

Question 2

- (a) Briefly describe a computational biology topic of your choice, its brief history and provide a motivation for choosing it. (1/2 page) (5 marks)
- (b) State some computational goal(s) or problem(s) most relevant to the topic. Aim to describe the computational goal(s) or formulate the problem(s) in a way understandable for computer scientists. (1/2 page) (10 marks)
- (c) Describe how molecular biology would benefit from achieving the computational goal(s) or solving the stated problem(s). (1/2 page) (5 marks)
- (d) Characterise some existing solutions (to date) to this problem available in the literature and clearly cite references.
- (i) Describe the underlying ideas or algorithms associated with these solutions by creating your own schemes, pseudocodes, equations, illustrations of model architectures (as applicable) to demonstrate your understanding of the technical background related to these solutions. (2 pages) (20 marks)
 - (ii) Quantify and compare advantages and limitations of existing solutions using figures and tables (as applicable). If available, you may use data (e.g. on quantifying the performance of a given method) in the literature but you have to create your own informative (and carefully labeled) figures and tables. (1 page) (10 marks)
- (e) Provide critical assessment of all solutions you consider and discuss their potential (future) impact. (1 page) (10 marks)
- (f) Based on your understanding of the state-of-art in the field, propose two specific suggestions for future development to complement, replace or improve on existing solutions. Provide clear arguments to support your suggestions and describe what expected advantages and limitations your proposed developments would bring when compared to existing methods available in the literature. (1 page) (10 marks)

Structure prediction of protein-nucleic acid complexes via deep learning

Weiqi Feng

(a) Topic

Our topic is to predict the structure of protein-nucleic acid complexes based on their sequences, through the means of deep learning.

This problem originates from the “protein folding problem” in Biology literature, which has been an important research topic for more than 50 years.¹ Recently, with the development in deep learning, new computational approaches, such as AlphaFold¹ and RoseTTAFold³, have been proposed to quickly predict the 3D structure of a protein using only its amino acid sequence. Achieving great success, researchers soon moved on to other biomolecules like nucleic acid, and more recently, their complexes. Given separate structures, some people tried to build up complexes using docking calculations, but AlphaFold3² and RoseTTAFoldNA⁵ demonstrate a more general and promising methodology.

Biomolecules are essential to life, and understanding their structure can facilitate a mechanistic comprehension of their function.¹ As a result, I chose this topic because I believe this is the crucial step to truly understand and thus make use of these fundamental molecules, which can benefit many fields such as biology and medicine, and ultimately, human kind. Moreover, as a computer science student, I am intrigued by how these networks incorporate iconic deep learning models, such as the ideas of transformer and RNN recycling, and combine them with knowledge in biology, chemistry and physics.

(b) Problem definition

The problem of choice is quite straightforward: given the sequence of a certain protein-nucleic acid complex, we want to use a network to predict its 3D structure quickly and accurately. In terms of computer scientists, the input is the sequence and chemical structures of a complex, such as ligands and covalent bonds. From this, we extract token-level and pair-level multi-dimensional features, pass them through a trained network to construct 3D structure from its output, and maybe calculate the confidence estimate.

(c) Benefit

Biological complexes are essential to life, and understanding their 3D structures is the key to understanding cellular functions and the rational design of therapeutics.³ Accurate computational approaches can largely accelerate the speed in determining their structure compared to experimental methods, and thus enable large-scale structural bioinformatics,¹ which greatly benefits molecular biology and relevant fields such as medicine and chemistry.

(d) Solutions

This essay focuses on two mainstream approaches, the AlphaFold family by DeepMind and the RoseTTAFold family.

(i) Algorithms

Rough architectures and pseudo codes of AlphaFold^{3,2} and RoseTTAFold All Atom⁶ are drawn on the following two pages. Black denotes model and code, blue denotes important illustrations, and purple denotes major modifications from their former models: AF2¹ and RF2⁴. Due to restriction on paragraph, I can not demonstrate all architecture details and pseudo codes, but here are some highlights worth noting.

① AlphaFold3^{1,2}

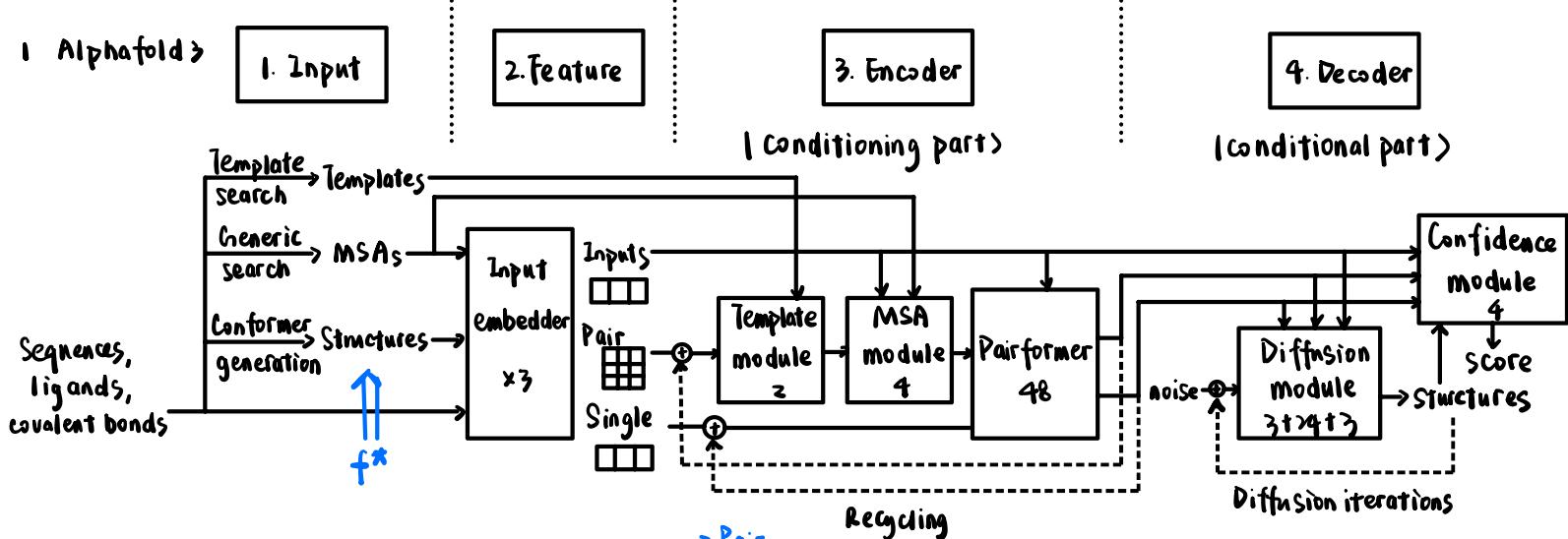
The key improvements of AF3 are that it is adapted to more general chemical structures and has higher learning efficiency. This is achieved through following modifications. Firstly, in the encoder, AF3 uses Pairformer instead of Evoformer, and isolates MSA module to reduce MSA-relevant dependency and calculations, more specifically, it replaces row-wise self-attention with pair-weighted average and single attention. Triangle update and self-attention remain the same. Secondly, in the decoder, it uses Diffusion Module on atoms and tokens instead of Structure Module, and ignores invariant point constraint in AF2.

② RoseTTAFoldNA^{3,4,5,6}

Because RF2⁴, RFNA⁵ and RFAA⁶ have very similar architecture and pipeline, this essay refers to all three papers. RFNA is the first network to predict the structures of protein – nucleic acid complexes without homology, a problem unsolved by AF and RF.⁵ Its key features are that it uses a 3D track architecture of iterative updates while sharing information among them, and uses SE(3)-Transformer with equivariance in rotation and translation.

Question 2 1d) (i)

1 AlphaFold3



def MainInferenceLoop($\{f^*\}$, Ncycle=4, Cs=384, Gs=128):

↳ Token: Inputs & Pair

$S_i^{inputs} = \text{InputFeatureEmbedder}(f^*)$ $S_i^{inputs} \in \mathbb{R}^{Cs}$ Inputs for every token

$S_i^{init} = \text{LinearNoBias}(S_i^{inputs})$

$S_i^{init} \in \mathbb{R}^{Cs}$ Single for every token

$z_{ij}^{init} = \text{LinearNoBias}(S_i^{inputs}) + \text{LinearNoBias}(S_j^{inputs})$

$z_{ij}^{init} = \text{RelativePositionEncoding}(\{f^*\})$ position and order information

$z_{ij}^{init} = \text{LinearNoBias} | f_{ij} \text{ token-bonds } \rangle$ bond info $z_{ij}^{init} \in \mathbb{R}^{G_s}$ Pair for every pair of tokens

$\{\hat{z}_{ij}\}, \{\hat{S}_i\} = \vec{0}, \vec{0}$

↳ Pair ↳ Single

for all $C \in \{1, \dots, Ncycle\}$ do Recycling 4 times

$z_{ij} = z_{ij}^{init} + \text{LinearNoBias} | \text{LayerNorm}(\hat{z}_{ij}) \rangle$

$\{z_{ij}\} = \text{TemplateEmbedder}(\{f^*\}, \{z_{ij}\})$

only 4 MSA blocks in recycle:

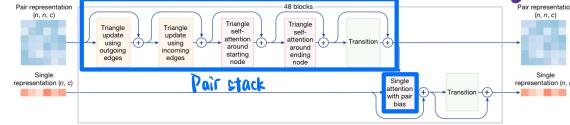
$\{z_{ij}\} = \text{MSAModule}(\{f_{S_i}^{MSA}\}, \{z_{ij}\}, \{S_i^{inputs}\})$

↳ MSA stack & Pair stack

$S_i = S_i^{init} + \text{LinearNoBias} | \text{LayerNorm}(\hat{S}_i) \rangle$

↳ row-wise self-attn → pair weighted averaging

$\{S_i\}, \{z_{ij}\} = \text{PairformerStack}(\{S_i\}, \{z_{ij}\})$



$\{\hat{S}_i\}, \{\hat{z}_{ij}\} \leftarrow \{S_i\}, \{z_{ij}\}$

~ Single row-wise

end for

$\{\hat{x}^{pred}\} = \text{SampleDiffusion}(\{f^*\}, \{S_i^{inputs}\}, \{S_i\}, \{z_{ij}\})$

Atom & token attn

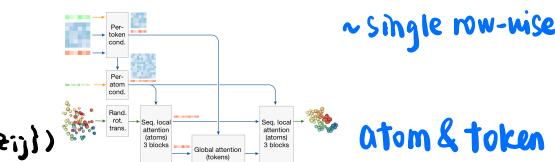
↳ Diffusion Transformer

4.

Decoder

$\{P_{ij}^{plddt}\}, \{P_{ij}^{pae}\}, \{P_{ij}^{rde}\}, \{P_{ij}^{resolved}\} = \text{ConfidenceHead}(\{S_i^{inputs}\}, \{S_i\}, \{z_{ij}\}, \{\hat{x}^{pred}\})$

$P_{ij}^{distogram} = \text{DistogramHead}(z_{ij})$



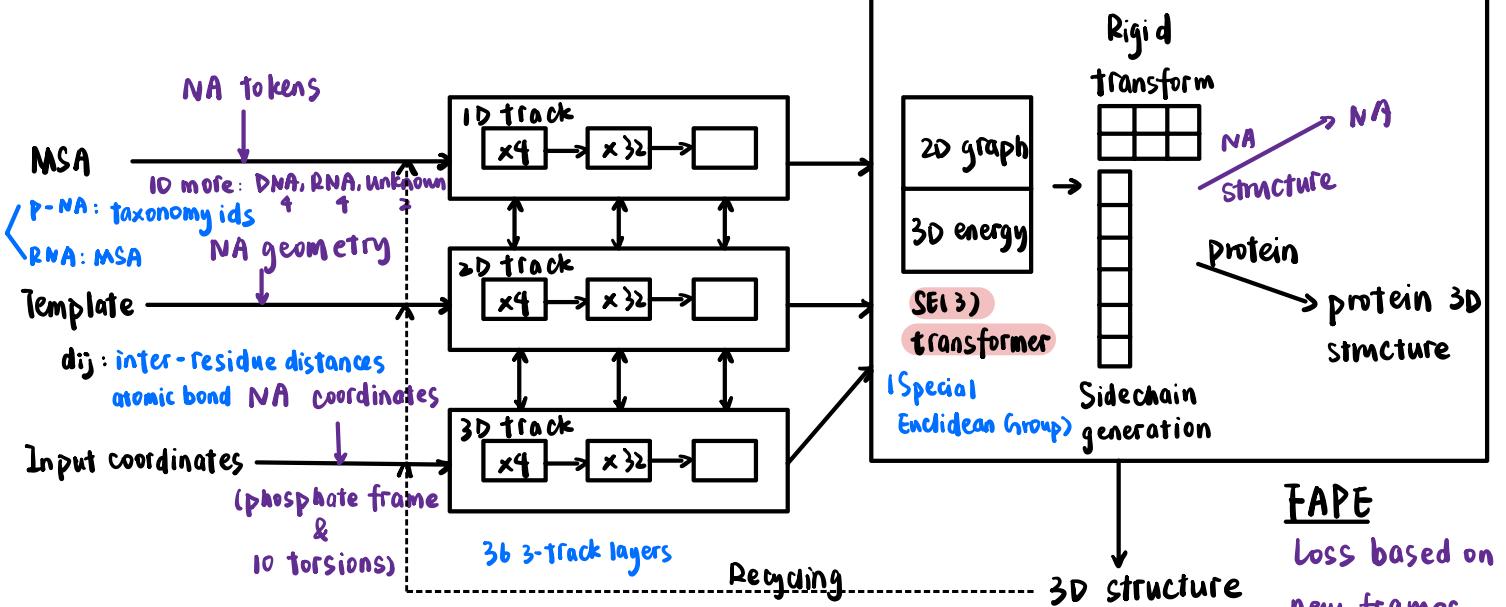
return $\{\hat{x}^{pred}\}, \{P_{ij}^{plddt}\}, \{P_{ij}^{pae}\}, \{P_{ij}^{rde}\}, \{P_{ij}^{resolved}\}, P_{ij}^{distogram}$

2 RosettaFold NA scarcity of data, physical terms in loss

Structure refinement module

x 4 layers

RosettaFold NA:



We refer to RF2 & RFAA for more details of the Simulator stage and recycling.

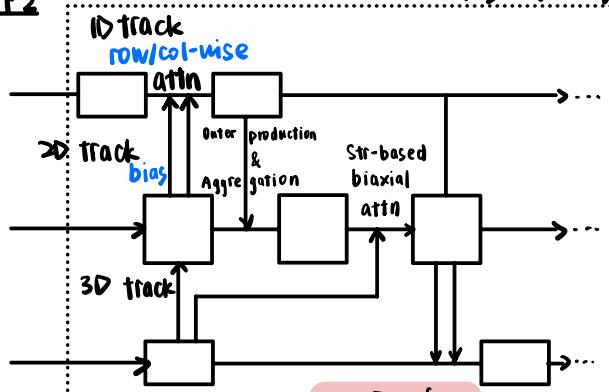
One cycle of the simulator stage, RFAA forward pass:

```

    ↗ special feature of RF   ↗ residue
def RFAA(msa-masked, msa-full, seq, idx, bond-feats, dist-matrix,
         recycled, iterations
         templates, msa-prev, pair-prev, state-prev, xyz-prev, sc-torsions-prev,
         ↳ recyclables in total J
         Nfull-blocks=4, Nmain-blocks=32, Nref-blocks=4):
    MSA_fmn = MSA.Full.Embed(msa-full)
    MSAcluster, pair, state = MSA.Cluster.Embed(msa-masked
        3D track           seq, idx, bond-feats, dist-matrix)
    pair_t = Bond.Embed(bond-feats)
    MSAcluster, pair, state_t = Recycle.Embed(recyclables)
    ↳ only update 1st row of msa
    pair, state_t = Template.Embed(templates, pair, state)
    for i=1 to Nfull-blocks:
        msa_full, pair, state, xyz, sc-torsions =
            Full_Block(msa-full, pair, state, xyz)  ↳ 1D, 2D, 3D
            iteratively
    for j=1 to Nmain-blocks:
        MSAcluster, pair, state, xyz, sc-torsions =
            Main_Block(msacluster, pair, state, bond-feats, dist-matrix)
    for i=1 to Nref-blocks:
        state, xyz, sc-torsions = only 3D track
        Ref_Block(msacluster, pair, state, xyz,
                  idx, bond-feats, dist-matrix, ...)
    return msacluster, pair, state, xyz, sc-torsions
  
```

Main block architecture and pseudo code:

RF2



RFAA

```

def Main-Block(msacluster, pair, state, xyz, idx, bond-feats, dist-matrix):
    Str-bias = Str-Bias(xyz)  ↳ equivariant: rotation & translation
    ↳ iteratively update
    rel-pos, bond-dist = Positional-Encoding(idx, bond-feats, dist-matrix)  ↳ 2D
    Str-bias_t = linear(rel-pos) + linear(bond-dist)
    MSAcluster = 1D-Update(msacluster, pair, state, Str-bias)  ↳ column-wise and row-wise attn
    pair = Aggregation(msacluster, pair)
    pair = 2D-Update(pair, state, Str-bias) } use all 3 tracks
    xyz, state, sc-torsions = 3D-Update(msacluster, pair,
                                         State, xyz, idx, bond-feats, dist-matrix)  ↳ use all 3 tracks
    return msacluster, pair, state, xyz, sc-torsions
  
```

(ii) Strengths and Limitations

Since quantifiable and comparable data from these papers mainly focuses on model results, we will analyse prediction and confidence accuracy in this part, and leave qualitative assessment to question (e). We only have data for prediction accuracy, so the following part contains some figures from the papers.

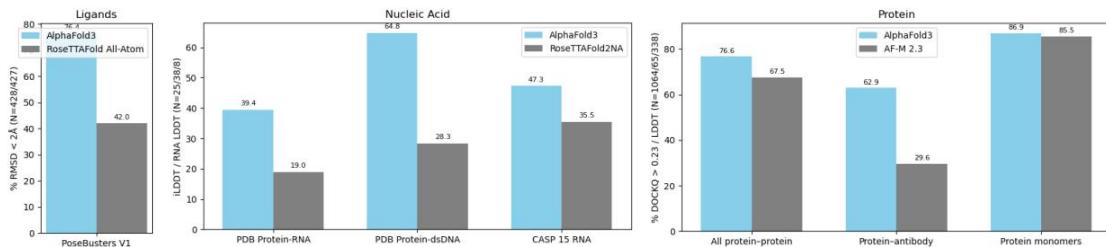


Figure 1 | Prediction accuracy Data extracted from the AF3 paper.²

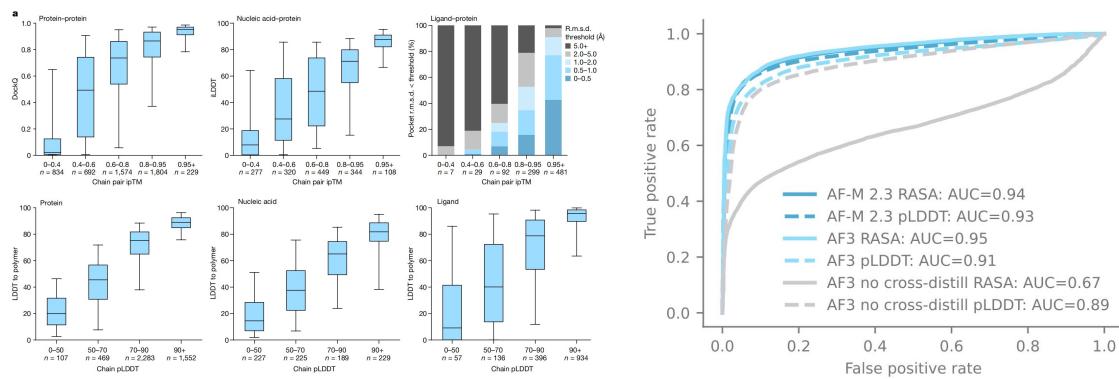


Figure 2 | AF3 confidences track accuracy

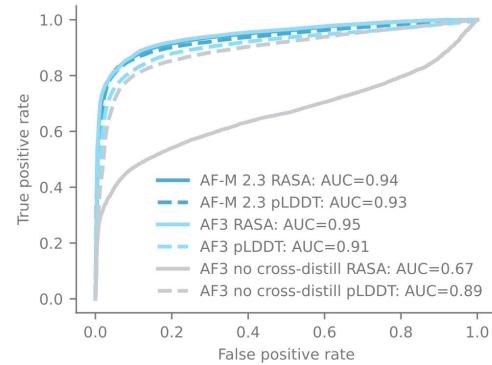


Figure 3 | Disorder prediction CAID 2

As shown in Fig.1, AF3 excels in versatility and accuracy across molecular types, especially protein–ligand (76.4% with $\text{RMSD} < 2\text{\AA}$) and protein–dsDNA (64.8% LDDT) interactions. It also has high performance on large and complex systems such as human coronavirus OC43 spike protein (4,665 residues), with LDDT=83.0.² Fig.2 shows the well-calibrated confidence metrics (ipTM and pLDDT) that correlate with model accuracy. Limitations include hallucinations (Fig.3) and stereochemical error (chirality violation rate of 4.4% in PoseBusters benchmark).

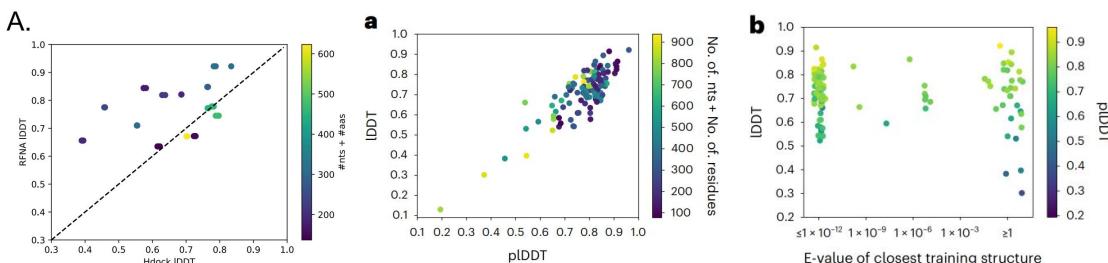


Figure 4 | Compared with docking⁵ **Figure 5 | IDDT vs plDDT⁵** **Figure 6 | IDDT of different structures⁵**

In comparison, RFNA only targets at protein–nucleic acid complexes, achieves moderate accuracy for RNA predictions, but greatly lags behind AF3 overall. RFNA was not available for baselining at the time of writing AF3 paper, but the RFNA paper indicates their accuracy is below specialist predictors in almost all categories.²

RFNA has specialized frame-torsion modeling tailored for protein–nucleic acid interactions, outperforming traditional docking methods (Fig.4). Fig.5 shows the model

correctly identifies inaccurate predictions, with 81% of high-confidence models successfully capture protein-nucleic acid interfaces by CAPRI metrics. In addition, Fig.6 shows the model generalizes well, with no clear performance difference between structures with and without sequence homologs in the protein-NA training set. For the 33 clusters with no sequence homologs, we have average IDDT of 0.68 and 24% of models > 0.8 IDDT.

One limitation of RFNA is limited protein-NA data: there are only 1,632 RNA clusters, 1,556 protein-nucleic acid clusters compared to 26,128 all protein clusters. One other is that RFNA has poor prediction of individual subunits, typically large multidomain proteins, large RNAs (>100 nt) and small single-stranded nucleic acids.

(e) Assessments

AlphaFold3 (AF3) and RoseTTAFoldDNA (RFNA) are state-of-the-art models for structural biology. While AF3 excels in versatility and accuracy, RFNA specializes in protein–nucleic acid complexes. Other solutions, such as AlphaFold-Multimer and Alchemy_RNA, also contribute significantly to specific niches within this domain. Here are some notable assessments:

Model	Targets	Performance Highlights	Strengths	Limitations
AF3	Unified model	Best in almost all categories	High versatility and accuracy, diffusion-based, robust confidence metrics	Computationally expensive, limited dynamic behavior modeling, hallucination and stereochemical error
RFNA	Protein-NA complexes	Moderate accuracy for RNA modeling	Tailored for protein-NA systems, efficient for RNA/DNA modeling	Limited training data, struggles with large RNAs and multi-domain proteins
AF-M	Protein interactions	Improved multimeric interface predictions	Accurate protein–protein interaction modeling	Less focus on nucleic acids or small molecules
Alchemy_RNA	RNA	Best performance for CASP 15 RNA dataset	Specialization in RNA modeling, superior for CASP RNA targets	Dependency on curated datasets, limited in broader applications

While each model provides a special approach to its respective problem, AF3 stands out as the SOTA model for almost all biological systems. Individually, people can dig deeper into alternatives for specialized tasks, for example fine-tune RFNA for protein-NA specific problems, use AF-M in protein engineering to aid synthetic biology, and better Alchemy_RNA for RNA-centric research.

AF3 is bound to be more significant in the history of bioinformatics. It demonstrates the possibility to accurately predict the structure of a wide range of biomolecular systems in a unified network with SOTA accuracy without homology and less dependence on MSAs. The success of AF3 shows that given the right network, we only need significantly less data and

can amplify existing data.

In the future, with the parallel developments of experimental and computational methods, more training data and improved deep learning models based on AF3 will propel us into an era of structurally informed biological understanding and therapeutic development.

(f) Suggestions

① Integrate dynamics

One limitation of protein structure prediction models is that they typically predict static structures as seen in the PDB, not the dynamical behaviour of biomolecular systems in solution,² where most biological activities and reactions take place. To address this problem, the development of hybrid models integrating molecular dynamics (MD) simulations or neural approximations of dynamics is a promising field. Such models would not only predict static structures but also capture conformational ensembles and dynamic states, essential for understanding biological processes.

A feasible idea is to integrate recent advancements in neural network-based simulation tools, such as Graph Neural Networks (GNNs), which can approximate dynamics. This would significantly enhance the biological realism of these models, enabling them to represent intermediate states. Such models can be applied to areas where dynamics and flexibility are paramount, such as drug discovery and RNA folding. Additionally, hybrid models could complement existing static prediction methods, refining or validating predictions with dynamic insights for greater accuracy and utility.

However, this approach is not without its challenges. Incorporating MD simulations would require substantial computational resources, particularly for large systems, posing a barrier to scalability. Furthermore, early implementations may face difficulties in balancing computational feasibility with prediction precision, potentially leading to accuracy trade-offs.

Despite these limitations, the integration of dynamic conformational states into predictive frameworks has the potential to revolutionize the field, bridging the gap between static and dynamic representations of biomolecular systems, and unlocking new possibilities for therapeutic and biotechnological applications.

② Mutual distillation

Given that structure prediction models are an application of the Transformer structure in Biology, we can draw inspiration for them from the more general models in computer science, such as Large Language Models (LLMs) and Large Multimodal Models (LMMs). In this case, we can improve model performance and reduce computational cost by bootstrapping and leveraging mutual distillation.

Unified models excel in generalization and specialized models achieve high accuracy in specific domains. Mutual distillation addresses these complementary strengths and weaknesses by iteratively refining each other's predictions through a synergistic process.

The workflow begins with the unified model generating broad-spectrum predictions, which the specialized model then refines in its domain of expertise. The improved results from the specialized model are distilled back into the unified model's training process, allowing it to learn domain-specific insights without requiring additional architecture modifications. This iterative bootstrapping ensures that both models benefit from each other's strengths.

This approach has several expected advantages. By combining the two models, mutual distillation enhances the overall accuracy of predictions without sacrificing versatility. It also improves the efficiency of training by reusing predictions as additional training data. Furthermore, this method is modular, enabling the flexibility to add new specialized models into the framework without retraining the entire unified model.

However, there are challenges to consider. Careful model design is essential to ensure convergence, avoiding issues like overfitting in specialized models or over-generalization in the unified model. Additionally, integrating predictions from models trained on different datasets and architectures may require sophisticated data harmonization techniques.

Despite these limitations, mutual distillation offers a powerful mechanism to improve accuracy, scalability, and utility across biomolecular modeling tasks, paving the way for more advanced predictive tools in structural biology and therapeutic design.

1. Jumper, J., Evans, R., Pritzel, A. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* 596, 583–589 (2021). <https://doi.org/10.1038/s41586-021-03819-2>
2. Abramson, J., Adler, J., Dunger, J. et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* 630, 493–500 (2024). <https://doi.org/10.1038/s41586-024-07487-w>
3. Minkyung Baek et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* 373, 871–876 (2021). DOI:10.1126/science.abj8754
4. Minkyung Baek et al. Efficient and accurate prediction of protein structure using RoseTTAFold2. *bioRxiv* 2023.05.24.542179; doi: <https://doi.org/10.1101/2023.05.24.542179>
5. Baek, M., McHugh, R., Anishchenko, I. et al. Accurate prediction of protein–nucleic acid complexes using RoseTTAFoldNA. *Nat Methods* 21, 117–121 (2024). <https://doi.org/10.1038/s41592-023-02086-5>
6. Rohith Krishna et al. Generalized biomolecular modeling and design with RoseTTAFold All-Atom. *Science* 384, eadl2528 (2024). DOI:10.1126/science.adl2528