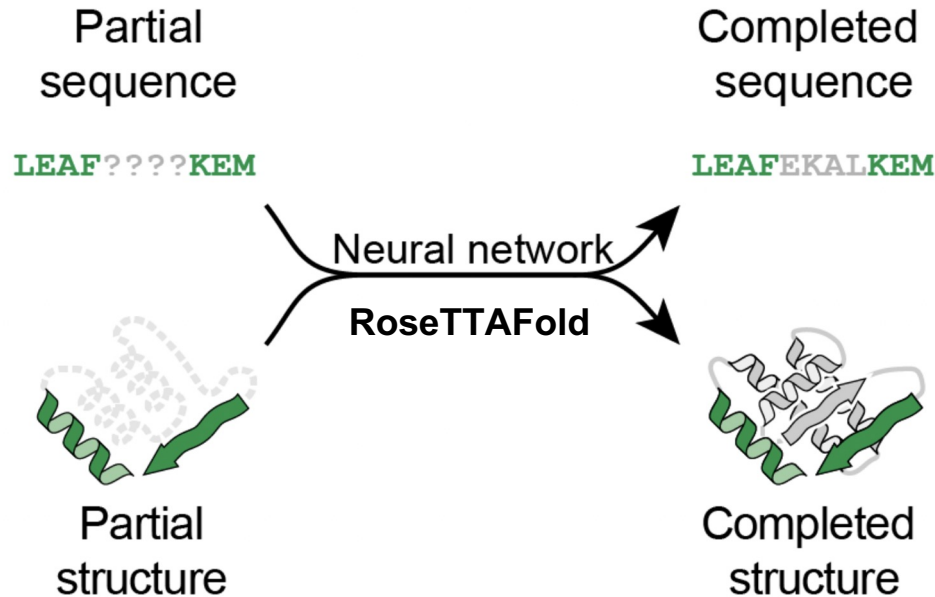
Hallucination



Native vs designed
motif

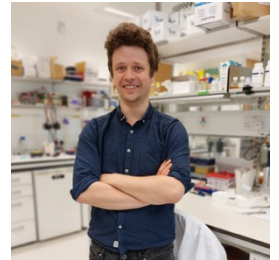
Protein Design via Inpainting

Formulate motif-based protein design as information completion (or inpainting)



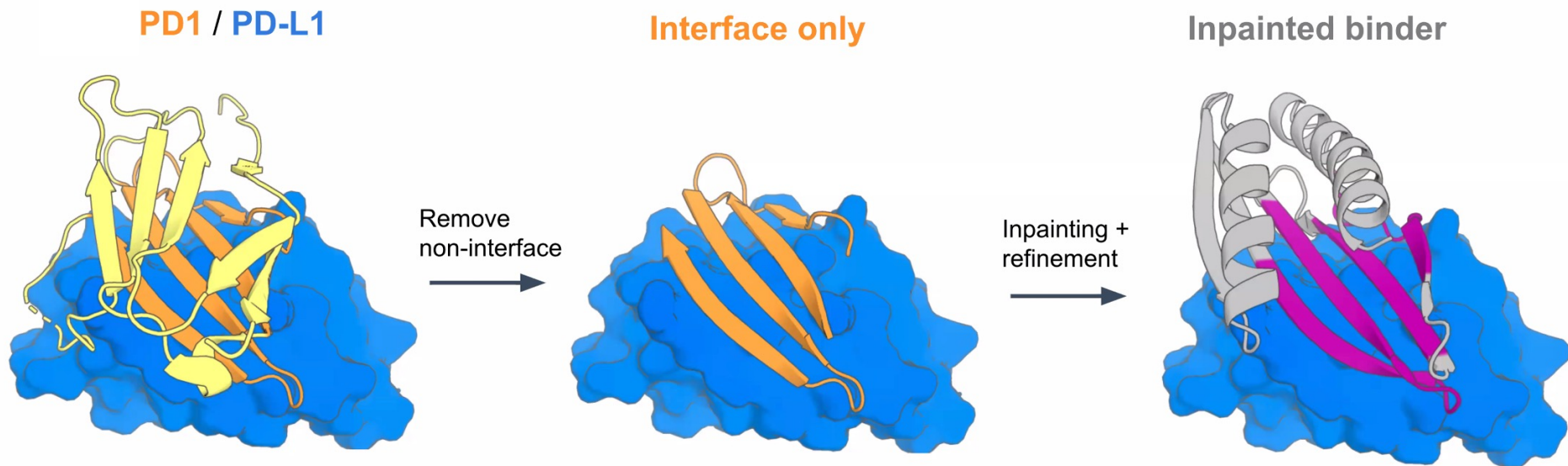
David Juergens

Wang, J., Lisanza, S., Juergens, D., Tischer, D., Watson, J., Science (2022)



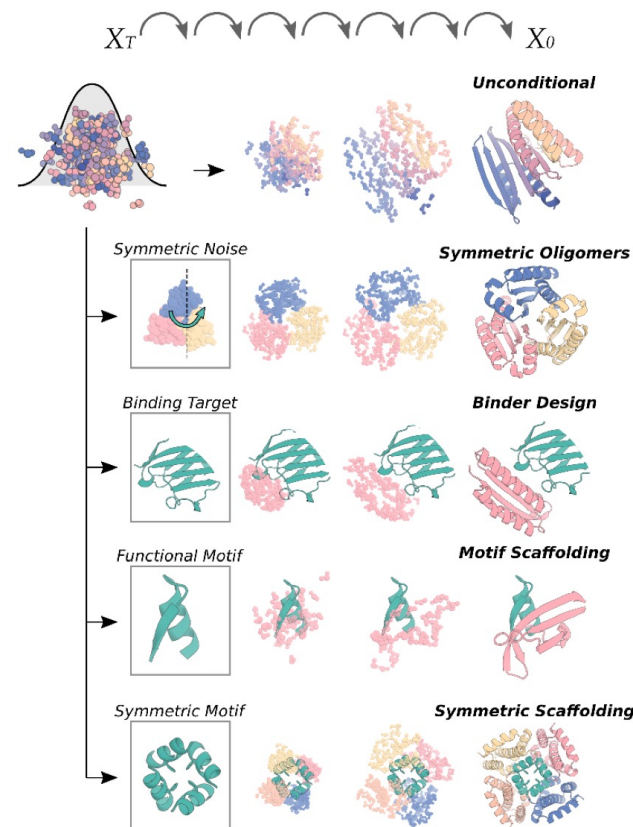
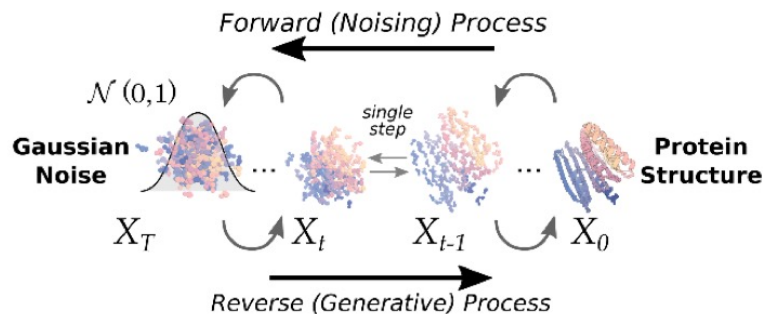
Joseph Watson

Design PD-L1/PD-1 binding inhibitor via inpainting

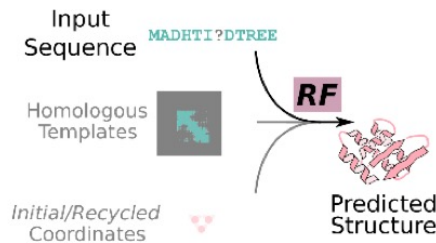


Generative model for protein design (RFdiffusion)

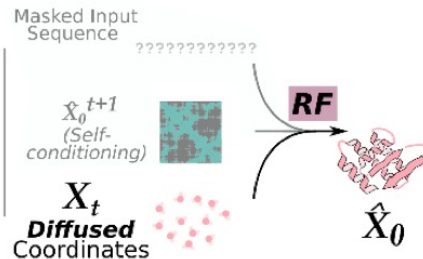
Diffusion Model



RoseTTAFold

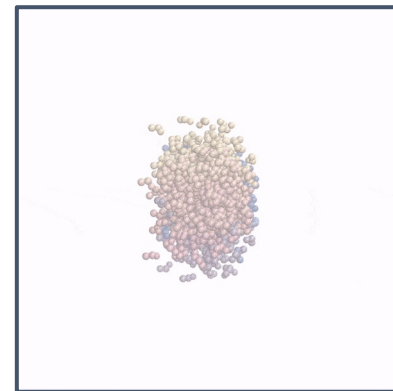
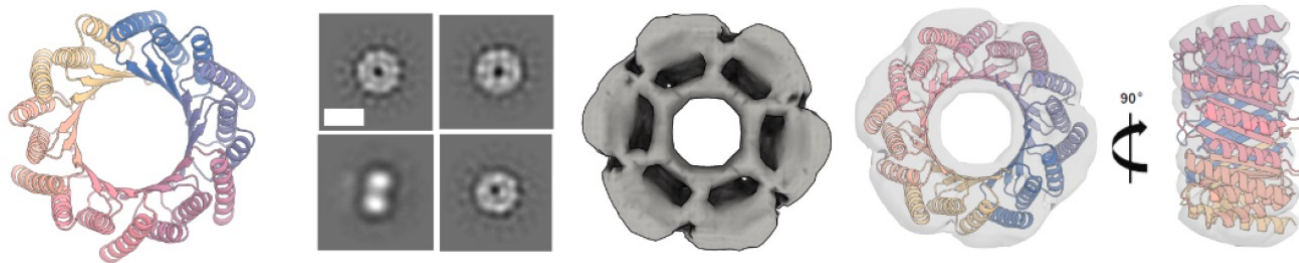


RFdiffusion

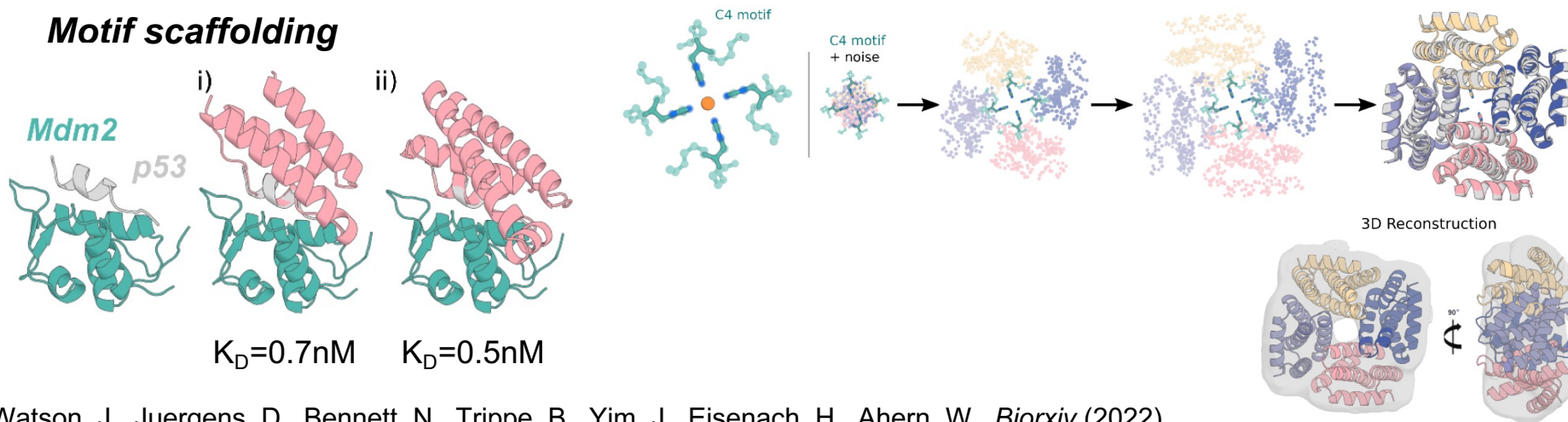


RFdiffusion can do various design tasks

High-order symmetric assemblies



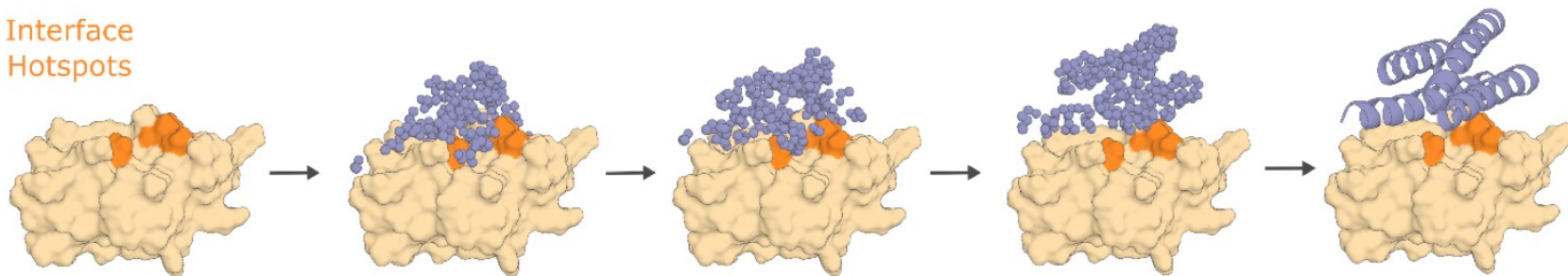
Motif scaffolding



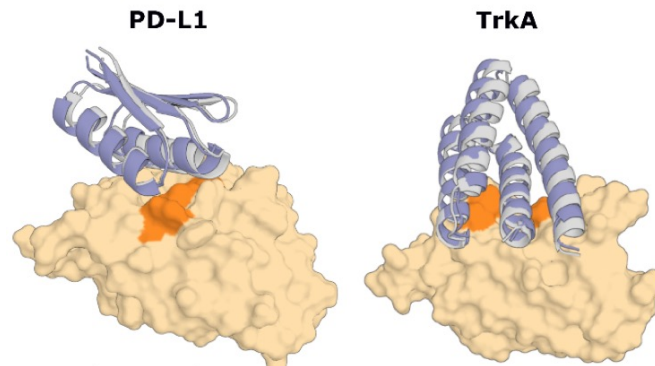
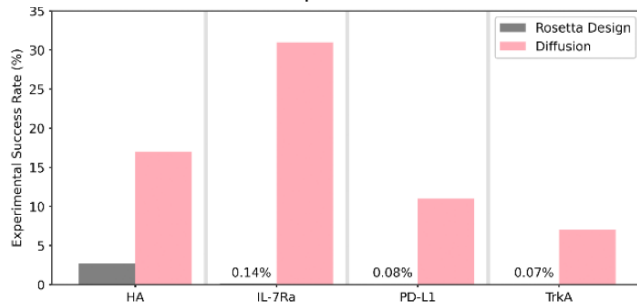
RFdiffusion can do various design tasks

Binder design based on a given interface

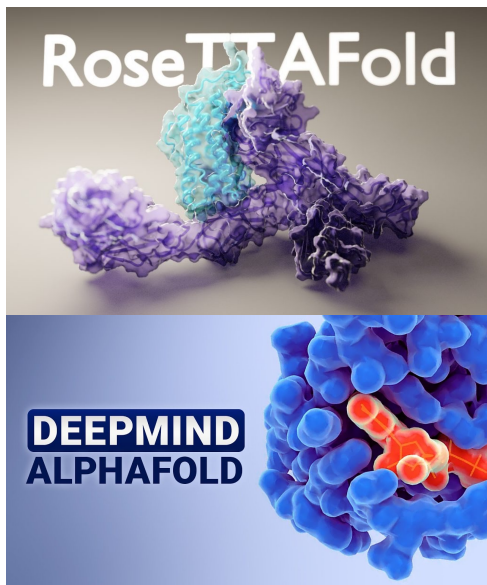
Interface
Hotspots



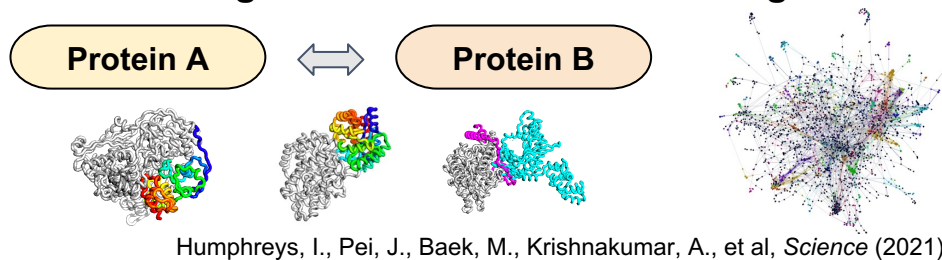
RFdiffusion has orders-of-magnitude
higher **experimental** success
rates than previous methods



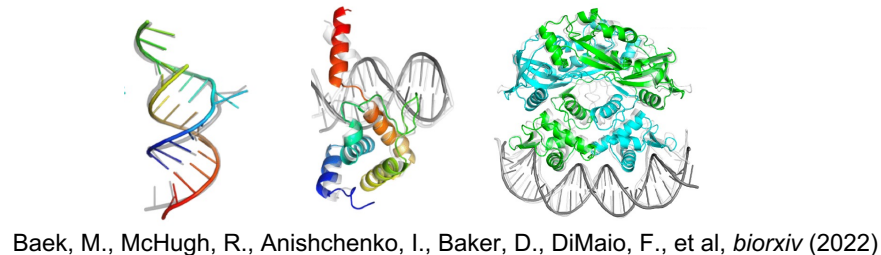
AI-based protein modeling



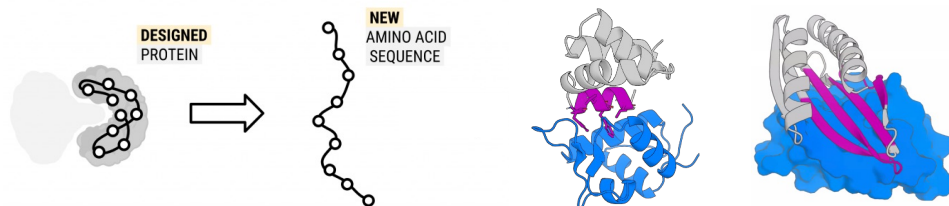
Large-scale *in silico* PPI screening



Nucleic acid structure & interaction prediction



De novo functional protein design



Acknowledgements

RoseTTAFold

Frank DiMaio

Ivan Anishchenko
Justas Dauparas
Sergey Ovchinnikov
Jue Wang

PPI screening

Qian Cong

Ian Humphreys

Aditya Krishnakumar
Jimin Pei

Hallucination

David Juergens

Joseph Watson

Inpainting & Diffusion

Joseph L. Watson

David Juergens

Nathaniel R Bennett

Brian L Trippe

Jason Yim

Helen E. Eisenach

Woody Ahern

Diffusion

Ivan Anishchenko

Doug Tischler

Jue Wang

Sidney Lisanza

Sam Pellock

Tamuka Chidyausiku

Sergey Ovchinnikov

Chris Norn

IT support

Luki Goldschmidt

David Kim

Microsoft

AWS

Discussion

David Baker

Lance Stewart

Eric Horvitz

Looking for graduate students & postdoc!!

If you're interested in, please e-mail me (minkbaek@snu.ac.kr)

