

# Highly Accurate Protein Structure Prediction with AlphaFold

---

报告人：陈星强



# 论文作者信息

---

## Affiliations

- **DeepMind, London, UK**  
John Jumper, Richard Evans, Alexander Pritzel, Tim Green,
- **Seoul National University, Seoul, South Korea**  
Martin Steinegger

Correspondence to [John Jumper](#) or [Demis Hassabis](#).

## Subjects

- [Computational biophysics](#)
- [Machine learning](#)
- [Protein structure predictions](#)
- [Structural biology](#)

# 摘要部分

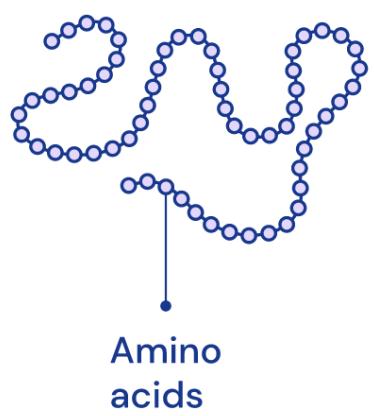
---

## Abstract

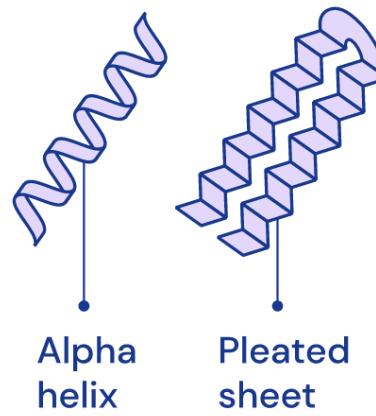
# Abstract

---

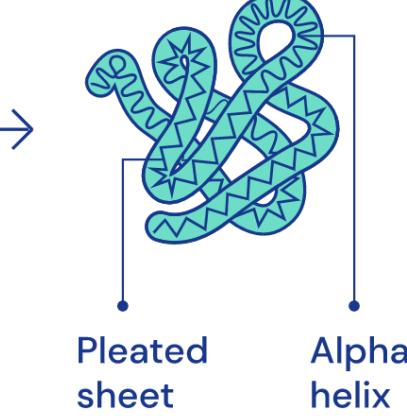
Every protein is made up of a sequence of amino acids bonded together



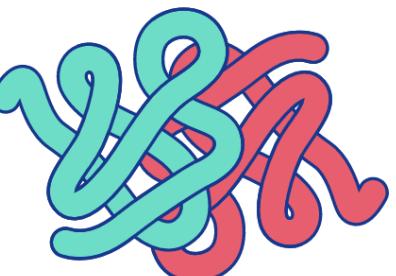
These amino acids interact locally to form shapes like helices and sheets



These shapes fold up on larger scales to form the full three-dimensional protein structure

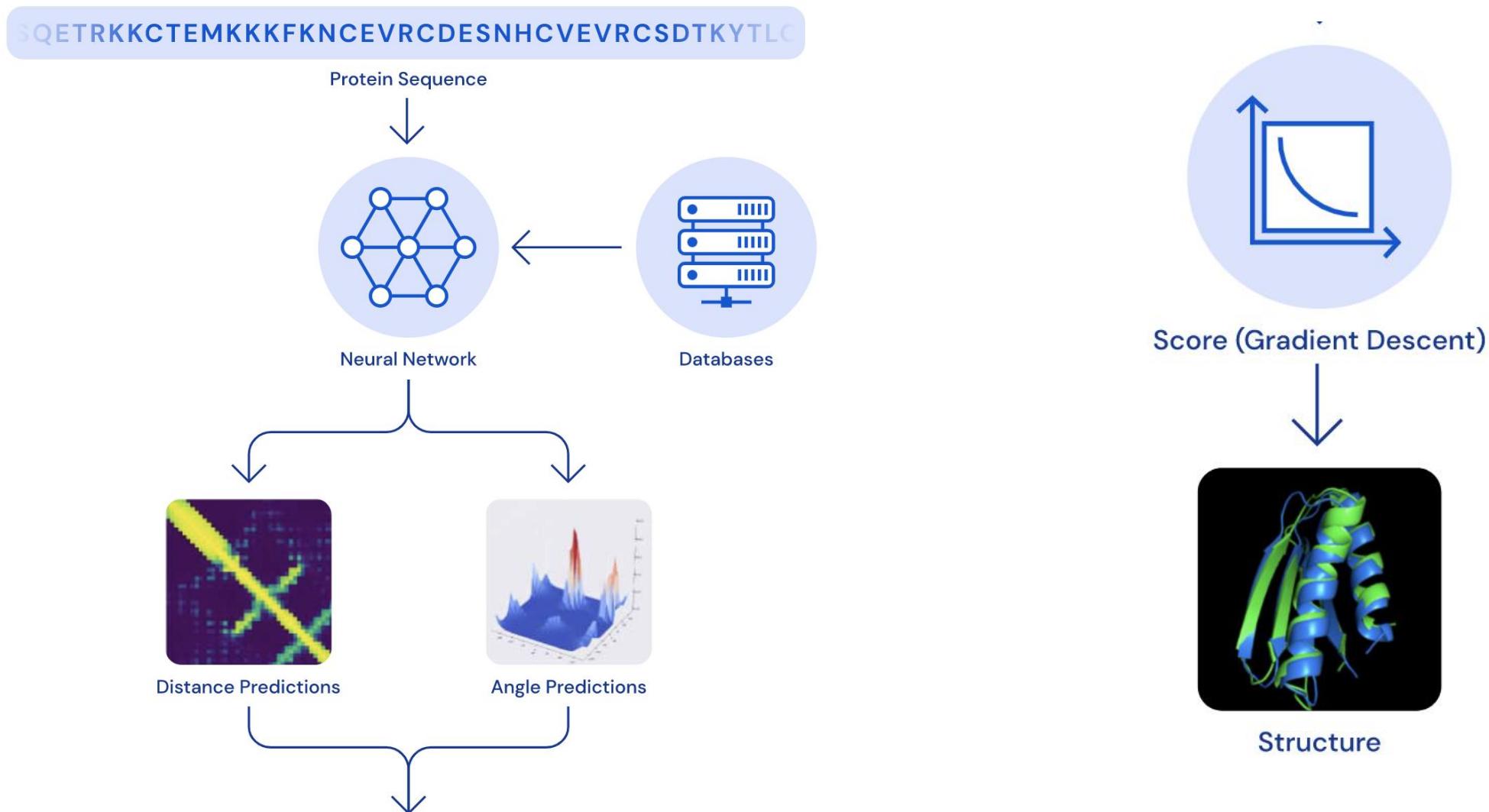


Proteins can interact with other proteins, performing functions such as signalling and transcribing DNA



# Abstract

---



A schematic of the architecture of the AlphaFold system predicting structure from protein sequence.

# 基础概念学习

---

## Notation

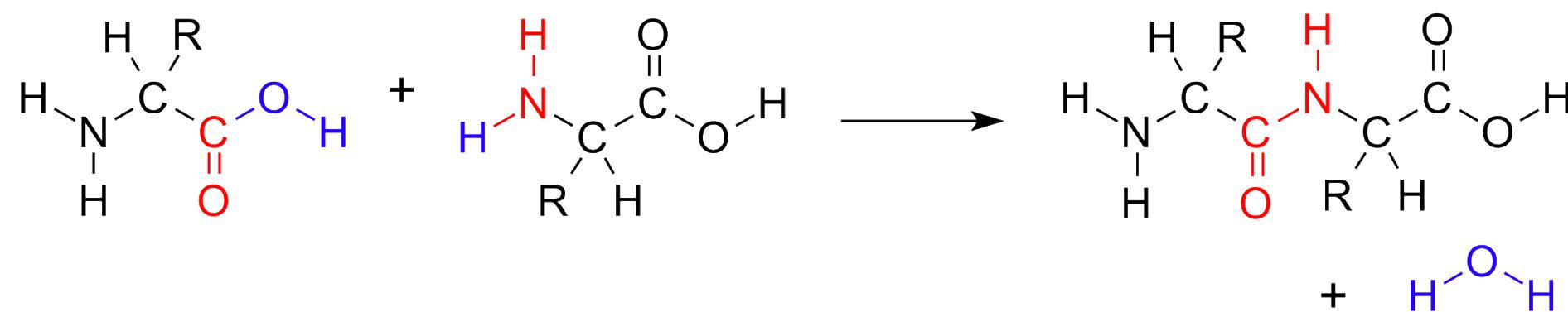
# 基础概念学习

---

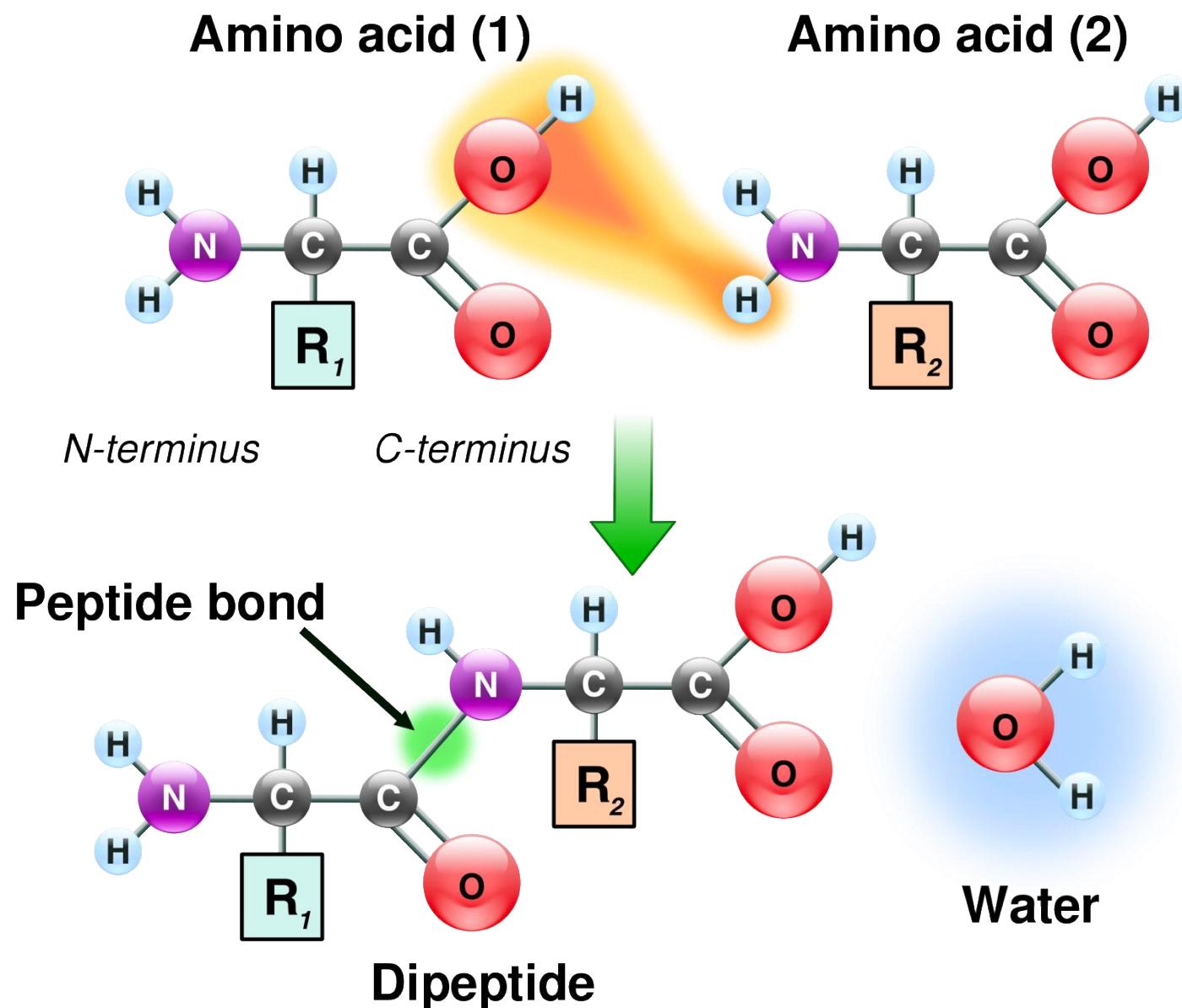
1. Amino acid
2. Peptide bond
3. Chi Angle
4. template modeling score (TM-score)
5. multiple sequence alignments (MSAs)
6. pairwise features
7. homologues
8. IPA
9. Evoformer
10. pLDDT per-residue accuracy of the structure
11. TM-score pTM
12. Pair-wise error prediction
13. Frame-aligned point error FAPE

## 2. 肽键

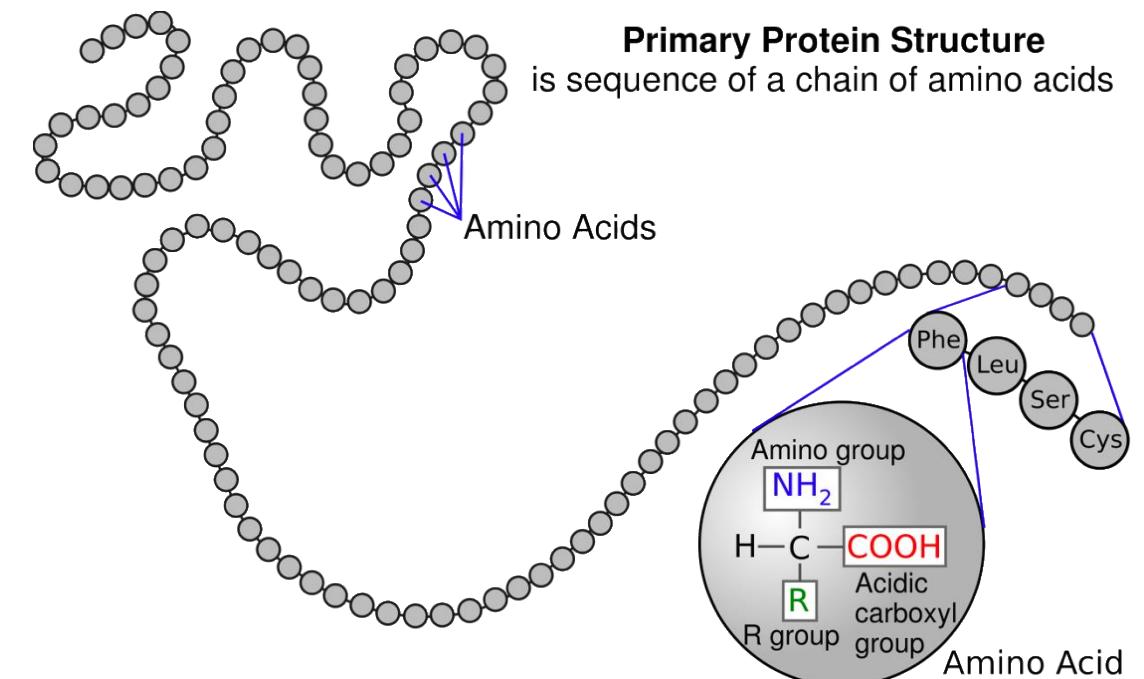
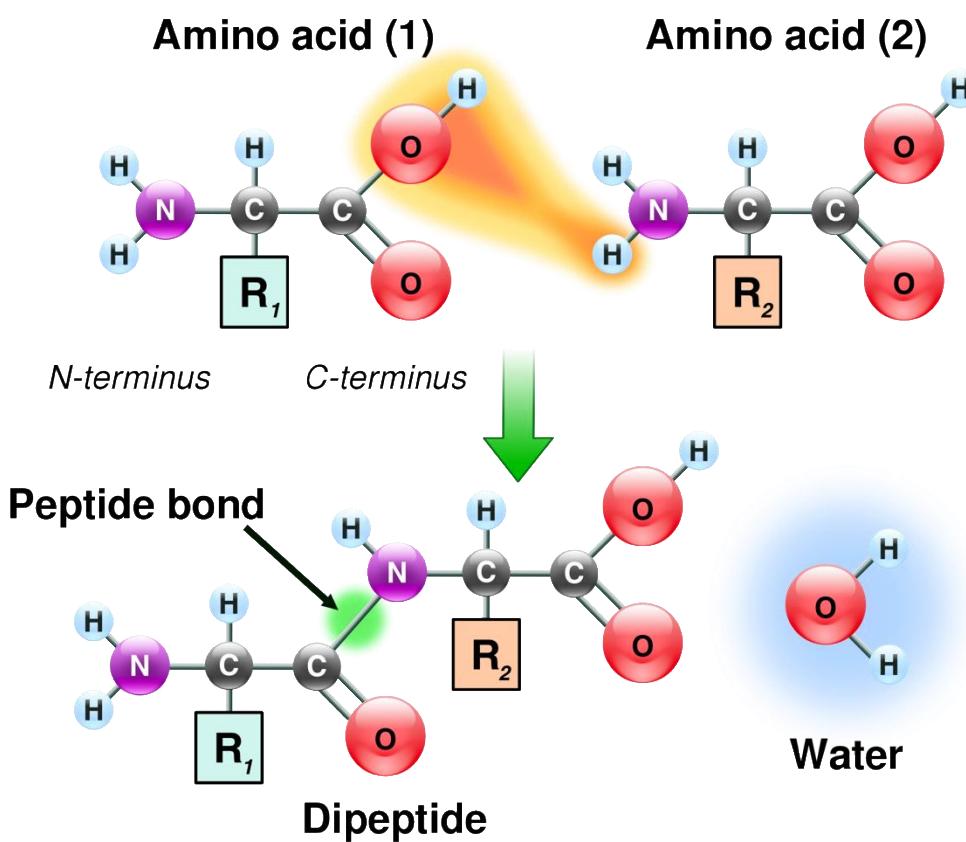
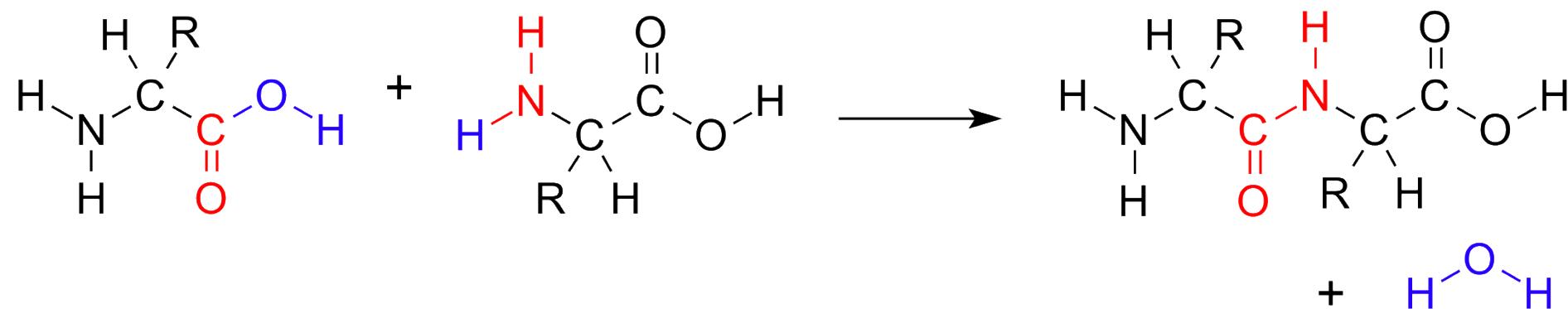
---



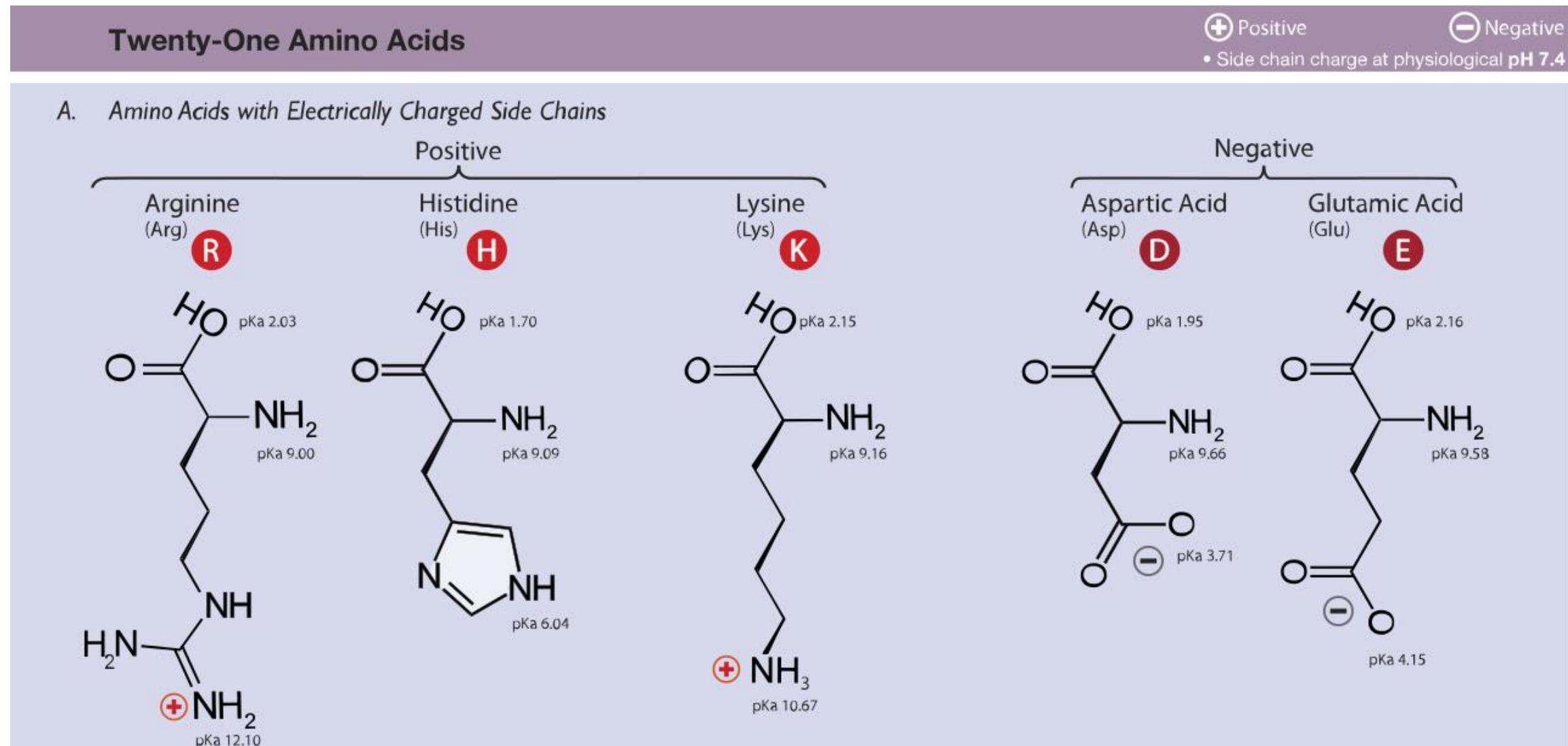
# 基础信息 概念解释 肽键



# 基础信息 概念解释 肽键



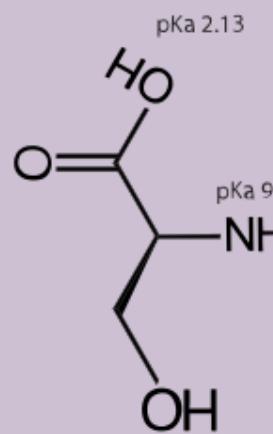
# 1. 氨基酸



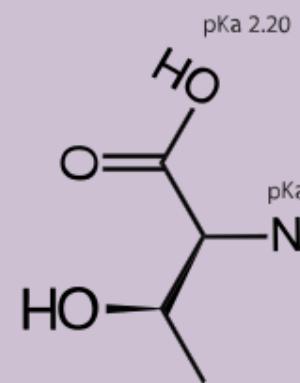
# 基础信息 概念解释 氨基酸

## B. Amino Acids with Polar Uncharged Side Chains

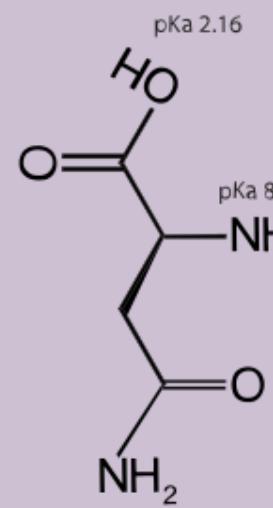
Serine  
(Ser) **S**



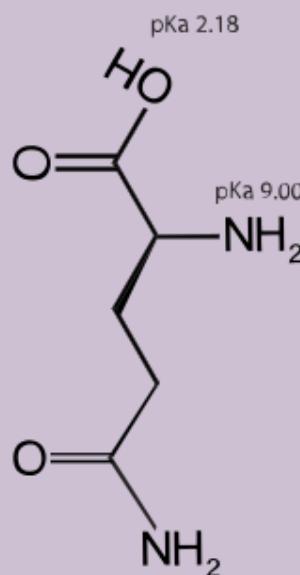
Threonine  
(Thr) **T**



Asparagine  
(Asn) **N**

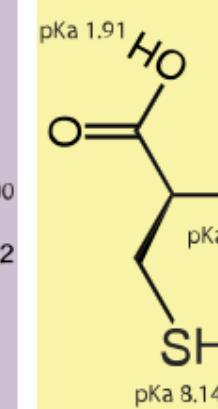


Glutamine  
(Gln) **Q**

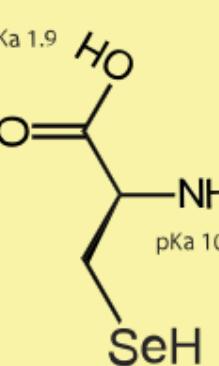


## C. Special Cases

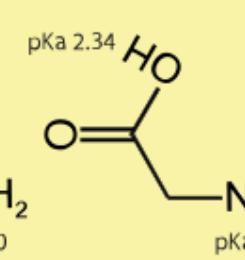
Cysteine  
(Cys) **C**



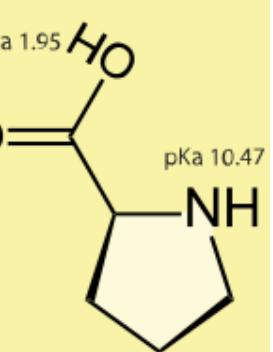
Selenocysteine  
(Sec) **U**



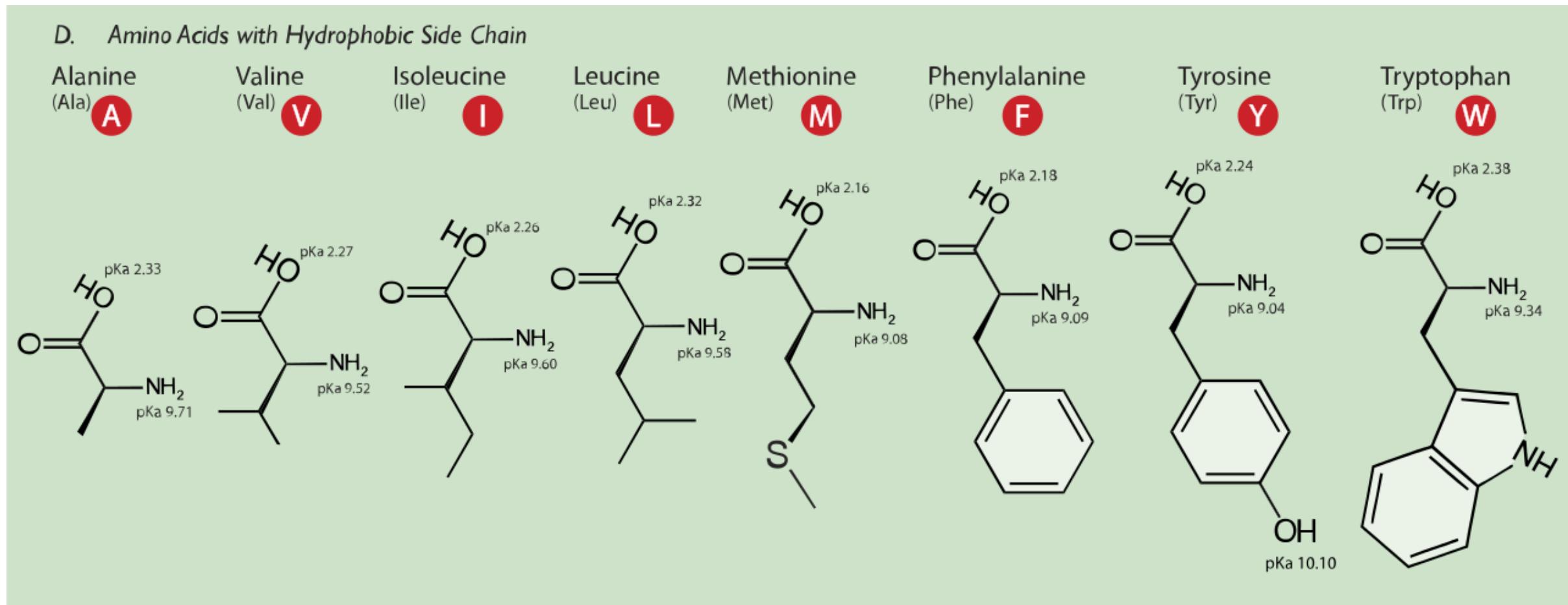
Glycine  
(Gly) **G**



Proline  
(Pro) **P**

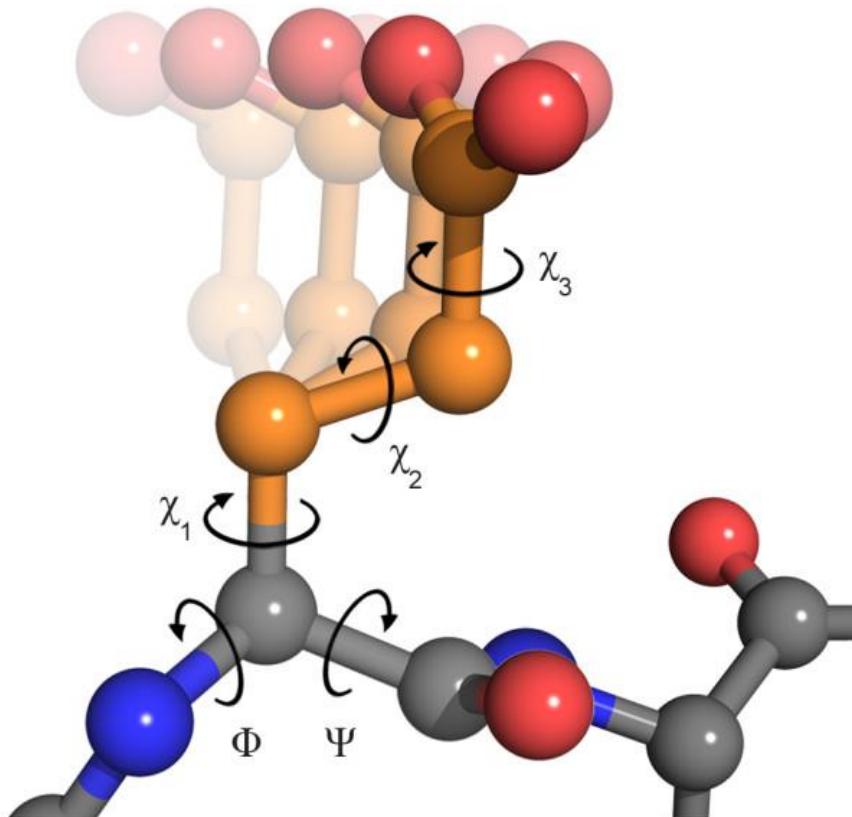


# 基础信息 概念解释 氨基酸



### 3. Chi Angle

---



Dihedral angles in glutamate

<http://www.biomedcentral.com/1471-2105/11/306>

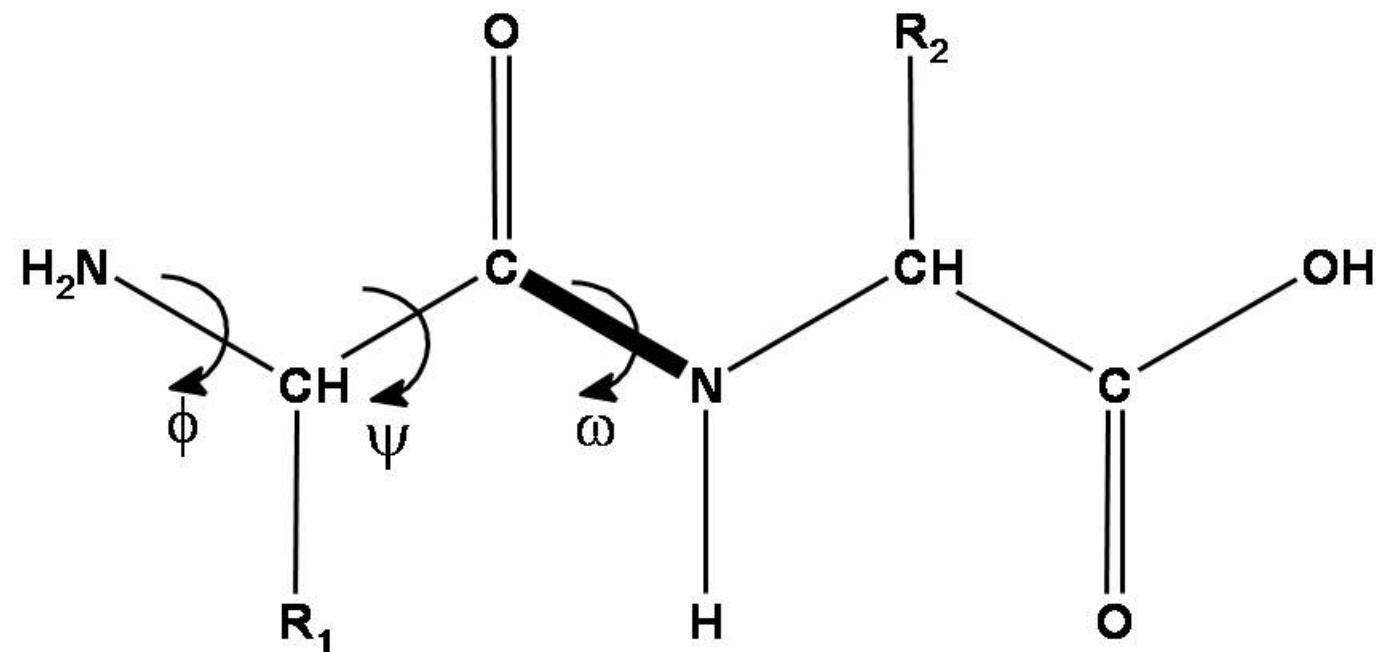
<https://swissmodel.expasy.org/course/text/chapter3.htm>

<http://www.mlb.co.jp/linux/science/garlic/doc/commands/dihedrals.html>

<http://www.mlb.co.jp/linux/science/garlic/doc/commands/dihedrals.html>

### 3. Chi Angle

---



#### The angle between two bonds

The backbone of a protein has three different torsion angles.

- The phi-angle ( $\phi$ ) - around the N-C $\alpha$  bond
- The psi-angle ( $\psi$ ) - around the C $\alpha$ -C bond
- The omega-angle ( $\omega$ ) - around the peptide bond between C and N.

### 3. Chi Angle

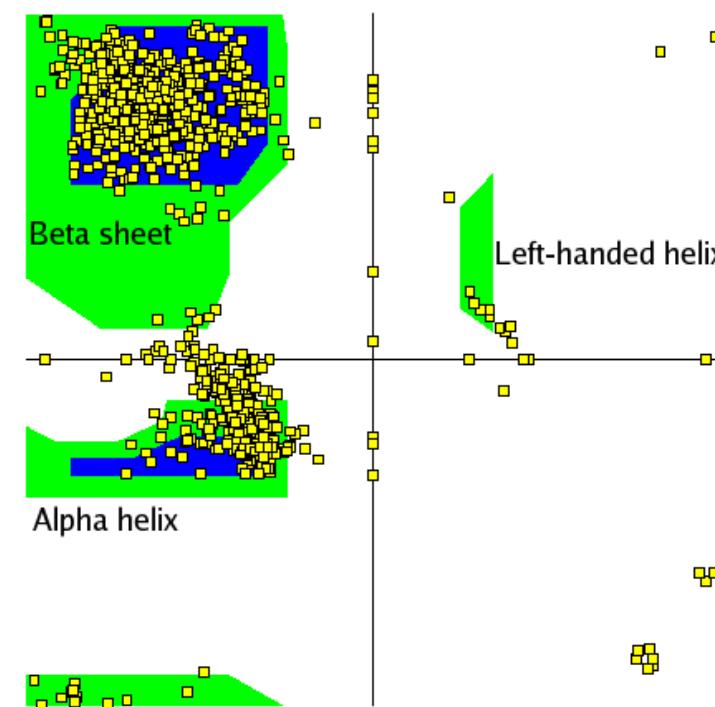
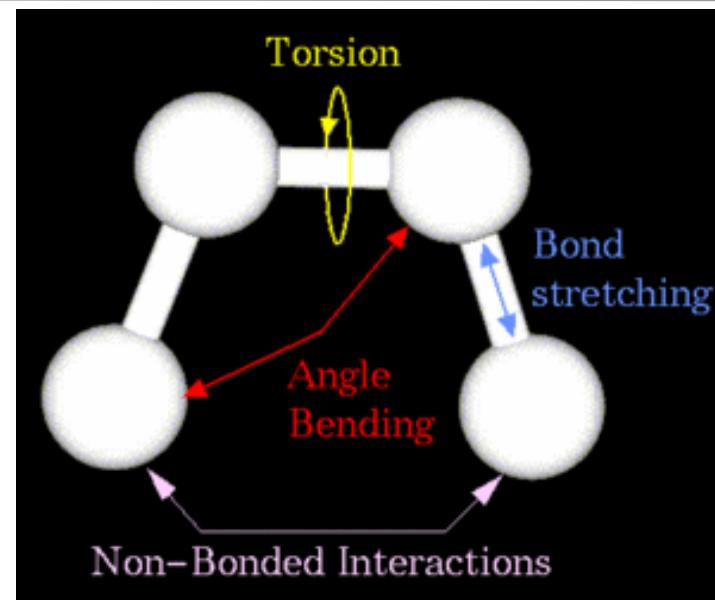
The  $\omega$ -bond has a slightly double-bond character and is therefore almost always 180 degrees.

The structure of a protein is mainly formed by the  $\phi$ - and  $\psi$ -angles.

Every [secondary structure](#) element forces the [backbone](#) into a specific range of torsion angles, this can be visualized in a [ramachandran plot](#).

The figure shows the location of the  $\phi$ ,  $\psi$  and  $\omega$ -angles.

[Side chains](#) can also contain torsion angles, they are indicated as [chi-angles](#).



# 3. Chi Angle

**Table 2** | Rigid groups for constructing all atoms from given torsion angles. Boxes highlight groups that are symmetric under 180° rotations.

aatype	bb	$\psi$	$\chi_1$	$\chi_2$	$\chi_3$	$\chi_4$
ALA	N, C $^{\alpha}$ , C, C $^{\beta}$	O	-	-	-	-
ARG	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta}$	N $^{\epsilon}$	N $^{\eta 1}$ , N $^{\eta 2}$ , C $^{\zeta}$
ASN	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	N $^{\delta 2}$ , O $^{\delta 1}$	-	-
ASP	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	O $^{\delta 1}$ , O $^{\delta 2}$	-	-
CYS	N, C $^{\alpha}$ , C, C $^{\beta}$	O	S $^{\gamma}$	-	-	-
GLN	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta}$	N $^{\epsilon 2}$ , O $^{\epsilon 1}$	-
GLU	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta}$	O $^{\epsilon 1}$ , O $^{\epsilon 2}$	-
GLY	N, C $^{\alpha}$ , C	O	-	-	-	-
HIS	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta 2}$ , N $^{\delta 1}$ , C $^{\epsilon 1}$ , N $^{\epsilon 2}$	-	-
ILE	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma 1}$ , C $^{\gamma 2}$	C $^{\delta 1}$	-	-
LEU	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta 1}$ , C $^{\delta 2}$	-	-
LYS	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta}$	C $^{\epsilon}$	N $^{\zeta}$
MET	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	S $^{\delta}$	C $^{\epsilon}$	-
PHE	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta 1}$ , C $^{\delta 2}$ , C $^{\epsilon 1}$ , C $^{\epsilon 2}$ , C $^{\zeta}$	-	-
PRO	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta}$	-	-
SER	N, C $^{\alpha}$ , C, C $^{\beta}$	O	O $^{\gamma}$	-	-	-
THR	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma 2}$ , O $^{\gamma 1}$	-	-	-
TRP	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta 1}$ , C $^{\delta 2}$ , C $^{\epsilon 2}$ , C $^{\epsilon 3}$ , N $^{\epsilon 1}$ , C $^{\eta 2}$ , C $^{\zeta 2}$ , C $^{\zeta 3}$	-	-
TYR	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma}$	C $^{\delta 1}$ , C $^{\delta 2}$ , C $^{\epsilon 1}$ , C $^{\epsilon 2}$ , O $^{\eta}$ , C $^{\zeta}$	-	-
VAL	N, C $^{\alpha}$ , C, C $^{\beta}$	O	C $^{\gamma 1}$ , C $^{\gamma 2}$	-	-	-

# 基础概念学习

---

**N<sub>res</sub>** : the number of residues in the input primary sequence (cropped during training).

**N<sub>templ</sub>** : the number of templates used in the model.

**N<sub>all\_seq</sub>** : the number of all available MSA sequences.

**N<sub>clust</sub>** : the number of clusters after MSA clustering.

**N<sub>seq</sub>** : the number of sequences processed in the MSA stack (where  $N_{seq} = N_{clust} + N_{templ}$ ).

**N<sub>extra\_seq</sub>** : the number of unclustered MSA sequences (after sub-sampling).

**N<sub>block</sub>** : the number of blocks in Evoformer-like stacks.

**N<sub>ensemble</sub>** : the number of ensembling iterations.

**N<sub>cycle</sub>** : the number of recycling iterations.

# 基础概念学习

---

**Table 4** | AlphaFold training protocol. We train each stage until convergence with the approximate timings and number of samples provided.

Model	Initial training	Fine-tuning
Number of templates $N_{\text{templ}}$	4	4
Sequence crop size $N_{\text{res}}$	256	384
Number of sequences $N_{\text{seq}}$	128	512
Number of extra sequences $N_{\text{extra\_seq}}$	1024	5120
Parameters initialized from	Random	Initial training
Initial learning rate	$10^{-3}$	$5 \cdot 10^{-4}$
Learning rate linear warm-up samples	128000	0
Structural violation loss weight	0.0	1.0
Training samples ( $\cdot 10^6$ )	$\approx 10$	$\approx 1.5$
Training time	$\approx 7$ days	$\approx 4$ days

# 算法列表

---

# List of Algorithms

---

### List of Algorithms

	Function Name	Instruction
1	MSABlockDeletion	<b>MSA block deletion</b>
2	Inference	<b>AlphaFold Model Inference</b>
3	InputEmbedder	<b>Embeddings for initial representations</b>
4	relpos	<b>Relative position encoding</b>
5	one_hot	<b>One-hot encoding with nearest bin</b>
6	EvoformerStack	<b>Evoformer stack</b>
7	MSARowAttentionWithPairBias	<b>MSA row-wise gated self-attention with pair bias</b>
8	MSAColumnAttention	<b>MSA column-wise gated self-attention</b>
9	MSATransition	<b>Transition layer in the MSA stack</b>
10	OuterProductMean	<b>Outer product mean</b>
11	TriangleMultiplicationOutgoing	<b>Triangular multiplicative update using “outgoing” edges</b>
12	TriangleMultiplicationIncoming	<b>Triangular multiplicative update using “incoming” edges</b>
13	TriangleAttentionStartingNode	<b>Triangular gated self-attention around starting node</b>
14	TriangleAttentionEndingNode	<b>Triangular gated self-attention around ending node</b>
15	PairTransition	<b>Transition layer in the pair stack</b>

# AlphaFold Model Inference

---

List of Algorithms

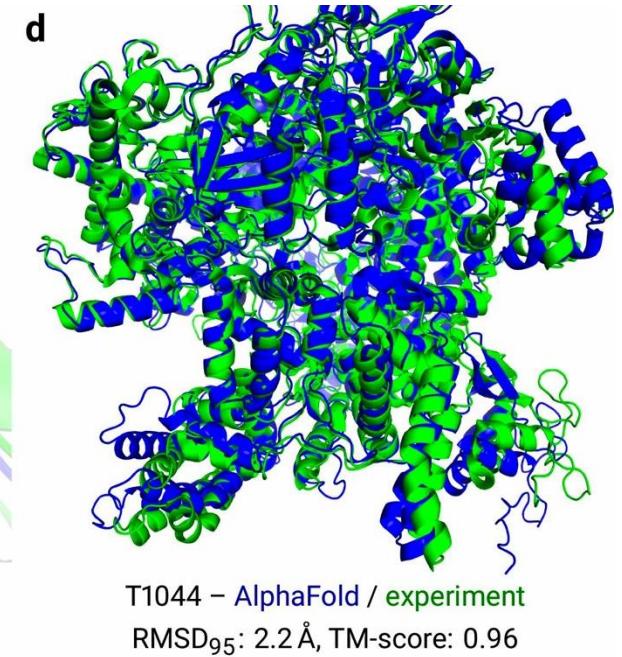
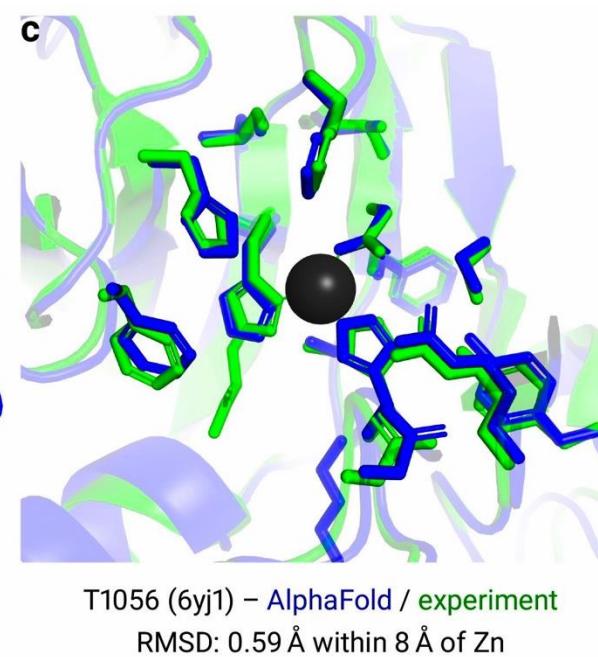
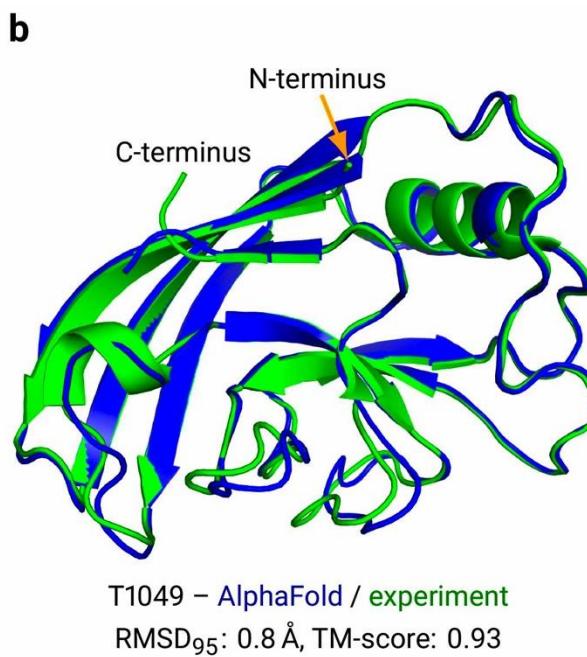
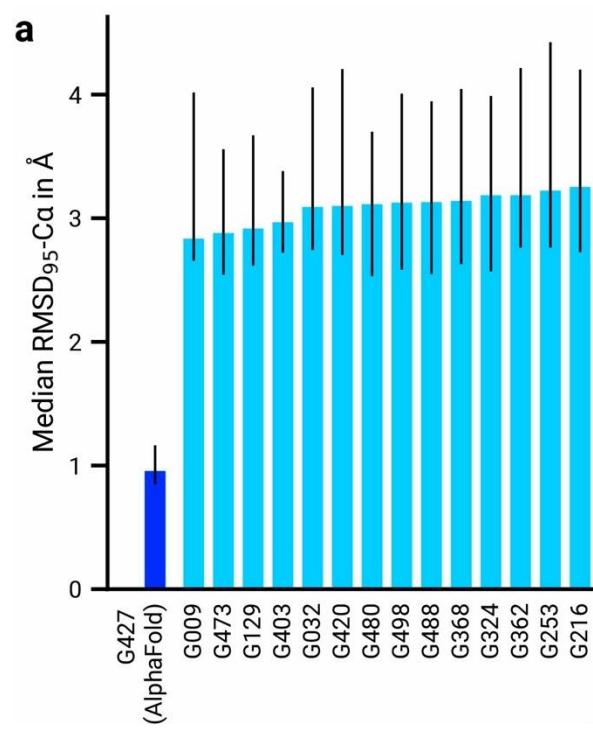
	Function Name	Instruction
16	TemplatePairStack	<b>Template pair stack</b>
17	TemplatePointwiseAttention	<b>Template pointwise attention</b>
18	ExtraMsaStack	<b>Extra MSA stack</b>
19	MSAColumnGlobalAttention	<b>MSA global column-wise gated self-attention</b>
20	StructureModule	<b>Structure module</b>
21	rigidFrom3Points	<b>Rigid from 3 points using the Gram–Schmidt process</b>
22	InvariantPointAttention	<b>Invariant point attention(IPA)</b>
23	BackboneUpdate	<b>Backbone update</b>
24	computeAllAtomCoordinates	<b>Compute all atom coordinates</b>
25	makeRotX	<b>Make a transformation that rotates around the x-axis</b>
26	renameSymmetricGroundTruthAtoms	<b>Rename symmetric ground truth atoms</b>
27	torsionAngleLoss	<b>Side chain and backbone torsion angle loss</b>
28	computeFAPE	<b>Compute the Frame aligned point error</b>
29	predictPerResidueLDDT	<b>Predict model confidence pLDDT</b>
30	RecyclingInference	<b>Generic recycling inference procedure</b>
31	RecyclingTraining	<b>Generic recycling training procedure</b>
32	RecyclingEmbedder	<b>Embedding of Evoformer and Structure module outputs for recycling</b>

# 论文背景介绍

---

## Introduction

# 基础信息 模型效果

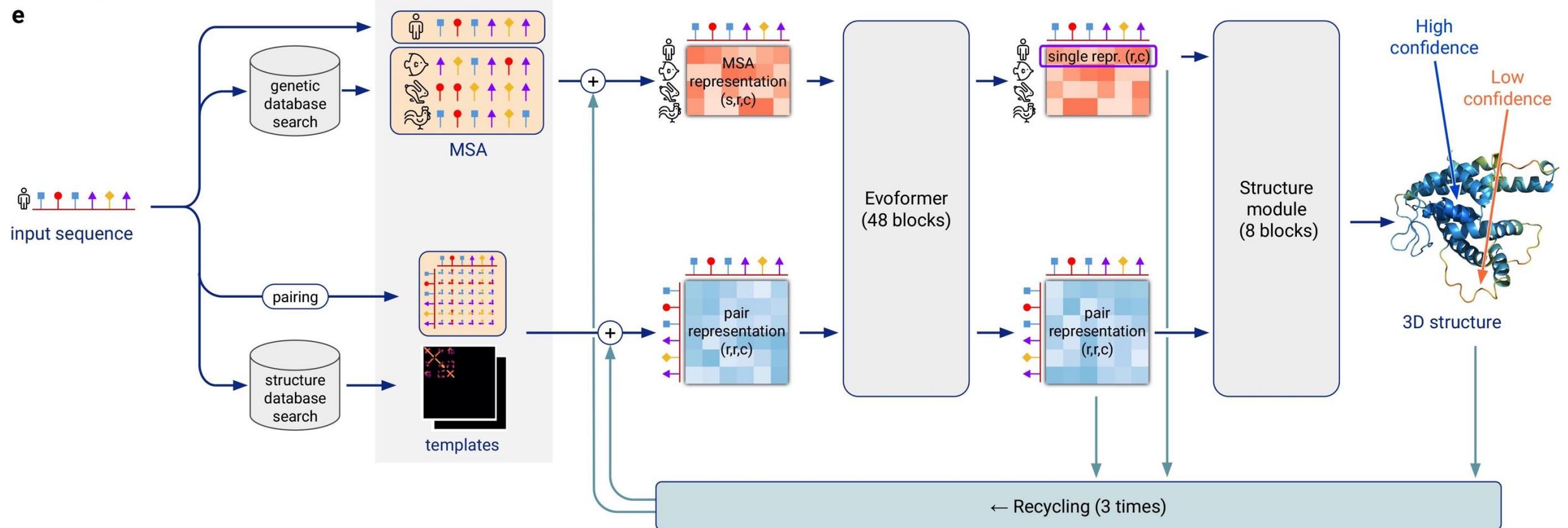


# 论文方法介绍

---

## AlphaFold network

# 基础信息 模型框架

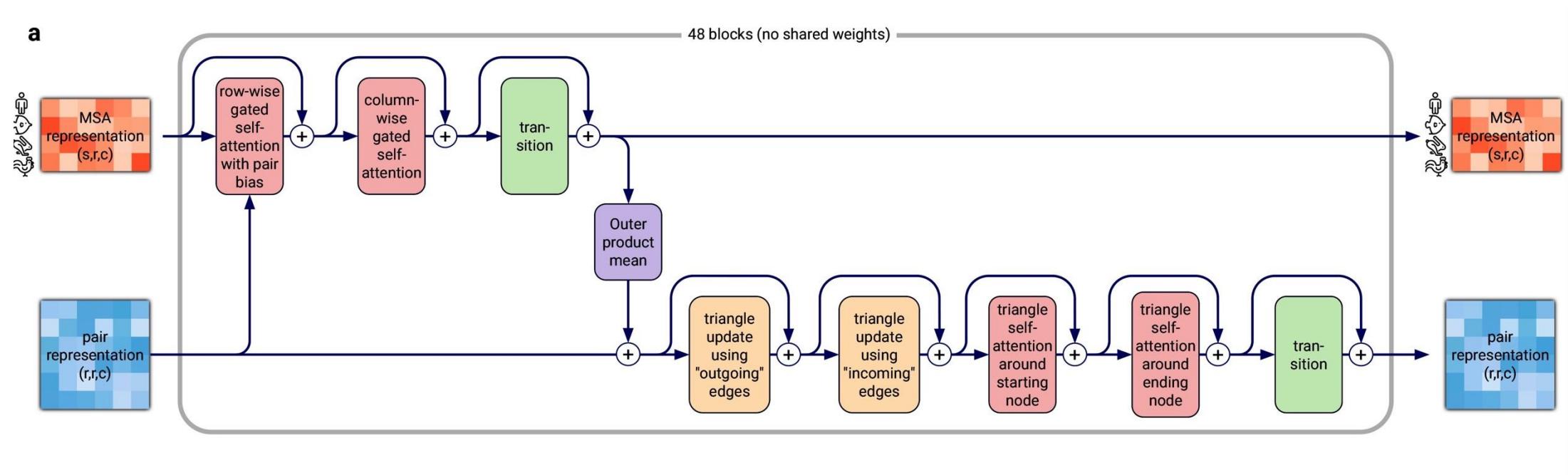


# 核心算法模块 Evoformer

---

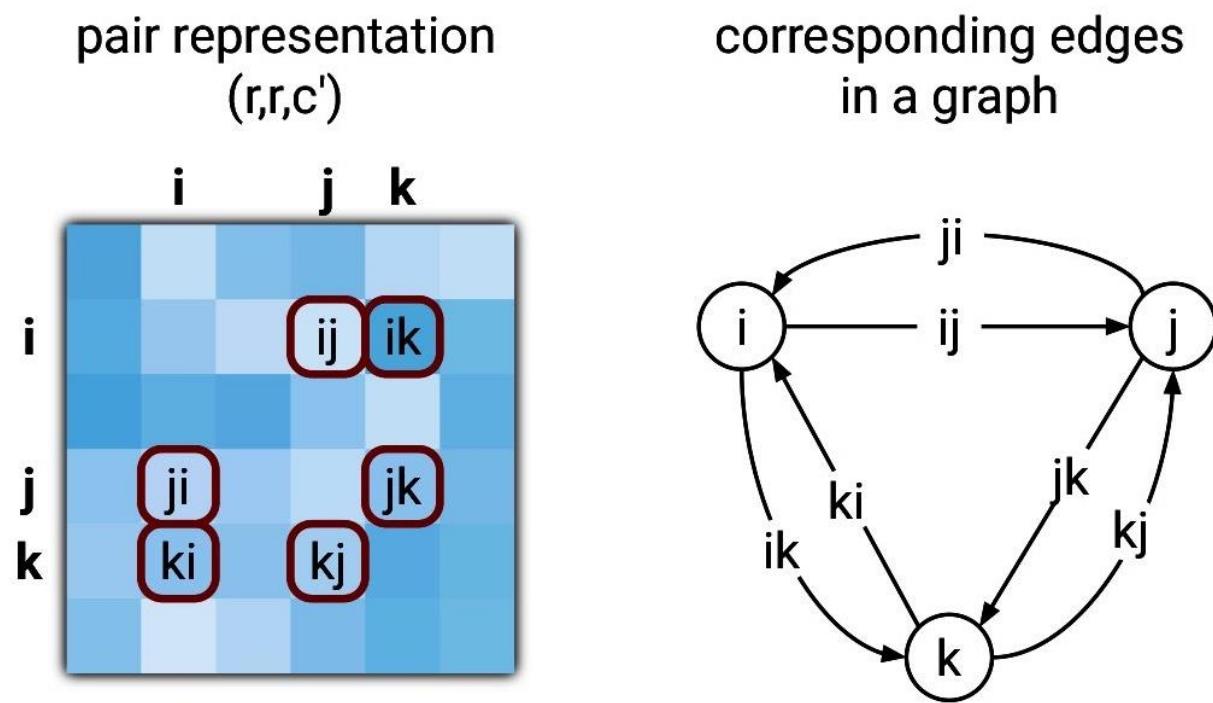
**Evoformer**

# 模型模块 Evoformer



# 模型模块 Evoformer

---

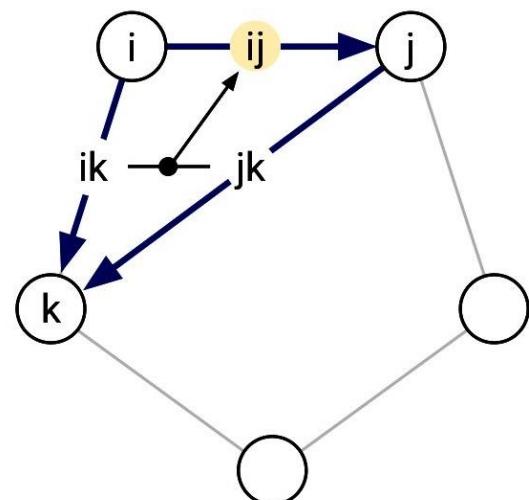


The pair representation interpreted as directed edges in a graph.

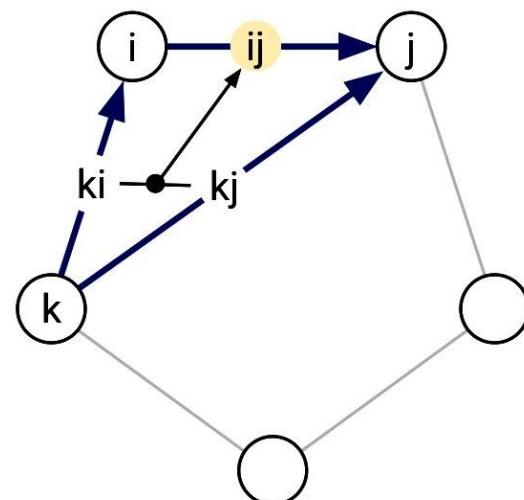
# 模型模块 Evoformer

c

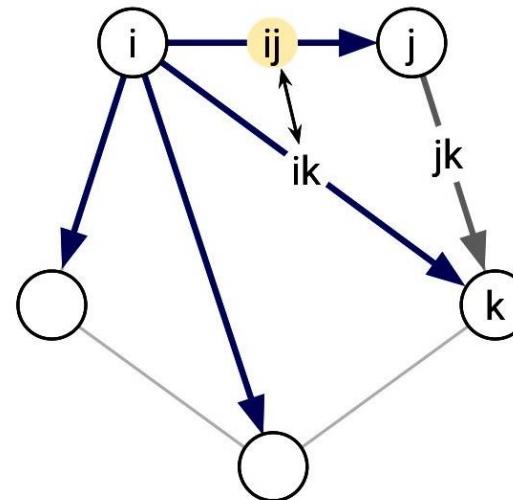
Triangle multiplicative update using "outgoing" edges



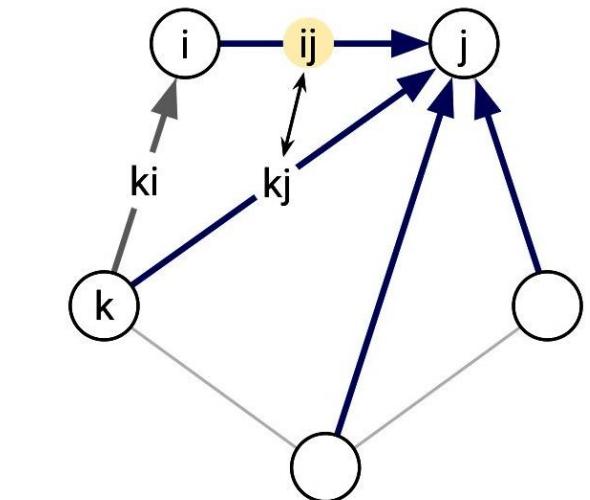
Triangle multiplicative update using "incoming" edges



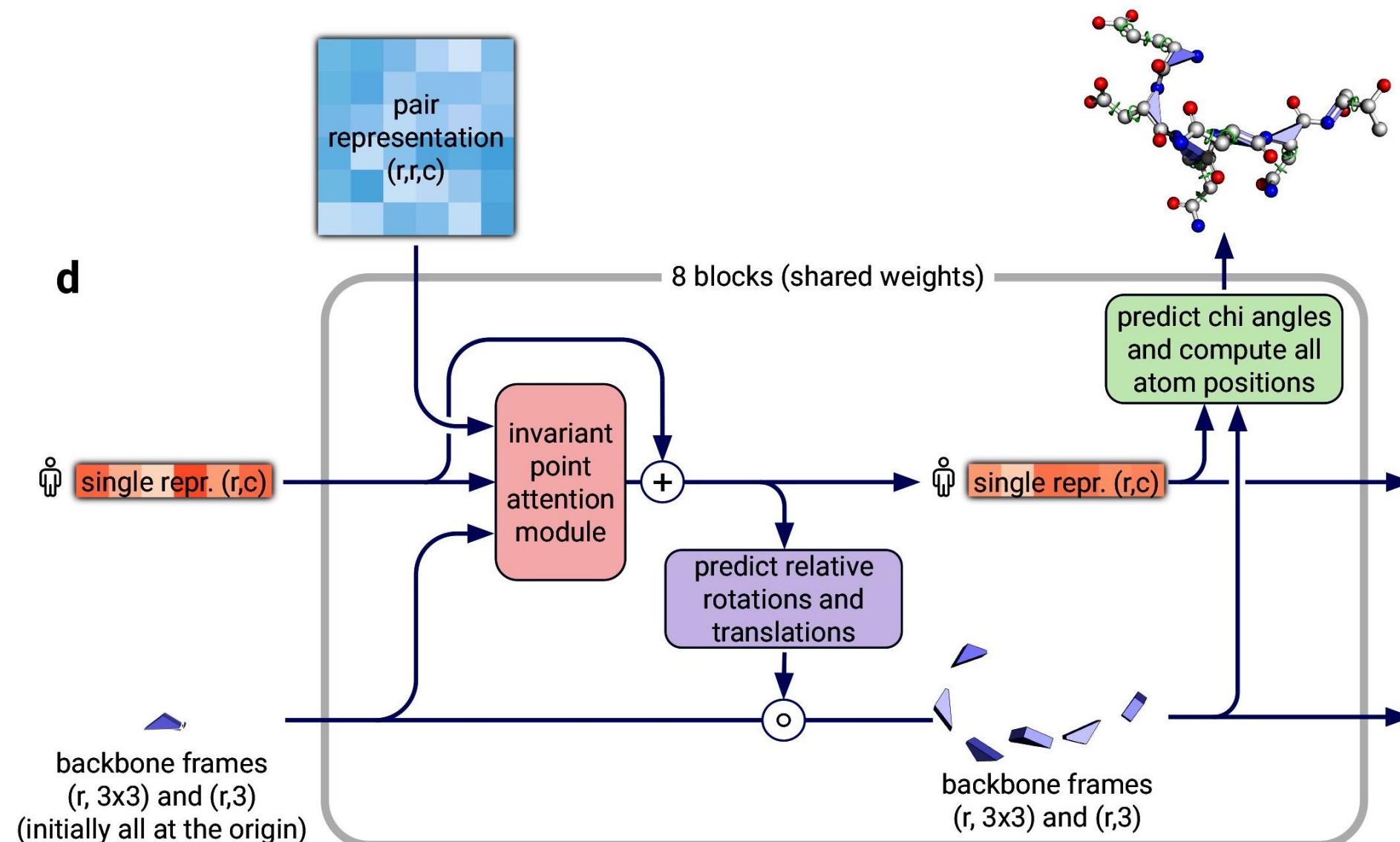
Triangle self-attention around starting node



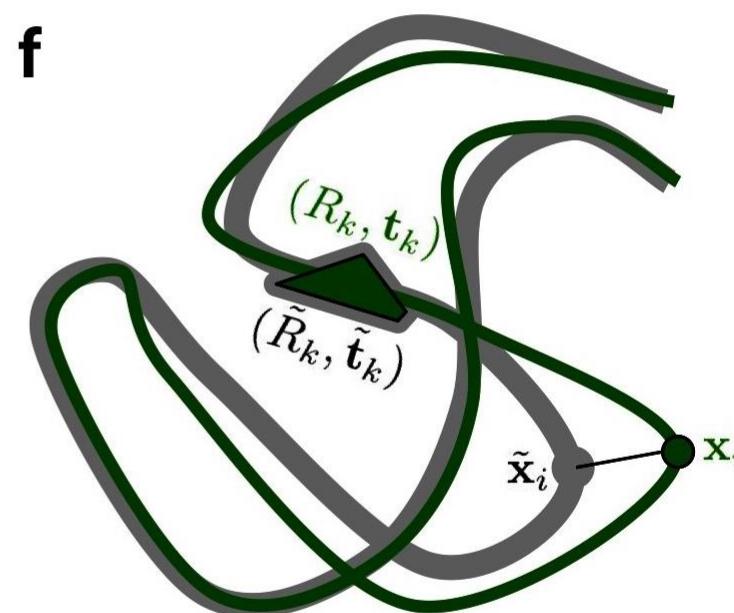
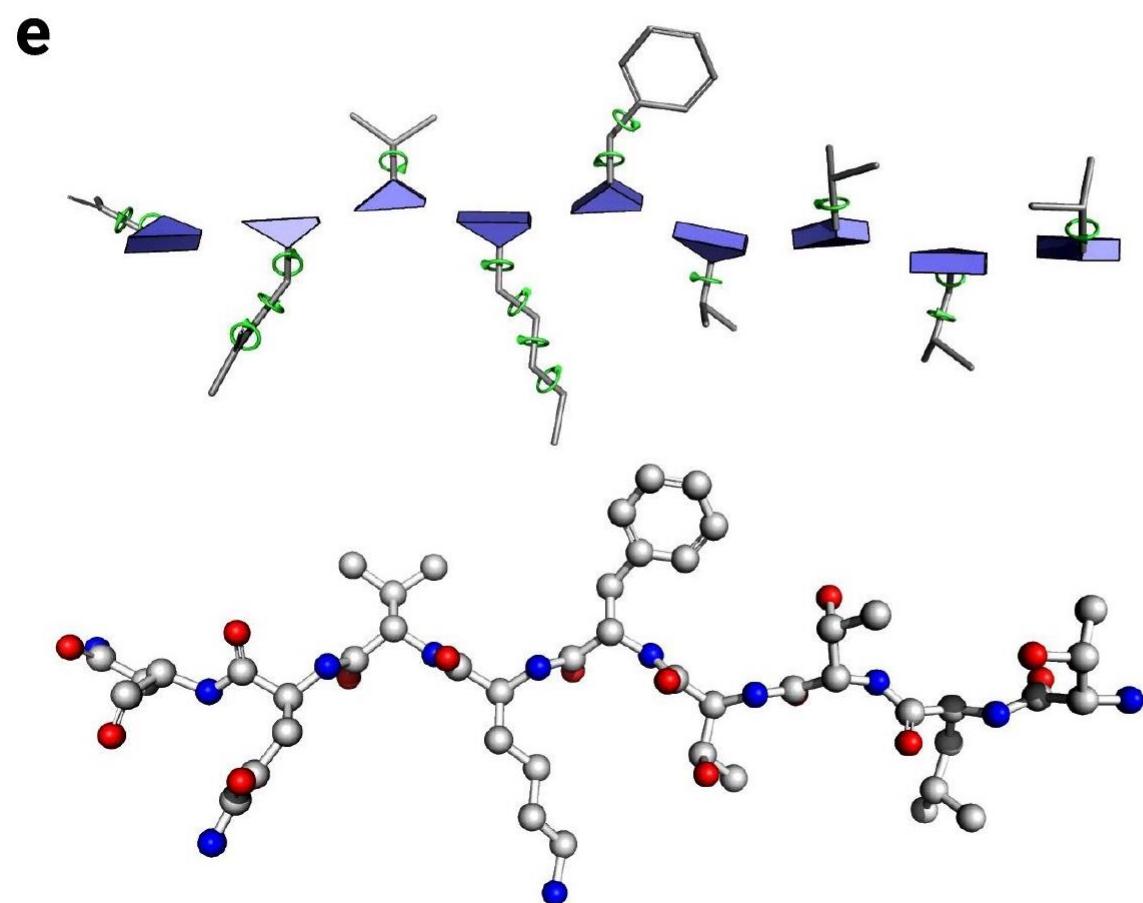
Triangle self-attention around ending node



# 模型模块 Evoformer



# 模型模块 Evoformer

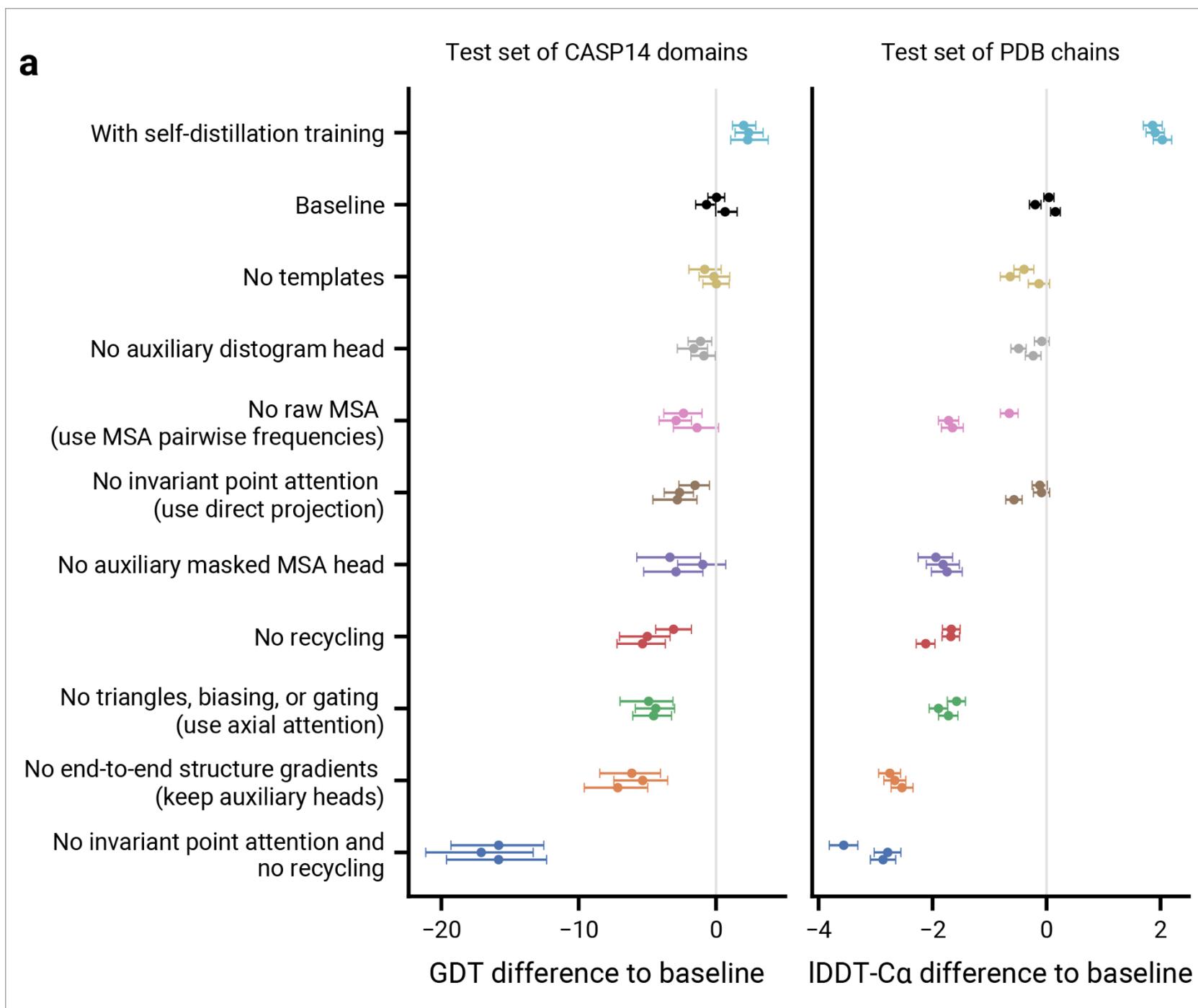


# 核心算法框架

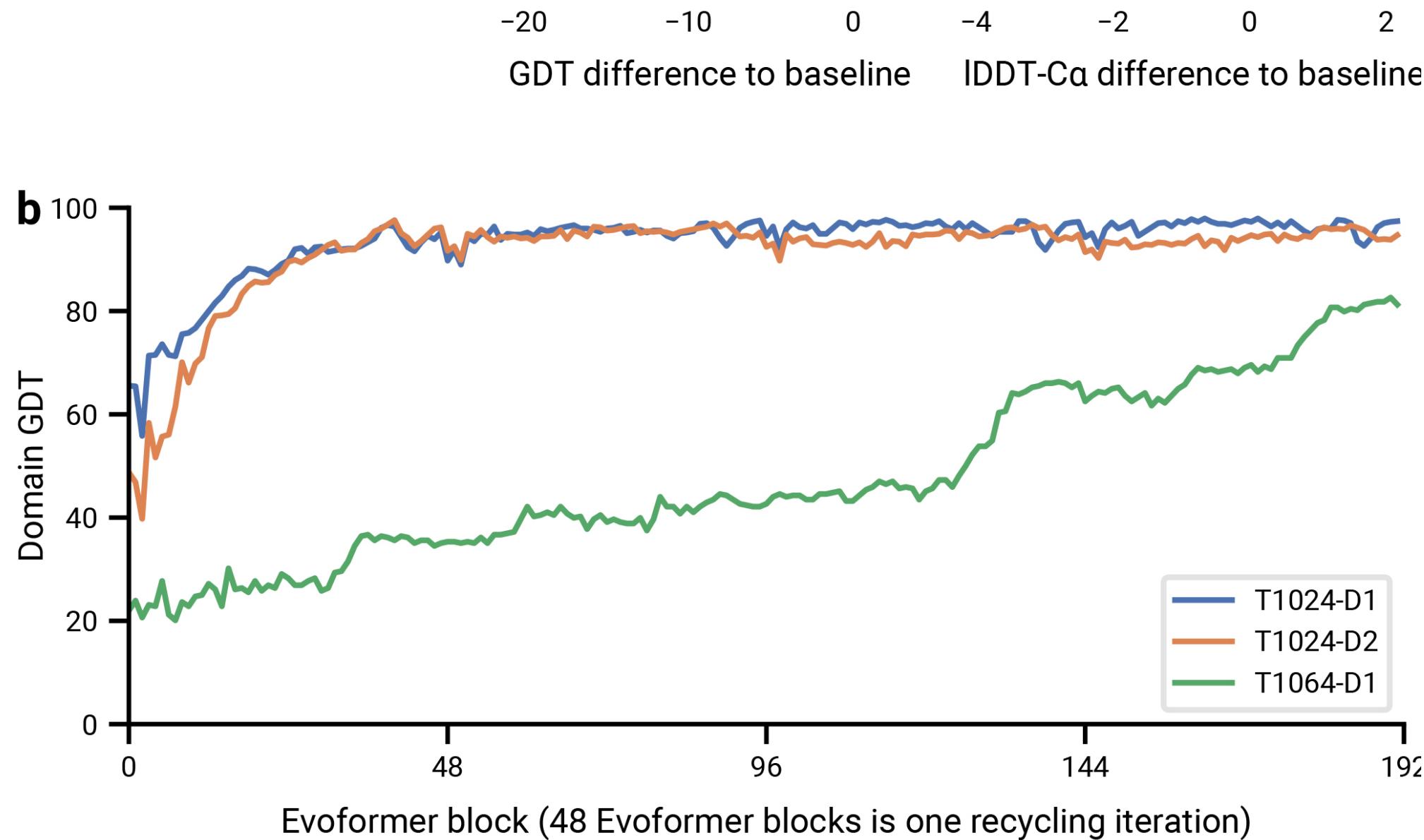
---

**End-to-end structure prediction**

# 核心算法框架



# 核心算法框架



# 模型训练方法

---

**Training with labelled and Unlabelled  
data**

# 模型训练方法

---

**Training with labelled and Unlabelled  
data**

# 模型训练方法

---

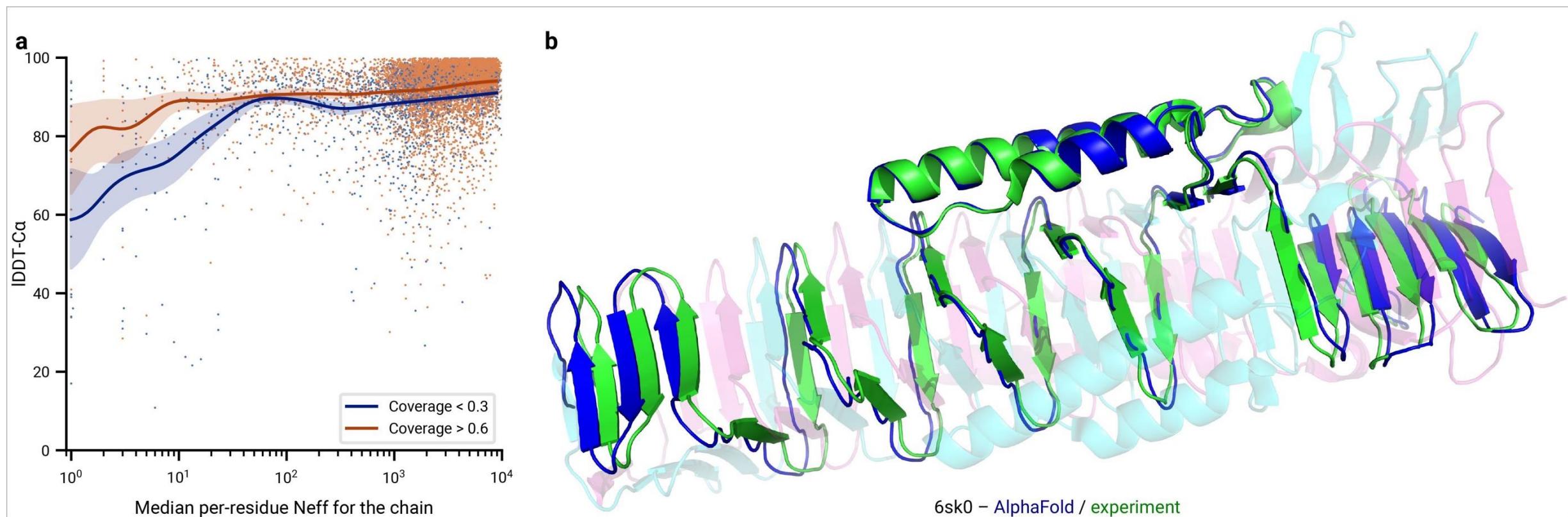
**Training with labelled and Unlabelled  
data**

# MSA 处理方法

---

**MSA depth and cross-chain contacts**

# MSA 处理方法



# MSA 处理方法

---

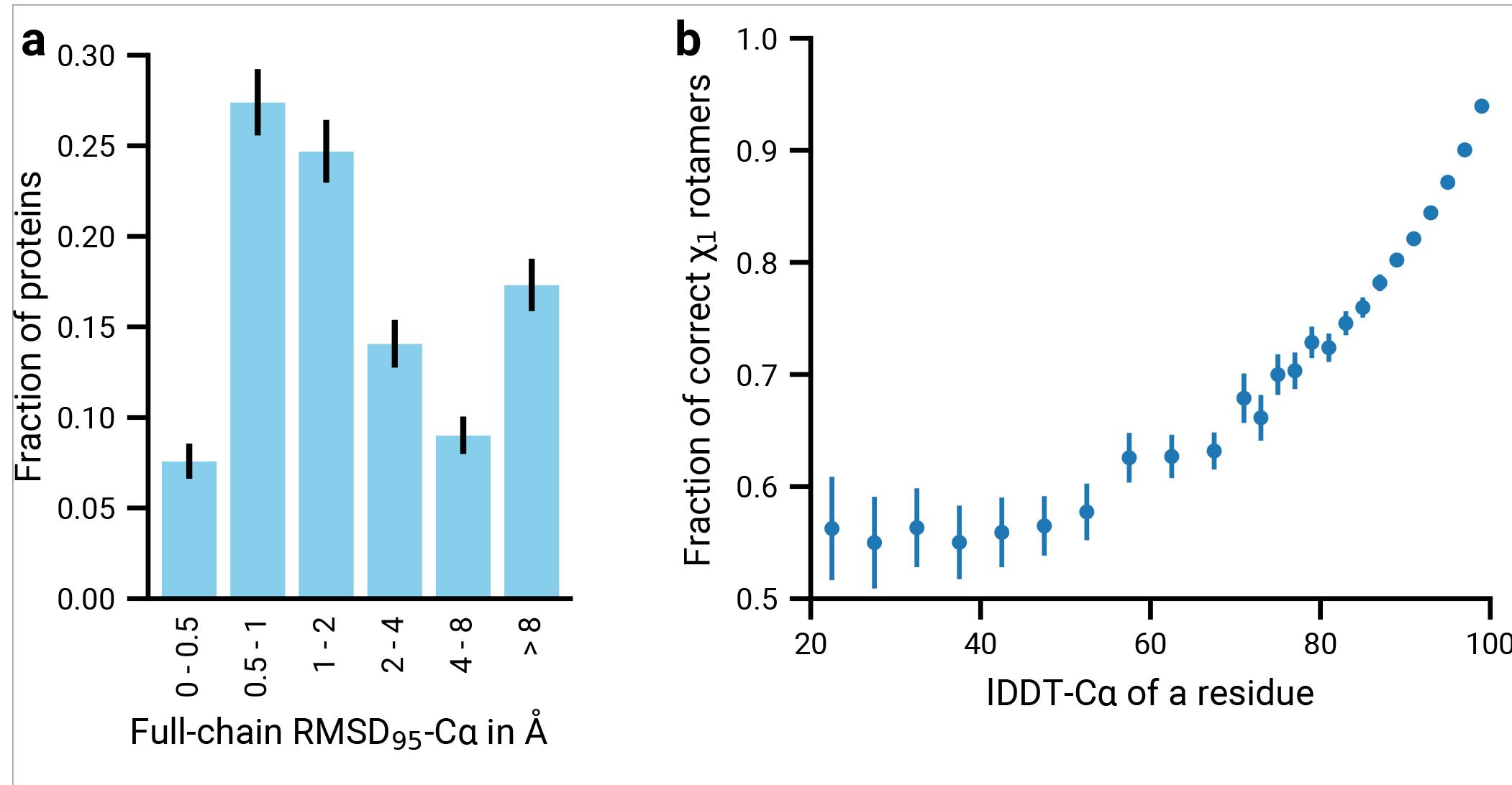
**MSA depth and cross-chain contacts**

# 模型效果

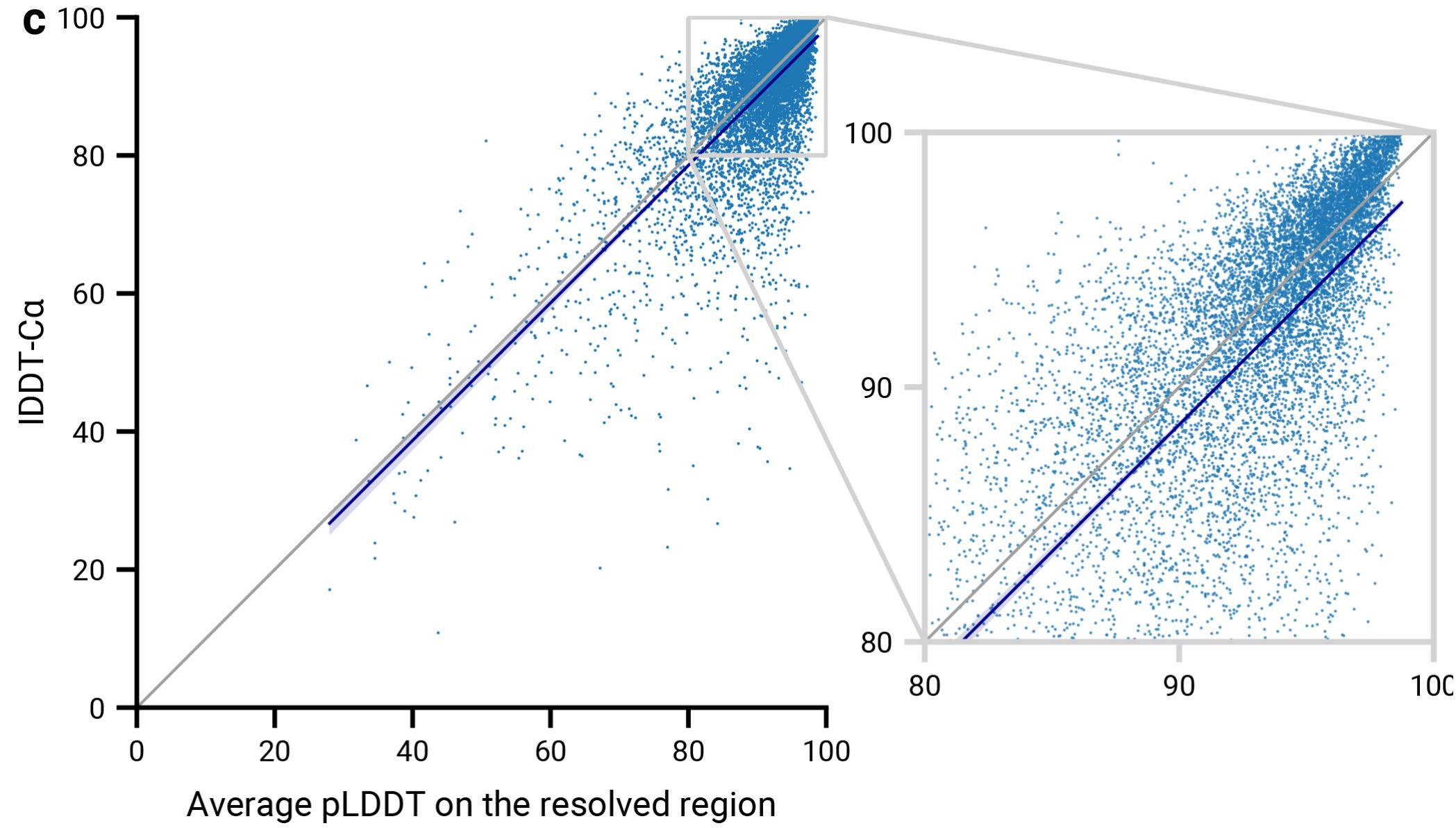
---

# Model Result

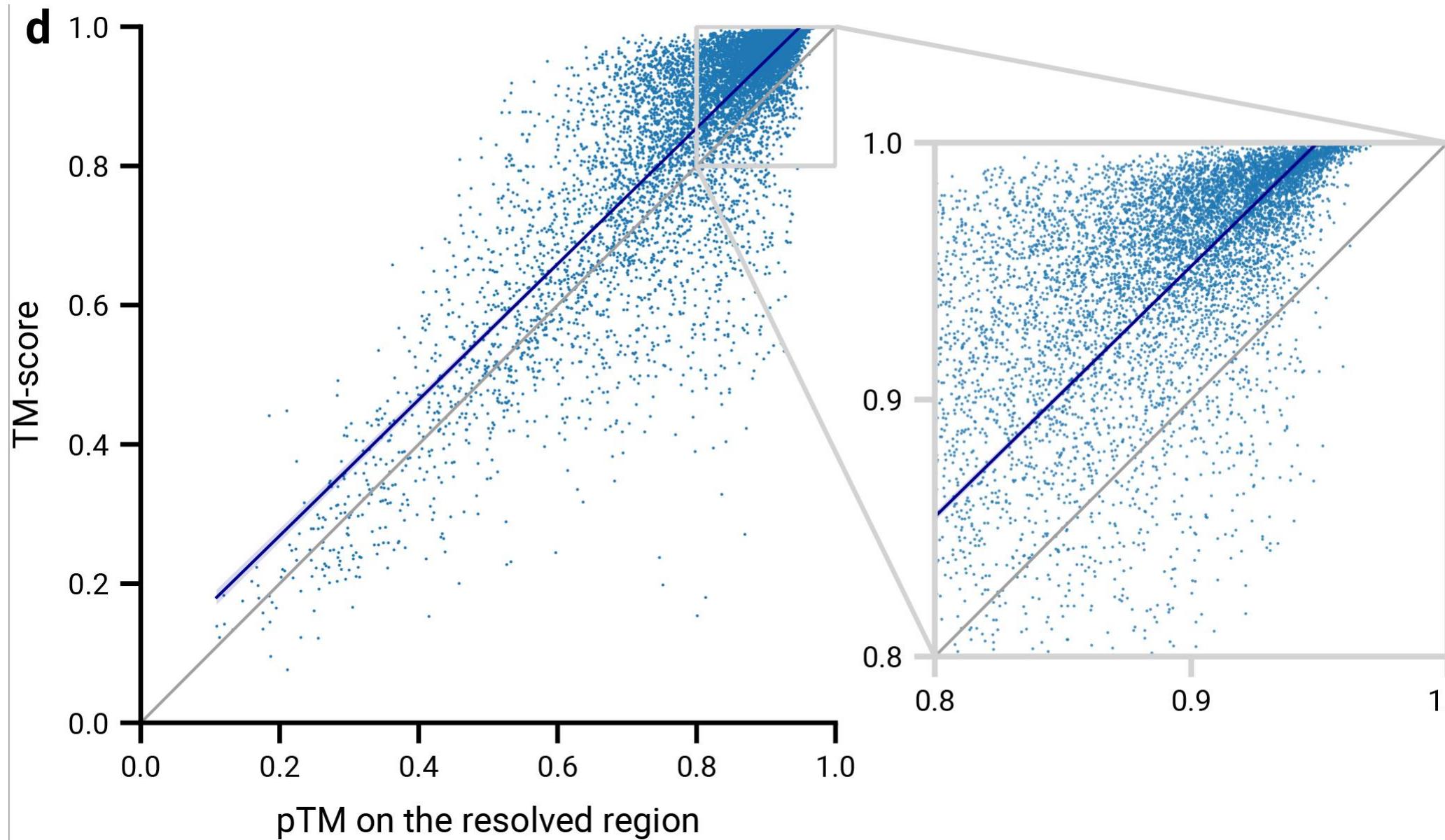
# 模型效果



# 模型效果



# 模型效果



# 模型效果

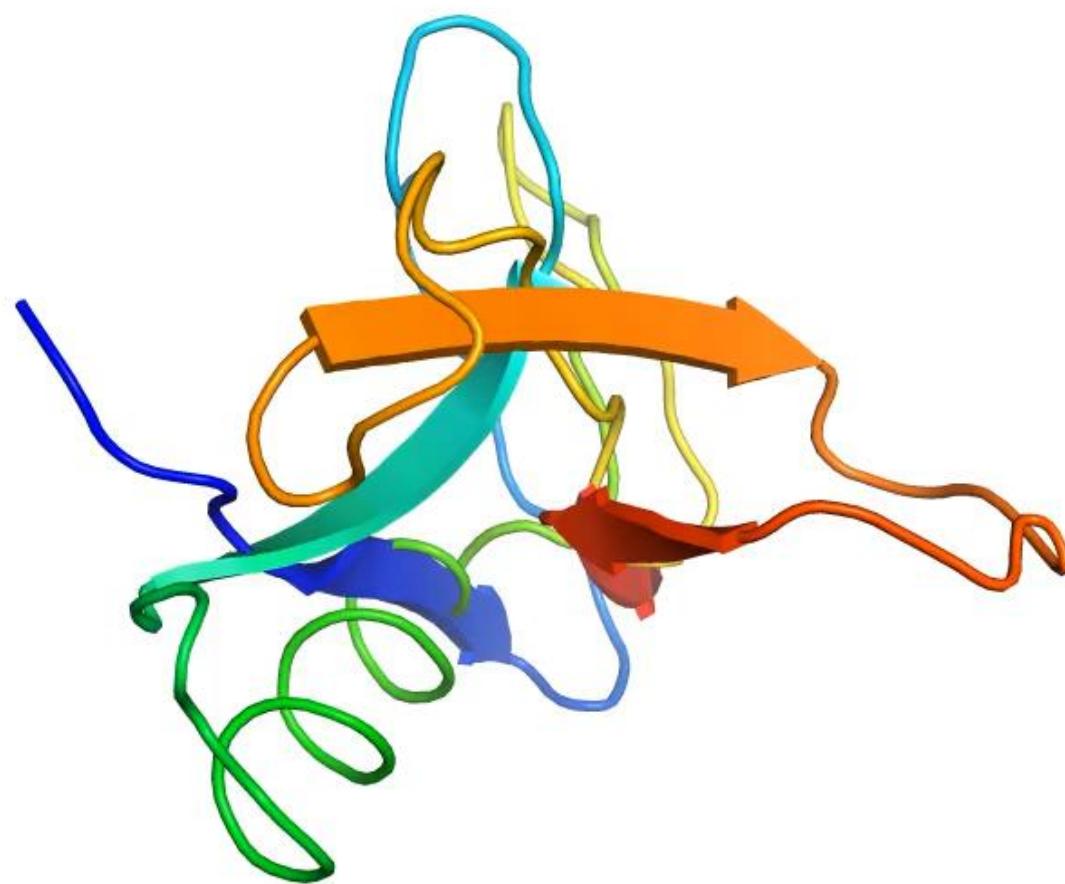
---



Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction

# 模型效果

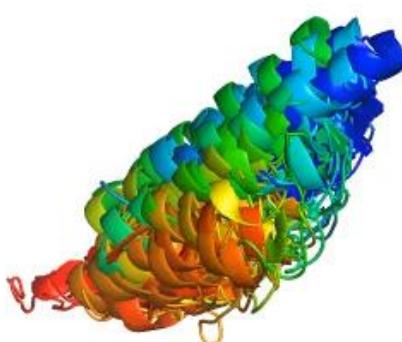
---



Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction

# 模型效果

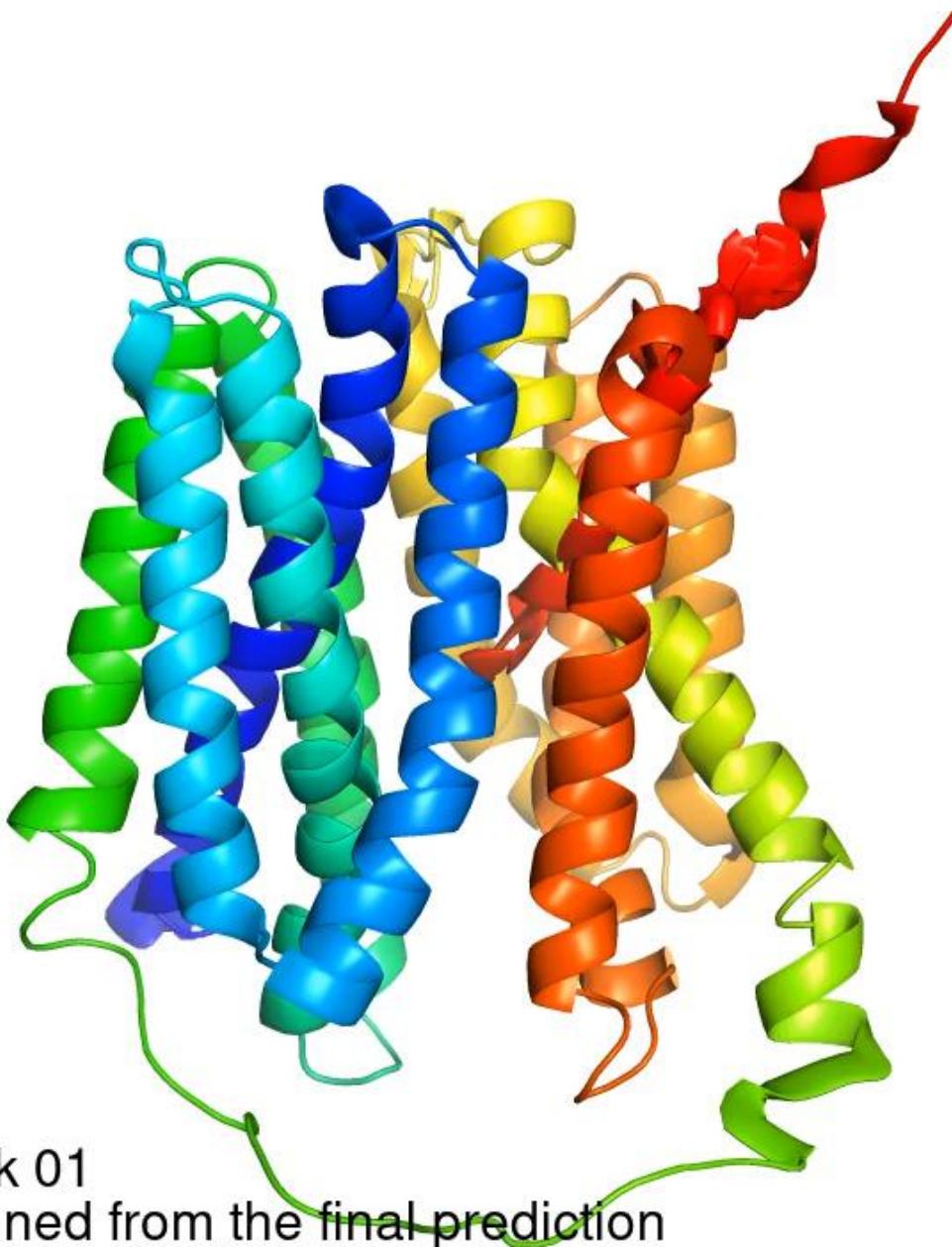
---



Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction

# 模型效果

---



# 相关工作

---

## Related work

## 相关工作

---

<https://github.com/chenxingqiang/ref-Alphafold-Code>

讨论

---

# Discussion

# 讨论

---

The methodology we have taken in designing AlphaFold is a combination of the bioinformatic and physical approaches: we use a physical and geometric inductive bias to build components that learn from PDB data with minimal imposition of handcrafted features (e.g. AlphaFold builds hydrogen bonds effectively without a hydrogen bond score function). This results in a network that learns far more efficiently from the limited data in the PDB but is able to cope with the complexity and variety of structural data.

# 讨论

---

In particular, AlphaFold is able to handle missing physical context and produce accurate models in challenging cases like intertwined homomers or proteins that only fold in the presence of an unknown heme group. The ability to handle underspecified structural conditions is essential to learning from PDB structures as the PDB represents the full range of conditions in which structures have been solved. In general, AlphaFold is trained to produce the protein structure most likely to appear as part of a PDB structure. In cases where a particular stoichiometry or ligand/ion is predictable from the sequence alone, AlphaFold is likely to produce a structure that respects those constraints implicitly

# 讨论

---

AlphaFold has already demonstrated its utility to the experimental community, both for molecular replacement , and interpreting cryogenic electron microscopy (cryo-EM) maps. Moreover, because AlphaFold outputs protein coordinates directly, AlphaFold produces predictions in graphics processing unit (GPU)-minutes to GPU-hours depending on the length of the protein sequence (e.g. around one GPU-minute per model for 384 residues, see Methods for details). This opens up the exciting possibility of predicting structures at the proteome-scale and beyond.

# 讨论

---

The explosion in available genomic sequencing techniques and data has revolutionized bioinformatics but the intrinsic challenge of experimental structure determination has prevented a similar expansion in our structural knowledge. By developing an accurate protein structure prediction algorithm, coupled with existing large and well-curated structure and sequence databases assembled by the experimental community, we hope to accelerate the advancement of structural bioinformatics that can keep pace with the genomics revolution. We hope that AlphaFold, and computational approaches that apply its techniques for other biophysical problems, will become essential tools of modern biology.

# 文献/资源列表

---

<https://github.com/chenxingqiang/alphafold2-codecs>