

UCLA

Samueli  
Computer Science

ScAI Machine Learning Group

# AlphaFold: AI Solution for Decade-long Protein Folding Challenge in Biology

Junheng Hao

University of California, Los Angeles (UCLA)

Tuesday, 10/19/2021



# Junheng Hao

PhD Candidate, University of California Los Angeles (2017-)

Advisor: Wei Wang, Yizhou Sun

Website: [Jeff's Home \(haojunheng.com\)](http://haojunheng.com)

## Bio

- 5th-year Ph.D. candidate at UCLA co-advised by Yizhou Sun and Wei Wang in UCLA Data Mining Group.
- My research interests include knowledge graph, graph representation learning, KG-empowered applications (NLP, Bioinformatics, recommender systems, etc.).

## Past Experiences

- Research Intern, Microsoft Research Redmond, 2021
- PhD Research Intern, IBM, 2020
- Applied Science Intern, Amazon Product Graph, 2019
- Research Intern, NEC Labs America, 2018

# Today's Agenda

- Background: Protein Structure Prediction
- AlphaFold v1: CNN
- AlphaFold v2: Transformer/Attention
- Science: Three-stack NN & SE(3)-Transformers
- Discussion: Network Science and Graph in Biology World

# Papers

- **AlphaFold:** Improved protein structure prediction using potentials from deep learning (Published on Nature, Jan 2020)
- **AlphaFold2:** Highly accurate protein structure prediction with AlphaFold (Published on Nature, July 2021)
- (Optional Reading) Accurate prediction of protein structures and interactions using a three-track neural network (Published on Science, July 2021)

# Background: What is protein folding and why is it important?

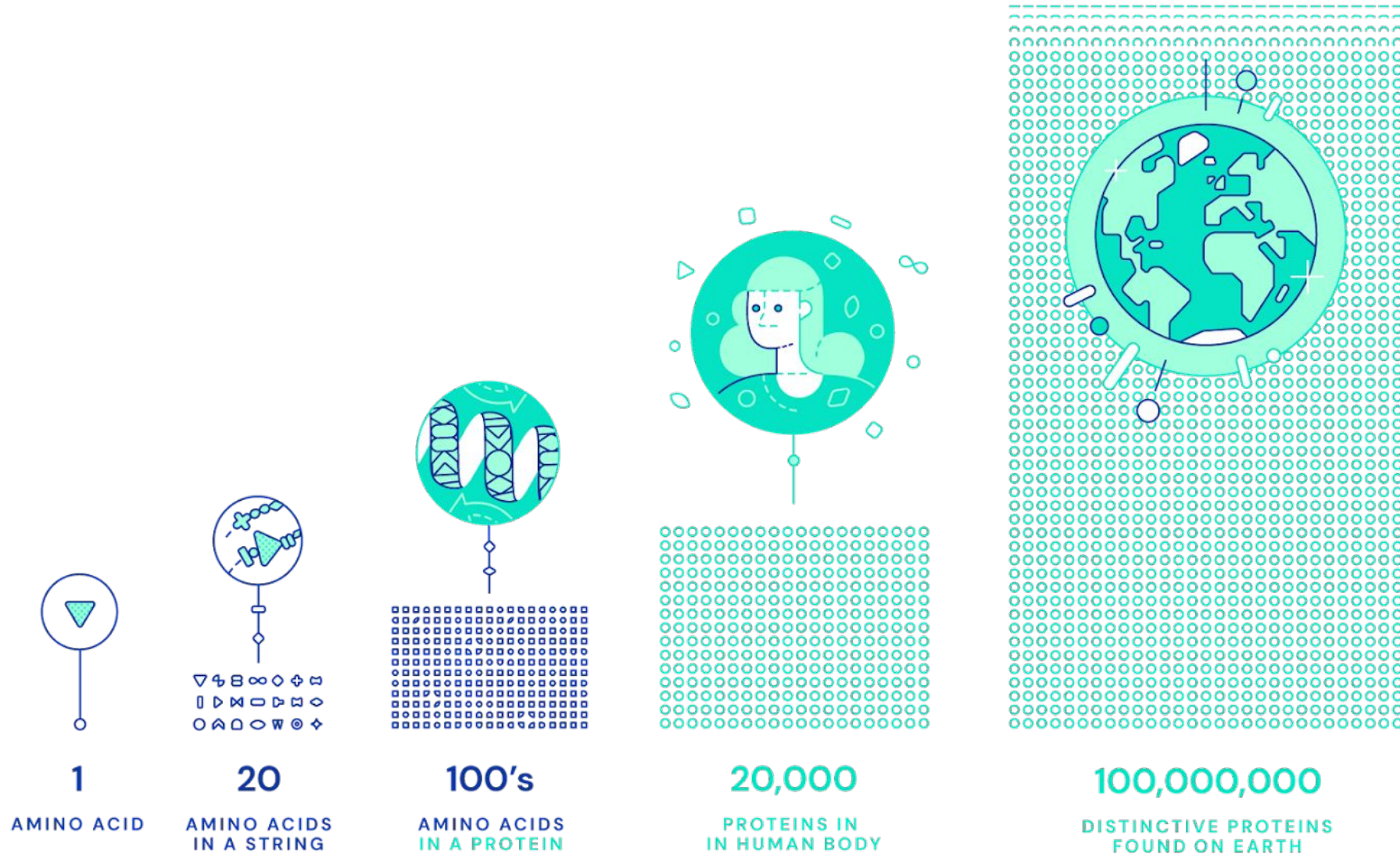
*A decade-long biology challenge for proteins, the building blocks of life in the planet.*

# Biology 101: Proteins

---

- Proteins are large, complex molecules essential to all of life. Nearly every function that our body performs (e.g. contracting muscles, sensing light, or turning food into energy), relies on proteins, and how they move and change.
- What any given protein can do (largely) depends on its unique 3D structure. Examples are:
  - Notorious “*spike proteins*” which stud coronavirus that allows the virus to enter our cells.
  - Antibody proteins utilized by our immune systems are *Y-shaped*, and form unique hooks.
  - Collagen proteins are shaped like cords, which transmit tension between cartilage, ligaments, bones, and skin.
- The recipes for those proteins, called genes, are encoded in our DNA and generated by Ribosome. Many diseases and deaths for an organism, are fundamentally linked to malformed proteins.
- Proteins are composed of **chains of amino acids** (also referred to as amino acid **residues**). But DNA only contains information about the sequence of amino acids, not how they fold into shape.

# Biology 101: Proteins



# Why is Protein Folding Important?

---

*“I think that we shall be able to get a more thorough understanding of the nature of disease in general by investigating the molecules that make up the human body, including the abnormal molecules, and that this understanding will permit...the problem of disease to be attacked in a more straightforward manner such that new methods of therapy will be developed.”*

-- Linus Pauling, 1960



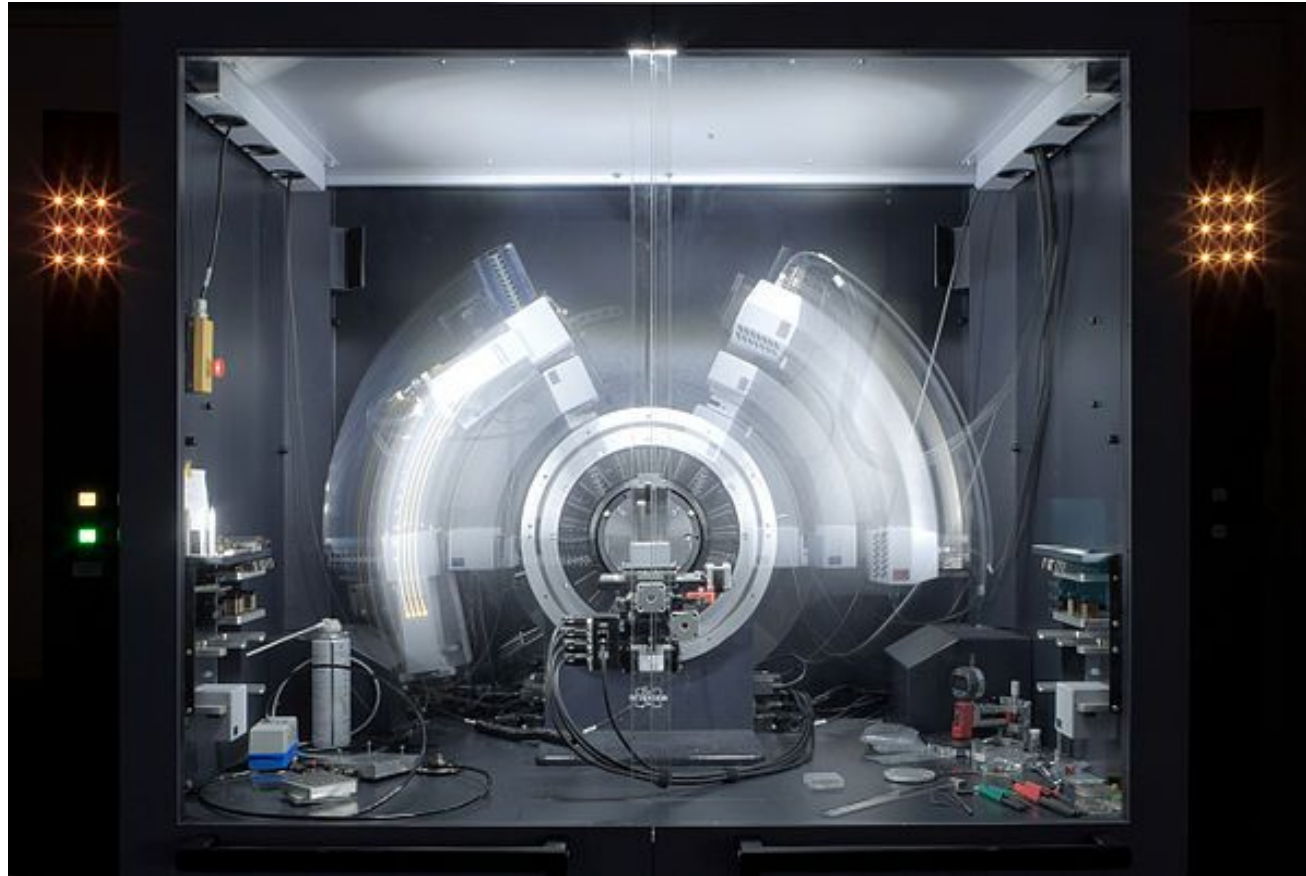
# Why is Protein Folding Important?

---

- Scientists have long been interested in determining the structures of proteins because a protein's form is thought to dictate its function.
- Once a protein's shape is understood, its role within the cell can be guessed at, and scientists can develop drugs that work with the protein's unique shape.
- Traditional methods: Experimental techniques like [cryo-electron microscopy](#), [nuclear magnetic resonance](#) and [X-ray crystallography](#)
  - *A lot of trial and error, time consuming, high cost*
  - *Tens or hundreds of thousands of dollars per protein*
- Motivation: Biologists are turning to AI methods as an alternative to this long and laborious process for difficult proteins.
- The ability to predict a protein's shape computationally from its genetic code alone could no doubt help accelerate research.

# X-ray crystallography

- Huge cost: Hundreds of thousands of dollars and about one year in duration for one protein → Only 170,000 protein folding structures have been identified

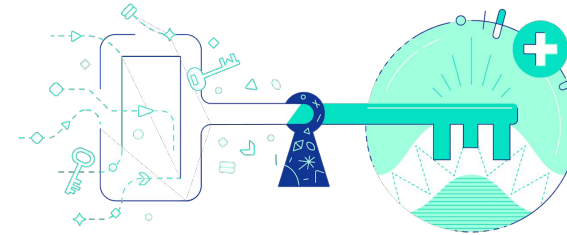


# Protein Folding: Take-away

**Uniqueness:** The sequence usually map 1-to-1 to a 3D structure.

**Problem:** Huge number ways and possibilities to fold.

**Cost:** X-ray crystallography costs \$120,000 and takes 1 year.



**Function:** 3D structure determines its function. Misfold → disease

**Dataset:** 200M proteins with sequences but only 170K with available 3D structures.

# Protein Folding: Promising Applications

## *Near-term*

**DNA → Function:** Learn unknown function of genes encoded in DNA

**Disease:** Understand the cause of disease as results of misfolded proteins.

**Treatment:** Design proteins to fix other misfolded proteins.

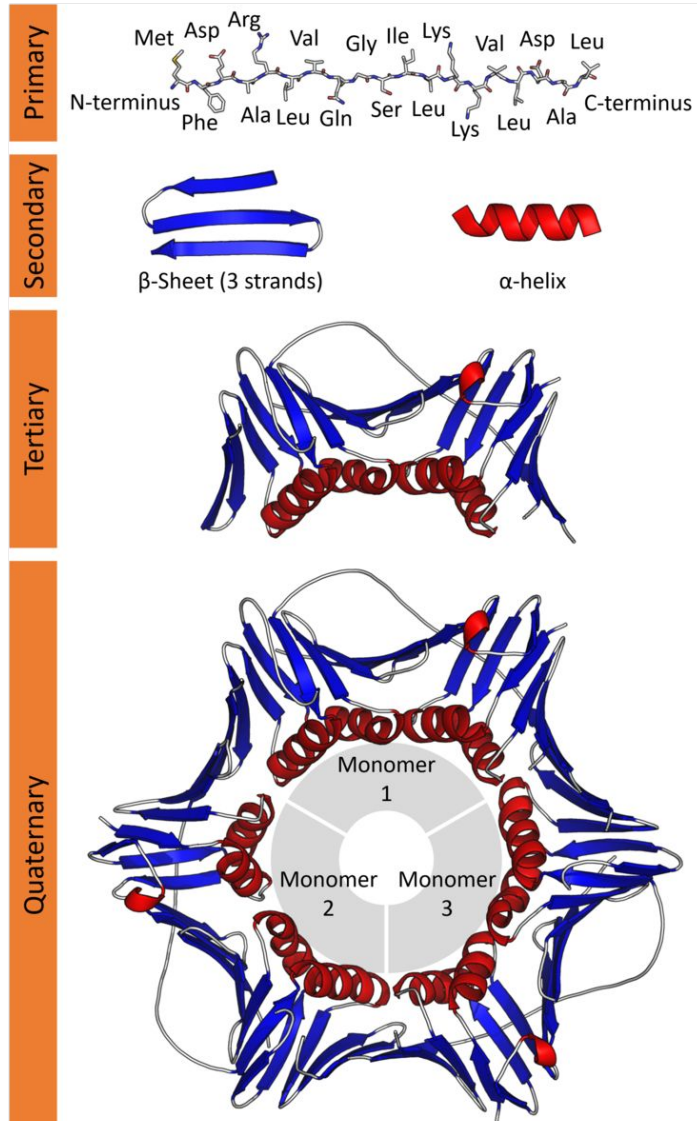
**Other applications:** Agriculture, Supplements and biomaterials.

## *Long-term*

**Physics-based stimulation of biological systems**

**Biological and artificial life**

# Four Levels of Protein Structures



← Level 1: What we mostly (and easily) know about!

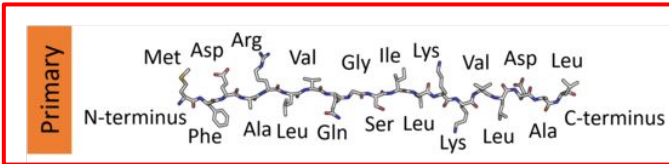
← Level 3: What we mostly care about! The Folding!

Credit:

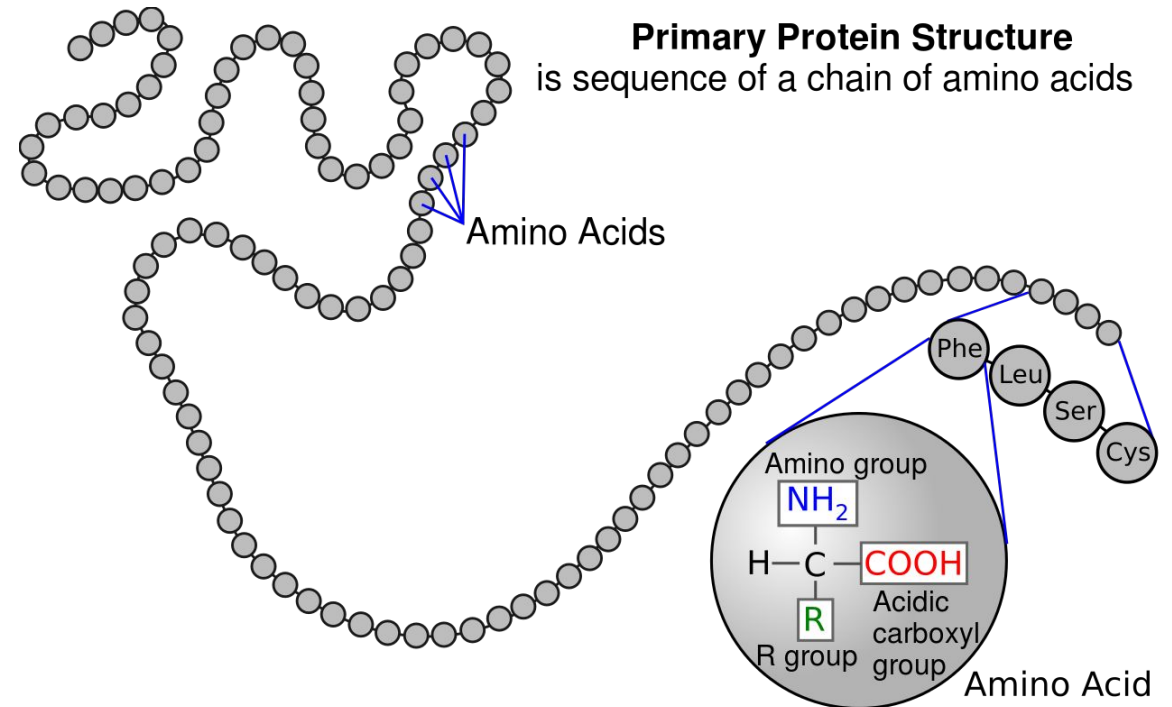
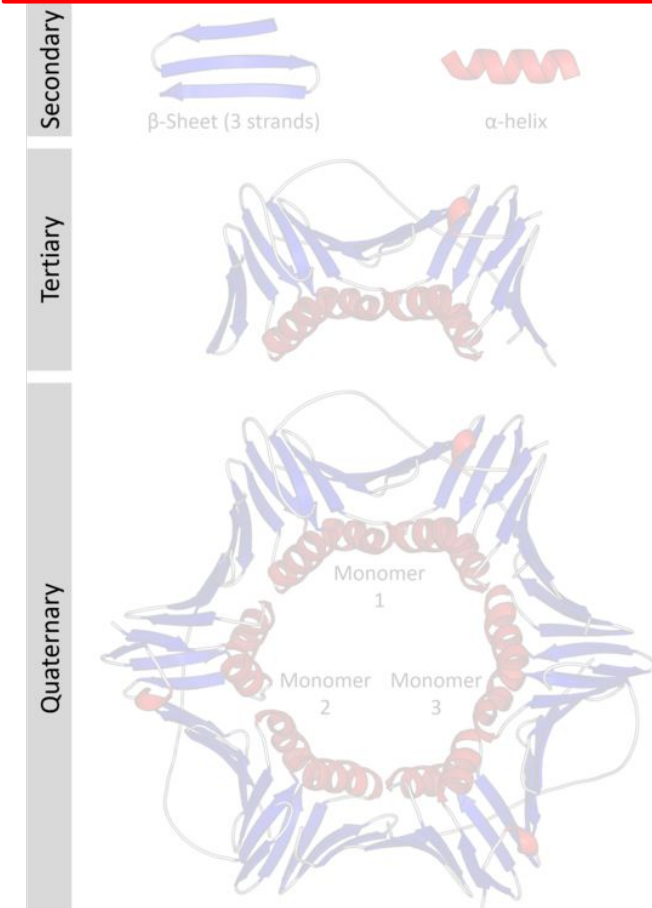
[1]<https://www.khanacademy.org/science/biology/macromolecules/proteins-and-amino-acids/a/orders-of-protein-structure>

[2][https://en.wikipedia.org/wiki/Protein\\_structure](https://en.wikipedia.org/wiki/Protein_structure)

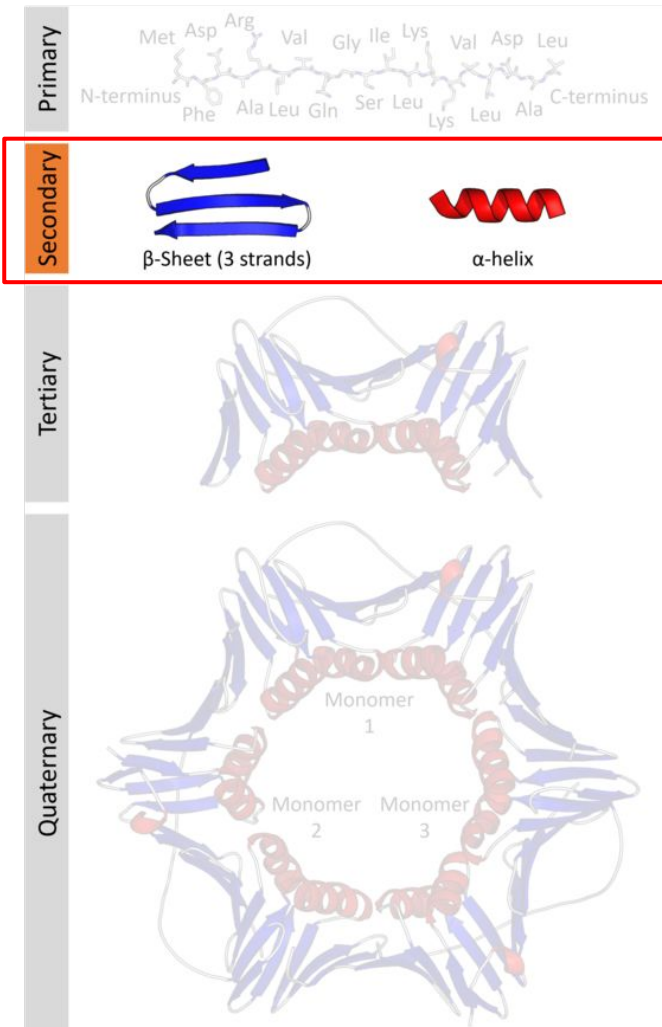
# Four Levels of Protein Structures: **Primary**



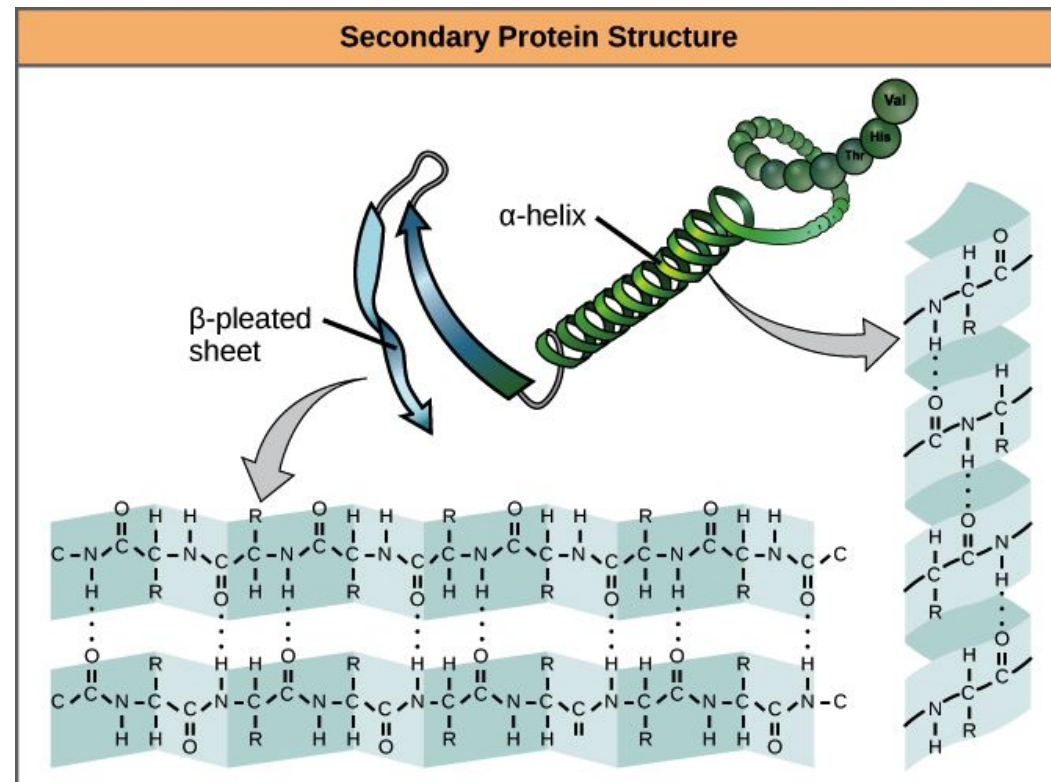
- A sequence of amino acids with the alphabet = “ARNDCQEGHILKMFPSTWYV”)
- Connected by **Peptide Bond** -CO-NH-



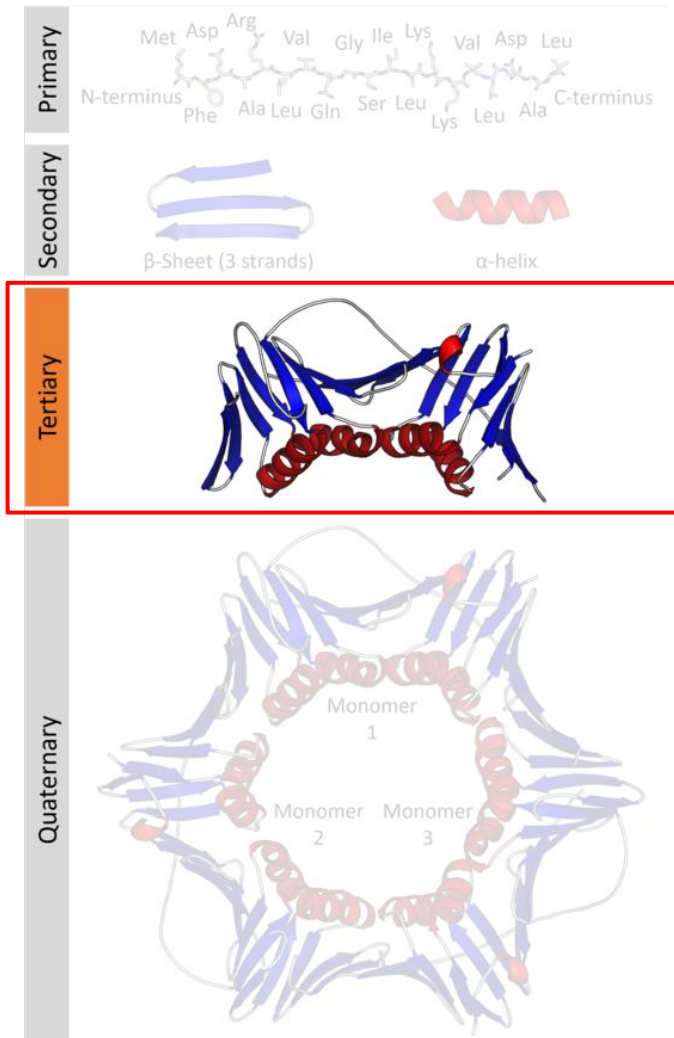
# Four Levels of Protein Structures: **Secondary**



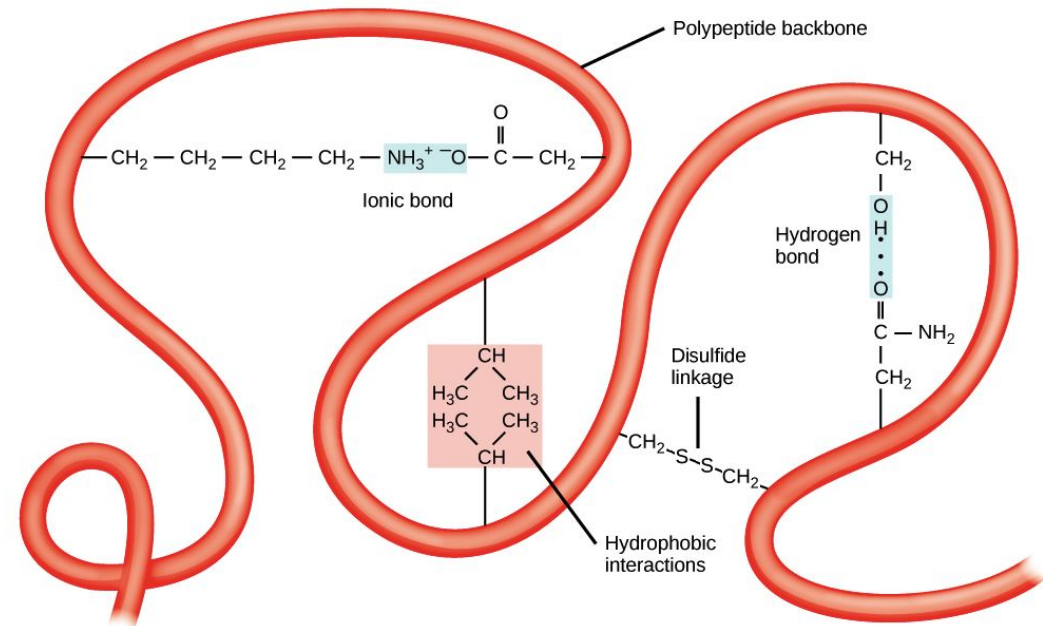
- Typical substructures: helices and sheets
- By Hydrogen Bonds



# Four Levels of Protein Structures: Tertiary

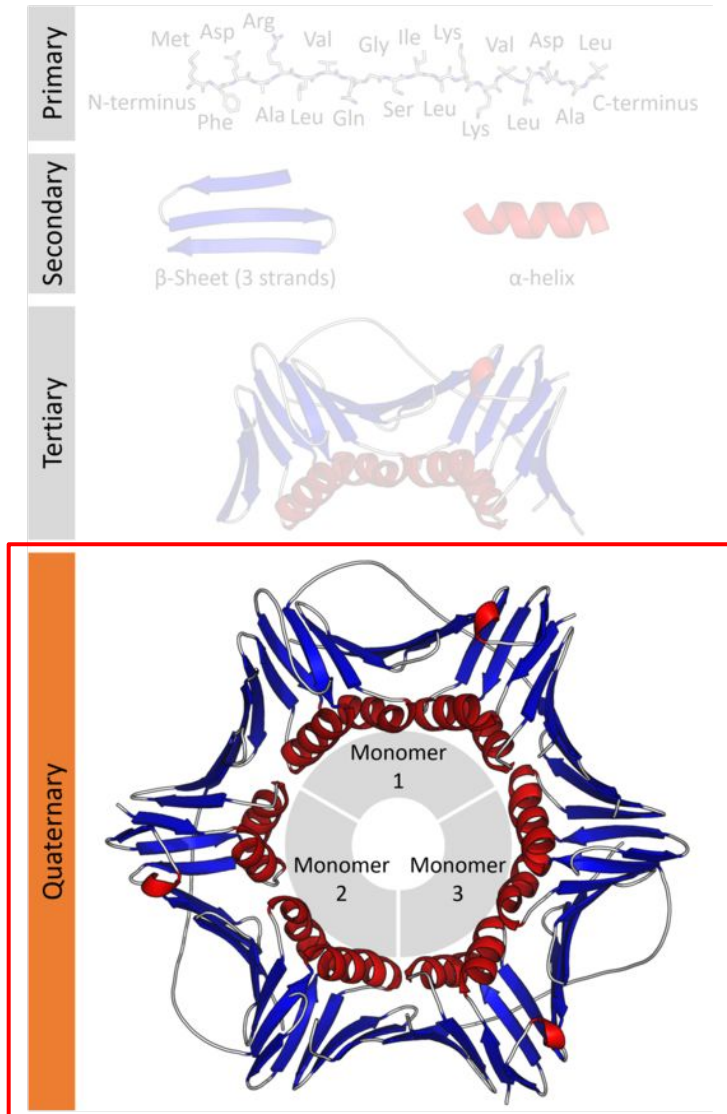


- The overall three-dimensional structure of a polypeptide.
- Typically require deep knowledge about stereochemistry and more advanced expertise.
- **This is the level of prediction where AlphaFold (and AlphaFold 2) focus.**



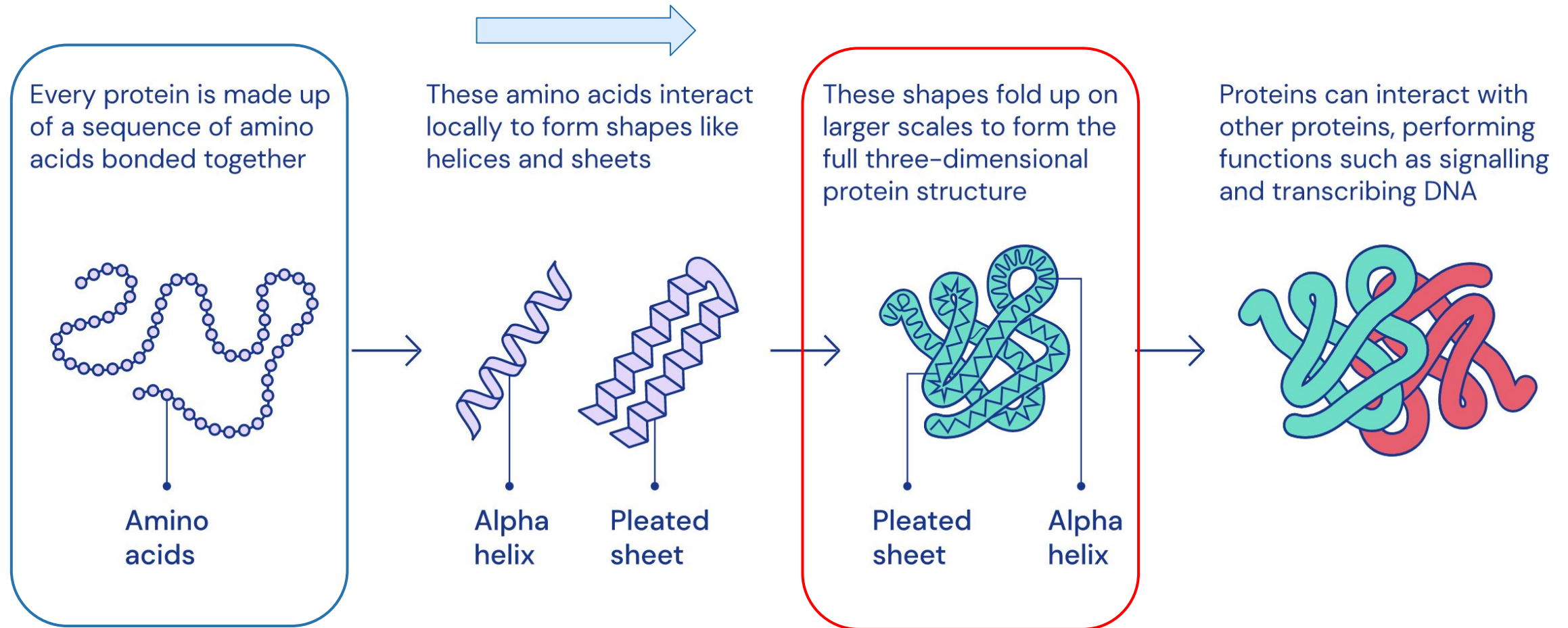


# Four Levels of Protein Structures: Quaternary



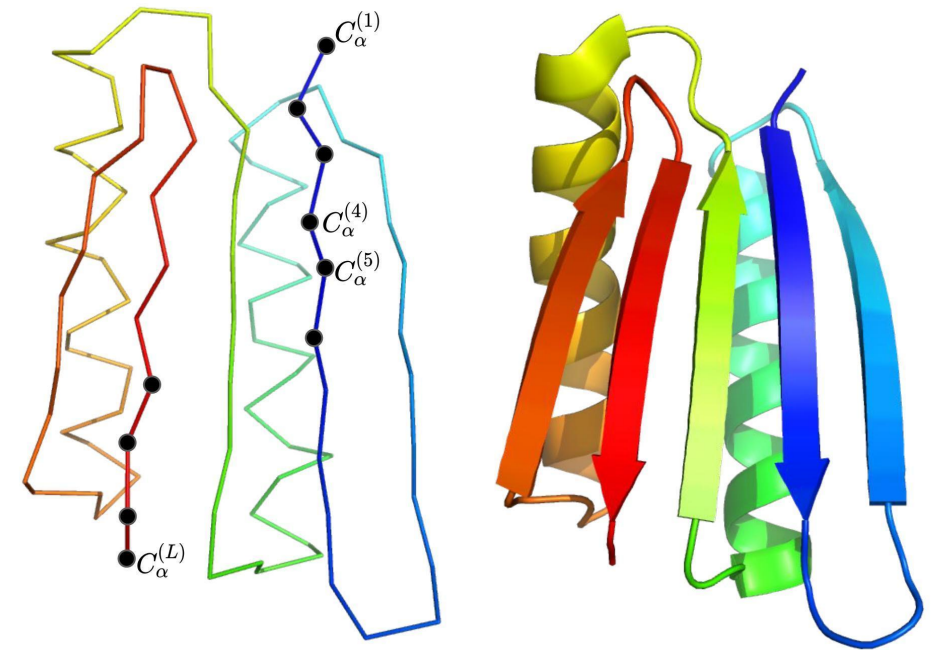
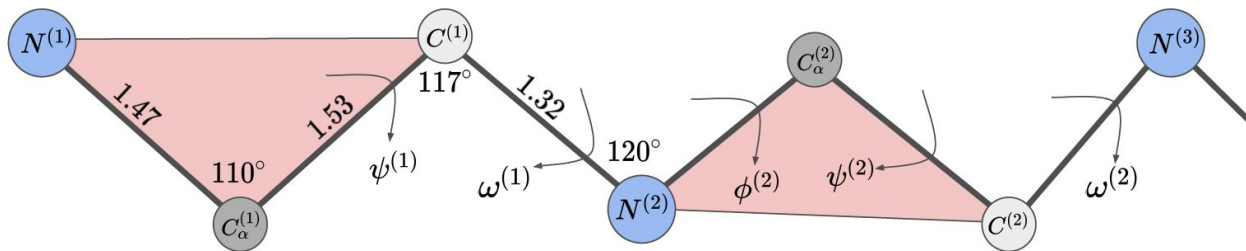
- Many proteins are made up of a single polypeptide chain and have only three levels of structure (the ones we've just discussed).
- However, some proteins are made up of multiple polypeptide chains, also known as subunits. When these subunits come together, they give the protein its **quaternary** structure.

# Protein Structures: High-level summary



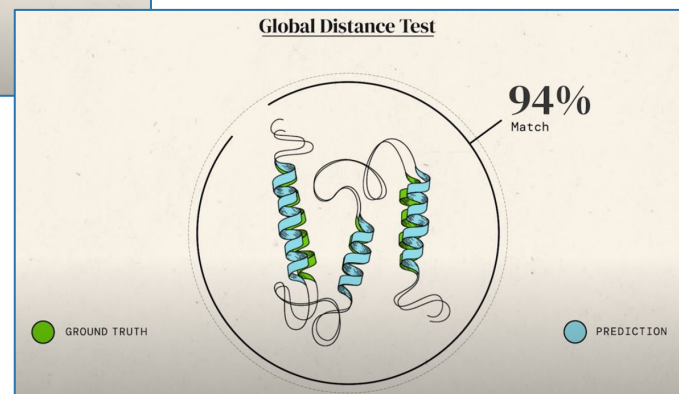
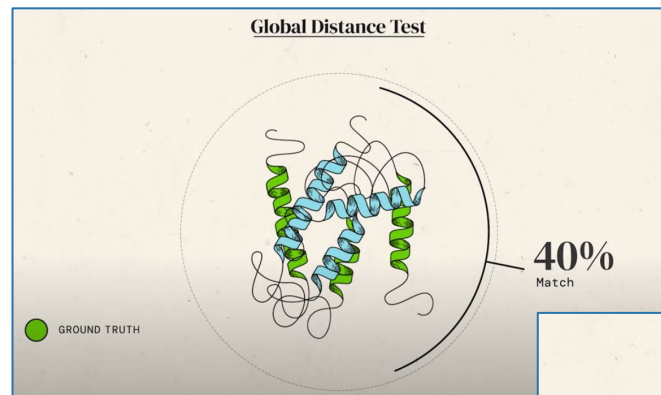
# Protein Backbone Geometry

A protein backbone is a repeating sequence (linear chain) of 3 atoms: nitrogen, carbon, and another carbon, namely  $\underbrace{N^{(1)}, C_{\alpha}^{(1)}, C^{(1)}}_{}, \underbrace{N^{(2)}, C_{\alpha}^{(2)}, C^{(2)}}_{}, \dots, \underbrace{N^{(L)}, C_{\alpha}^{(L)}, C^{(L)}}_{}$



# CASP: “Kaggle” on Protein Structure Prediction

- Critical Assessment of protein Structure Prediction [\[Main Page\]](#)
  - Known as “Protein Structure Prediction Center”
- Evaluation
  - Global Distance Test (GDT)
  - TM-Score, RMSD



**CASP14**

14th Community Wide Experiment on the Critical Assessment of Techniques for Protein Structure Prediction

Menu

- [Home](#)
- [PC Login](#)
- [PC Registration](#)
- ▼ [CASP Experiments](#)
  - [CASP14 \(2020\)](#)
    - [CASP Commons \(COVID-19, 2020\)](#)
    - [CASP13 \(2018\)](#)
    - [CASP12 \(2016\)](#)
    - [CASP11 \(2014\)](#)
    - [CASP10 \(2012\)](#)
    - [CASP9 \(2010\)](#)
    - [CASP8 \(2008\)](#)
    - [CASP7 \(2006\)](#)
    - [CASP6 \(2004\)](#)
    - [CASP5 \(2002\)](#)
    - [CASP4 \(2000\)](#)
    - [CASP3 \(1998\)](#)
    - [CASP2 \(1996\)](#)
    - [CASP1 \(1994\)](#)
  - [Initiatives](#)
  - [Data Archive](#)
  - [Proceedings](#)
  - [CASP Measures](#)
  - [Feedback](#)
  - [Assessors](#)
  - [People](#)
  - [Community Resources](#)
  - [Job Fair](#)

**CASP14**

CASP provides an independent mechanism for the assessment of methods of protein structure modeling. From May through August 2020, CASP organizers have been posting on this website sequences of unknown protein structures for modeling. Protein models have been collected from May through mid-September, and evaluated as the experimental coordinates become available. In the summer and fall, the tens of thousands of models submitted by approximately 100 research groups worldwide are processed and evaluated. Independent assessors in each of the prediction categories bring independent insight into their assessment. Tools for viewing, comparison, and analysis of submitted models are available from this website.

Targets	Predictors	Conference	Results	CASP14 in news
<a href="#">Target List</a>	<a href="#">Groups Info</a>	<a href="#">Abstracts</a>	<b>AUTOMATIC EVALUATION</b>	<a href="#">CASP Press Release</a>
<a href="#">Domain Definition</a>		<a href="#">Program</a>	CASP14 results will be published in a special edition of Proteins in 2021.	<a href="#">Nature</a>
		<a href="#">Presentations</a>		<a href="#">Science</a>
		<a href="#">Recordings</a>		<a href="#">New York Times</a>
		<a href="#">CASP14 Conference Platforms</a>		<a href="#">BBC news</a>
			<a href="#">Parseable Data</a>	<a href="#">Fortune</a>
			<a href="#">Rankings: Regular targets (T)</a>	<a href="#">CNBC news</a>
			<a href="#">Rankings: Multimeric targets (H,To)</a>	<a href="#">Bloomberg</a>
			<a href="#">Rankings: Inter-domain prediction</a>	<a href="#">Financial Post</a>
			<a href="#">Rankings: Refinement targets (R)</a>	<a href="#">MIT Technology Review</a>
			<a href="#">Rankings: Contact predictions</a>	<a href="#">The Guardian</a>
				<a href="#">The Telegraph</a>
				<a href="#">Daily Mail</a>
				<a href="#">Tech Crunch</a>
				<a href="#">Venture Beat</a>
				<a href="#">New Scientist</a>
				<a href="#">SciTech Daily</a>
				<a href="#">Eureka Alert</a>
				<a href="#">News Medical</a>
				<a href="#">MedCity News</a>

# Dataset: What do we know about proteins?

- **Sequence databases** → **200M+**
  - UniRefA (JackHMMER)
  - BFD (HHblits)
  - MGnify clusters (JackHMMER)
- **Structural databases** → **Around 170K**
  - PDB (training)
  - PDB70 clustering (hhsearch)

## References:

- [1] Berman et al., Nature Structural Biology (2003) doi:10.1038/nsb1203-980
- [2] Mitchell et al., Nucleic Acids Research (2019) doi:10.1093/nar/gkz1035
- [3] Potter et al., Nucleic Acids Research (2018) doi:10.1093/nar/gky448
- [4] Steinegger et al., BMC Bioinformatics (2019) doi:10.1186/s12859-019-3019-7
- [5] Steinegger et al., Nature Methods (2019) doi:10.1038/s41592-019-0437-4
- [6] Suzek et al., Bioinformatics (2015) doi:10.1093/bioinformatics/btu739

**Visualization:** PyMol <https://pymol.org/2/>

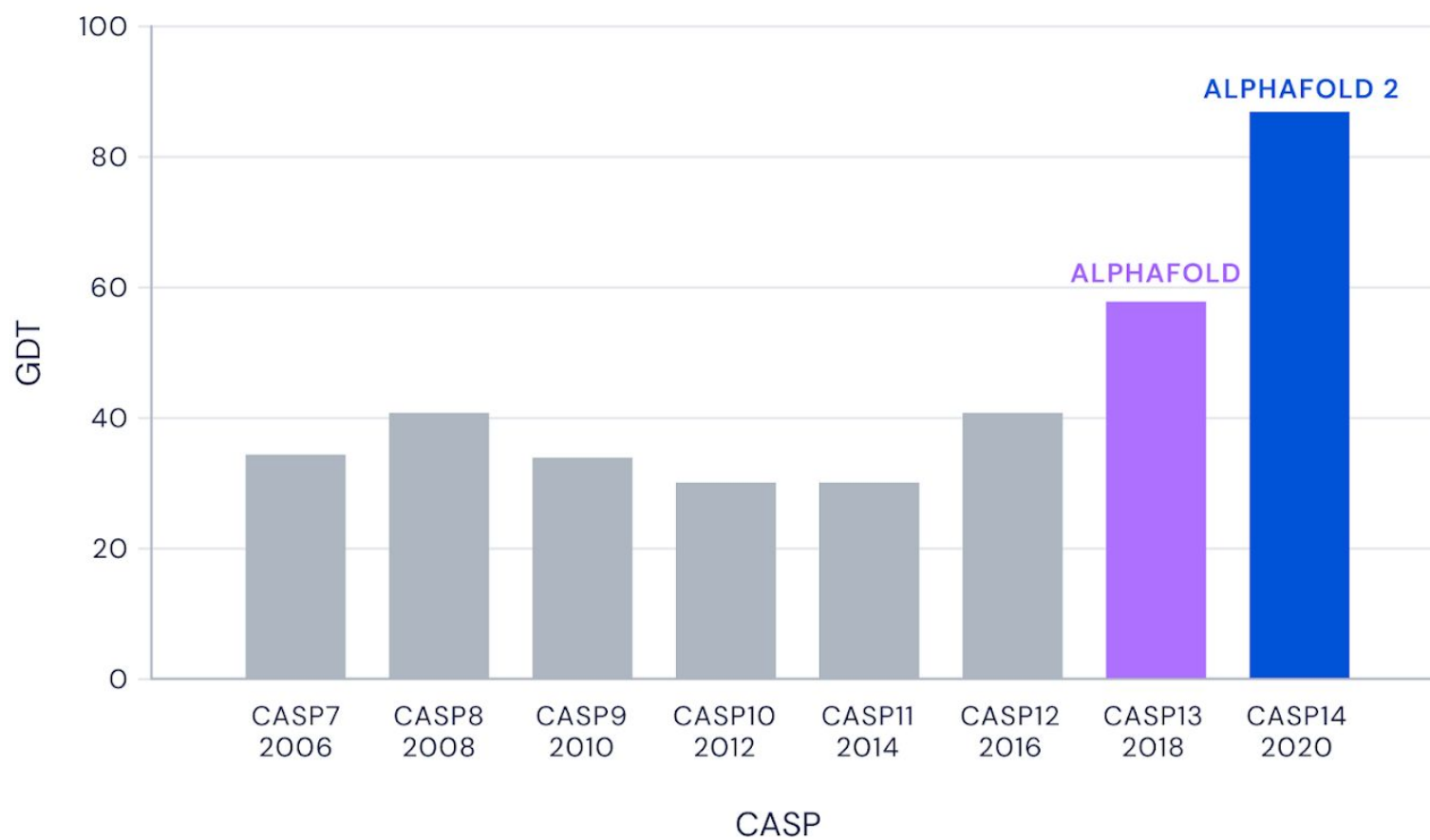


# AlphaFold v1: Improved protein structure prediction using potentials from deep learning

*One CNN-supported Protein Folding Model*

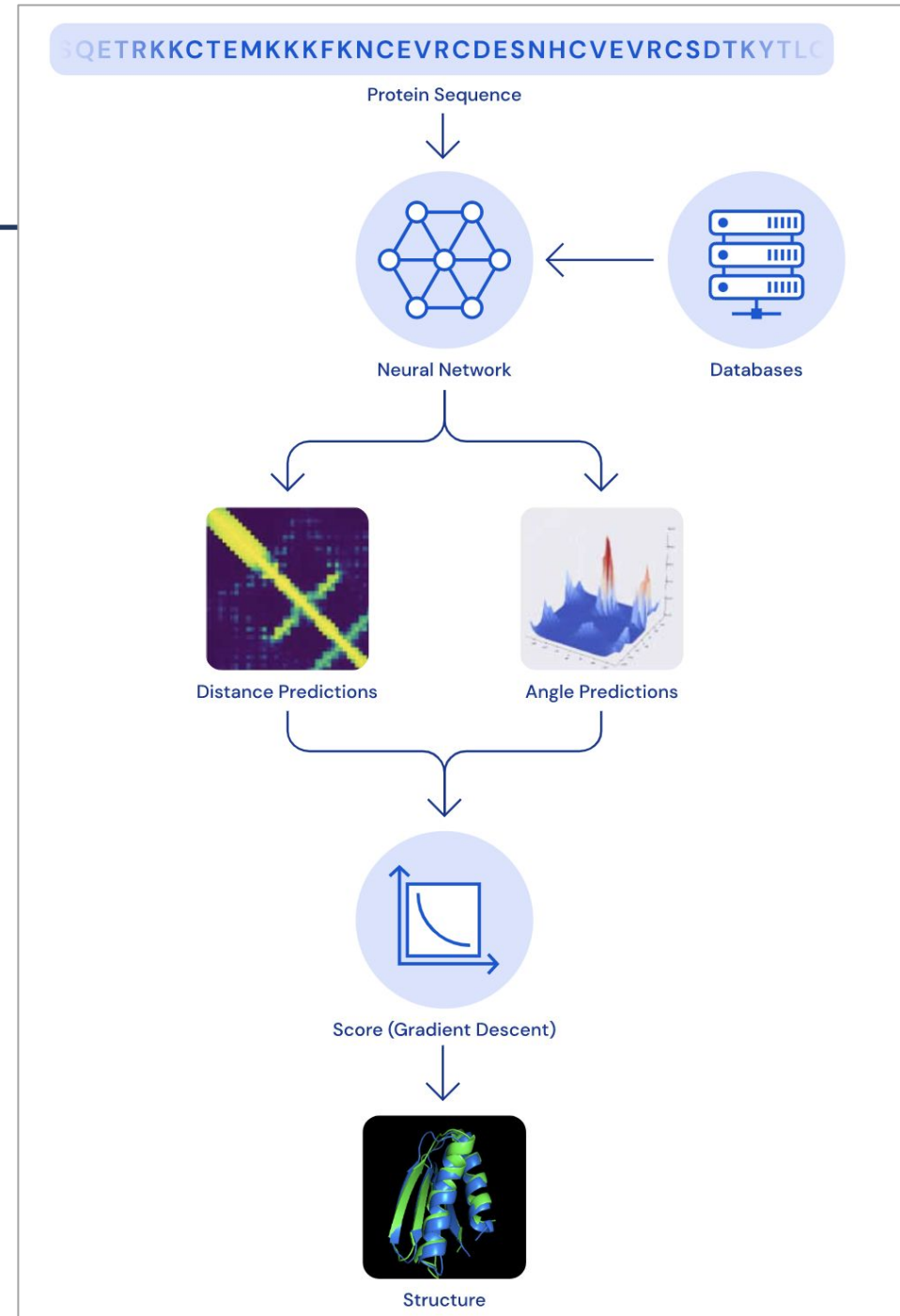
# AlphaFold and AlphaFold on CAPS14 Challenge

Median Free-Modelling Accuracy



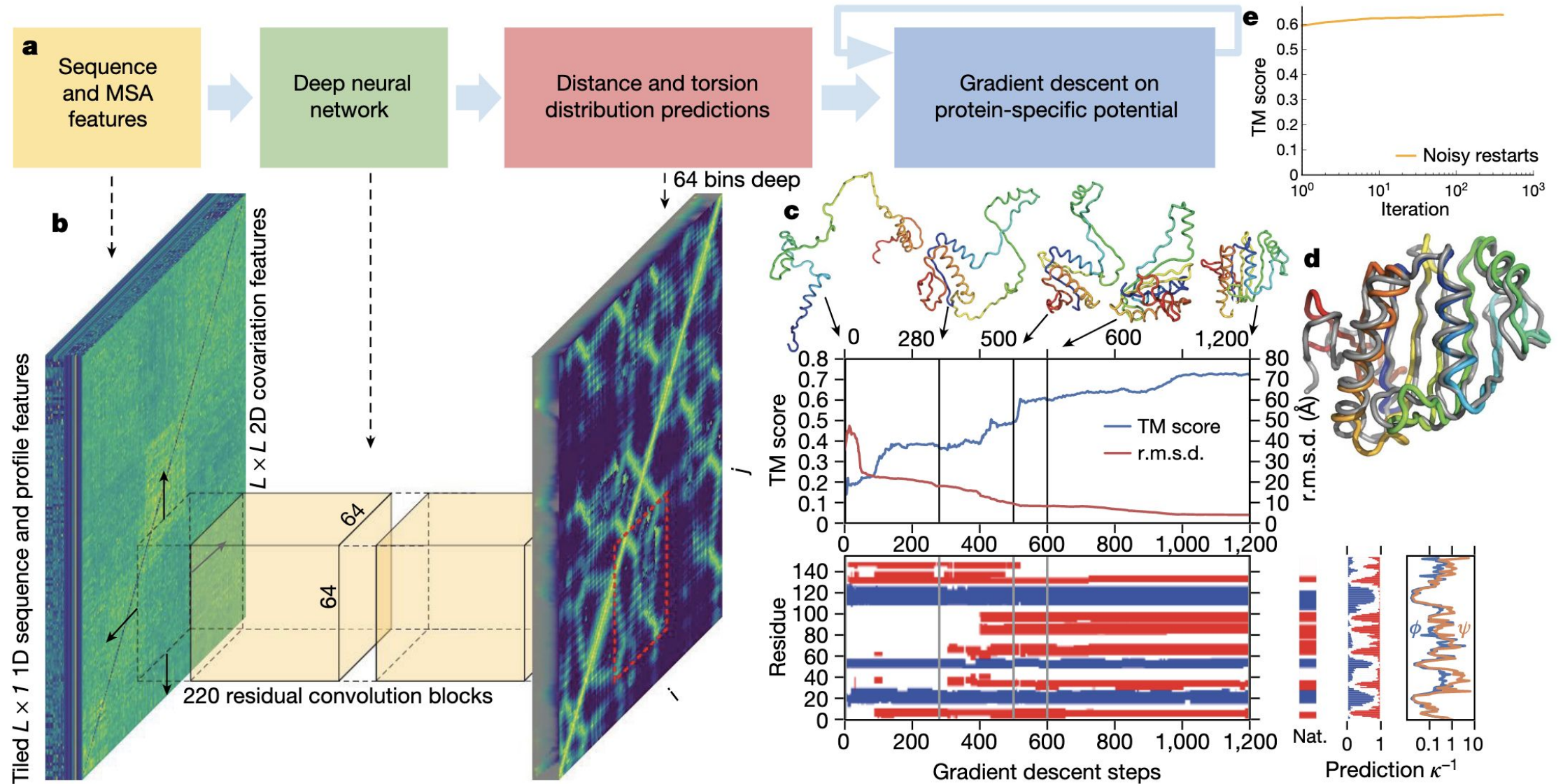
# AlphaFold v1: Schematic Architecture

- Residual CNN as core model to predict distance and angle to create final structure output
- Using Multiple Sequence Alignment (MSA) from databases for feature generation

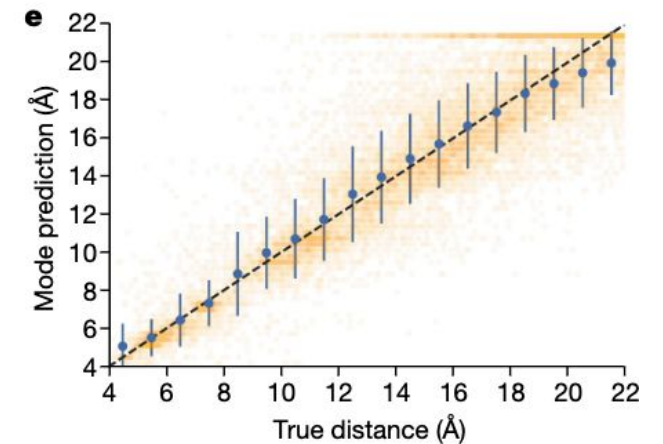
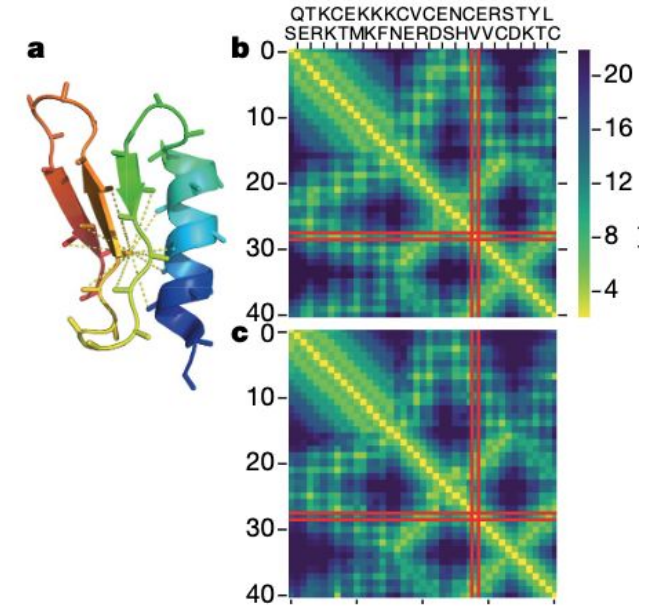
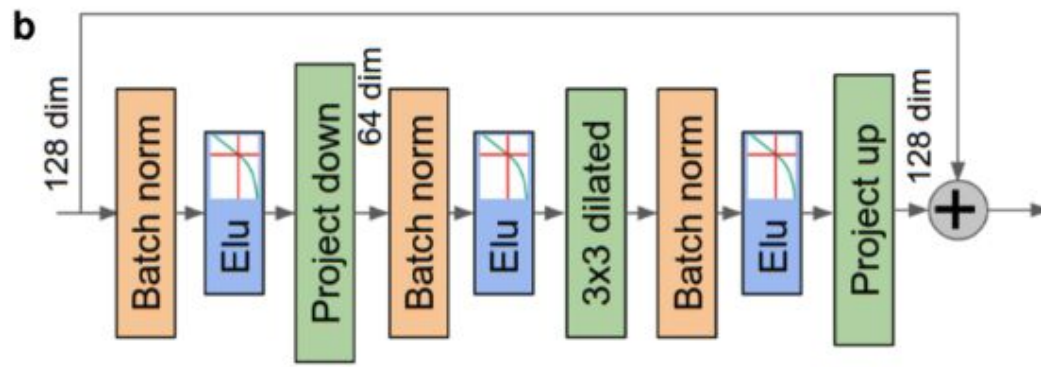
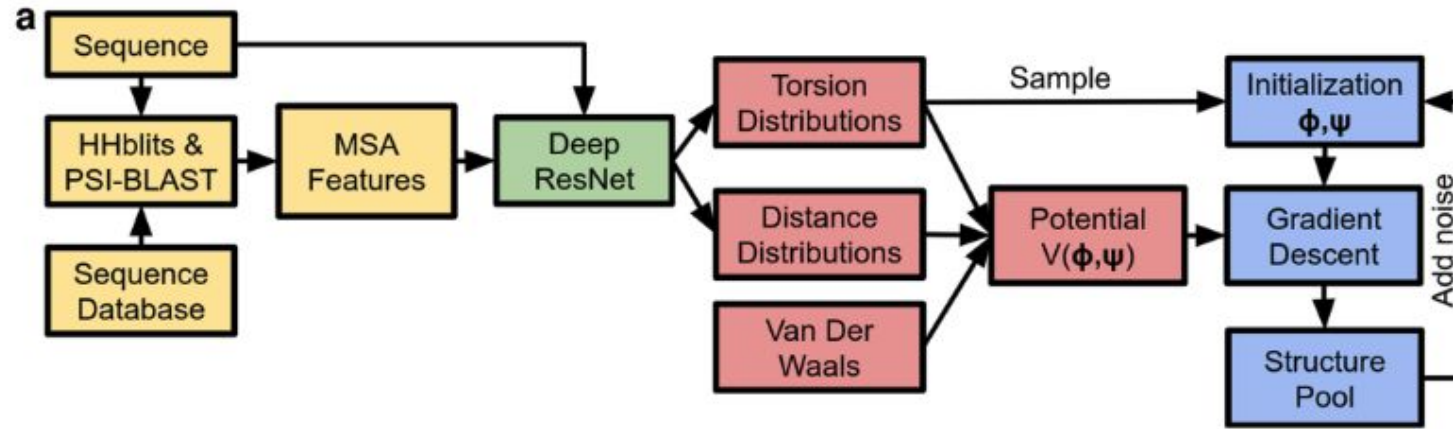




# AlphaFold v1: Model Overview

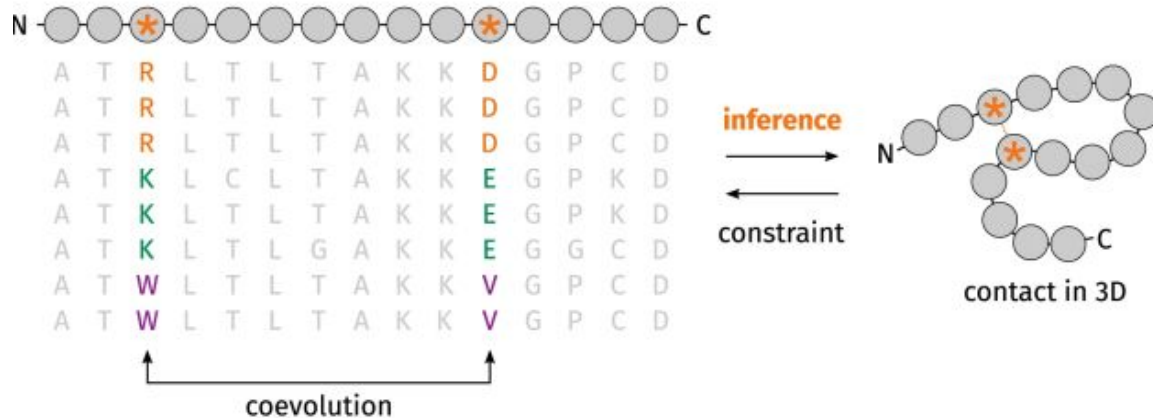


# AlphaFold v1: Model Details



# Multiple Sequence Alignment (MSA)

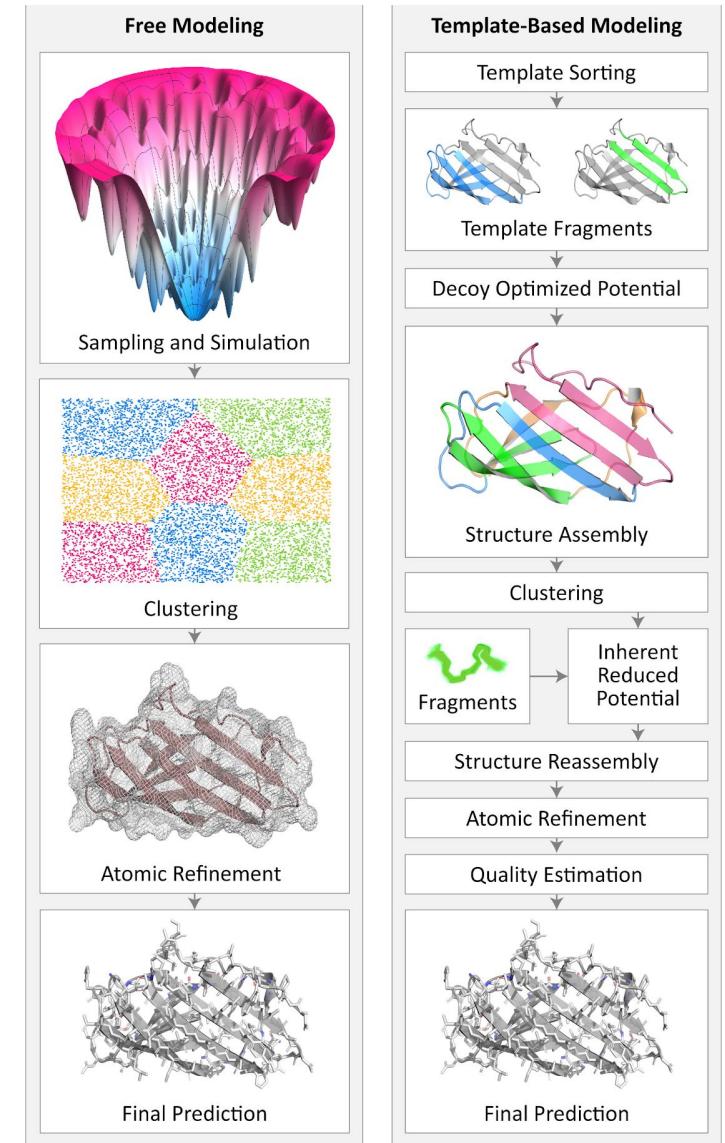
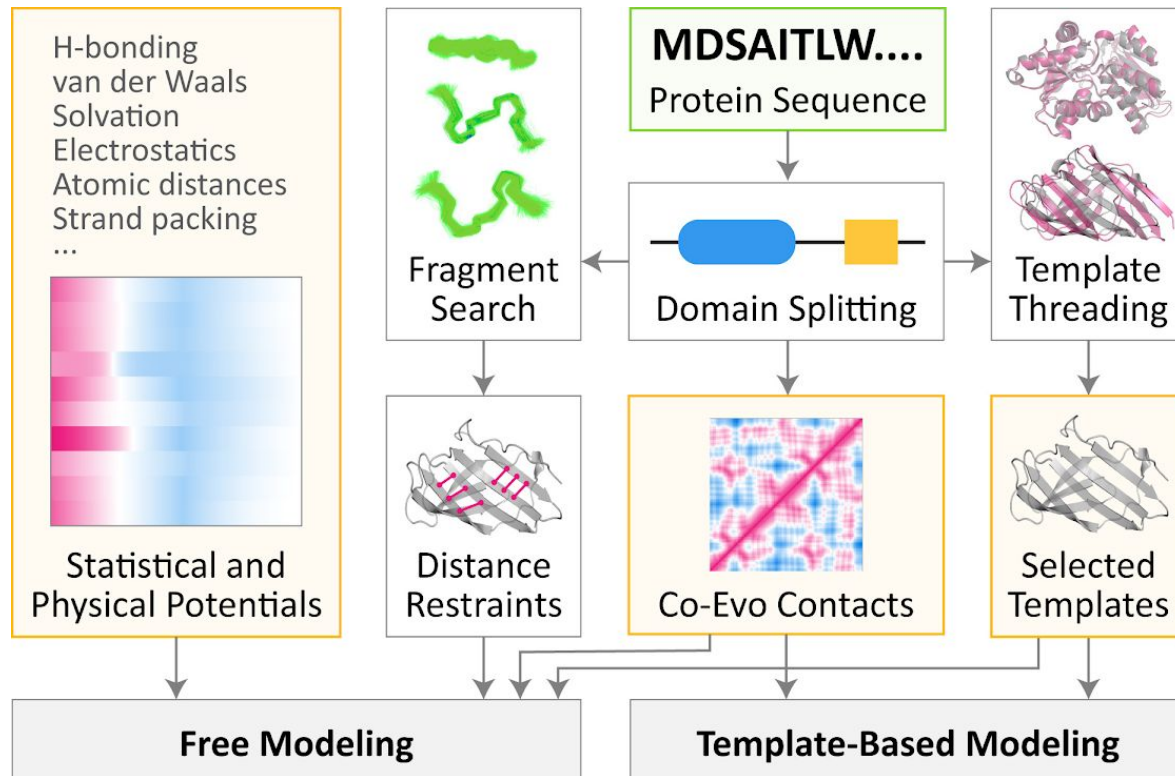
- Refer to the process or the result of sequence alignment of three or more biological sequences
- In AlphaFold, MSA is used to generate feature maps.
- Important indicator for structure information



Q5E940_BOVIN	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_HUMAN	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_MOUSE	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_RAT	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_CHICK	-----	MPREDRATWKSNYFMKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_RANSY	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	SALE	76	
Q7ZUG3_BRARE	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0 ICTPU	-----	MPREDRATWKSNYFLKIIQLDDYPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76	
RLA0_DROME	-----	MVRENKAAWKAQYFIKVVLELDFEFPKCFIVGADNVGSKMQQIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	POLE	76	
RLA0_DICDI	-----	MSGAGSKRRKLFIEKATKLFITVDMVIAEADVFVGSOLQKIRKSIIRGI	GAVLMGKNTMIRKIVIRDLADSK	PELD	75	
Q54LPO_DICDI	-----	MSGAGSKRRKVFIEKATKLFITVDMVIAEADVFVGSOLQKIRKSIIRGI	GAVLMGKNTMIRKIVIRDLADSK	PELD	75	
RLA0_PLAF8	-----	MAKLSKQQKQMYTEKLSLQQYKILIVHVDNVGSKMQQIRMSLRGK	AVVLMGKNTMIRKIVIRDLADSK	PALE	76	
RLA0_SULAC	-----	MIGLAVTTTKKIAKWKVEVAELTCKLKTHTIIIANIEGFPADKLHEIRKCLRK	ADIKVTKNLNFNIAKKNAG	VDTK	79	
RLA0_SULTO	-----	MRIMAVITQERKIAKWKIEVKELECKLREYHTIIIANIEGFPADKLHEIRKCLRK	ADIKVTKNLNFNIAKKNAG	LDVS	80	
RLA0_SULSO	-----	MKRLALALKQRKVASWKLIEYKELTELKNSNTLLIGNLEGFADKLHEIRKCLRK	ADIKVTKNLNFNIAKKNAG	IDIE	80	
RLA0_AERPE	-----	MSVSVLVGOMYKREKPIDEWKTLMLELELFSKIRVVFADLTGPIFFVYRVKRLWKK	YDMMVAKRILLRANKAAGLE	LDNN	86	
RLA0_PYRAE	-----	MMLAIGKRRYVRTQYPAKRVKIYSEATLFLQKPYVYFLFDLHGLSRIIHEIYRYLRRY	GVIKIIPKLFKIAATKYVGG	IPAE	85	
RLA0_METAC	-----	MAEERHHTHEIPQWKKDEIENIKELIQSHKVFQMVREGLATKMKIRRDLDV	AVLKVSRNTLIERALNQLG	ETIP	78	
RLA0_METMA	-----	MAEERHHTHEIPQWKKDEIENIKELIQSHKVFQMVREGLATKMKIRRDLDV	AVLKVSRNTLIERALNQLG	ESIP	78	
RLA0_ARCFU	-----	MAAVRGS---PPEYKVRAVEEIKRMISSKPVVAIVSFRNVPAGOMKIRREFRKG	AEIKVVTNLLERKDALG	GDYL	75	
RLA0_METKA	-----	MAVKAQKQPSGQYKPAVWKKRREVKELKLMDEYENVGLVDLEGLPAPOLQEI	TRAKLRERDTIIRMSRNTLMR	TALEEKLEDER	PELE	88
RLA0_METTH	-----	MAHVAEWKKKEVQELHDLKQYEVVGTANLADIPAROLQKMRQTLRDS	ALIRMSKNTLISLAEKAGREL	ENVD	74	
RLA0_METFL	-----	MTAESEHKIAPWKIEEYVNLKLELKNQIIVAVDMMVPPAROLQEIIRDKIR	GTMLLKMSRNTLIERAIEVAEETGNP	PEFA	82	
RLA0_METVA	-----	MIDAKSEHKIAPWKIEEYVNLKLELKNQIIVAVDMMVPPAROLQEIIRDKIR	DOMLLKMSRNTLIERAYEEVAEETGNP	PEFA	82	
RLA0_METJA	-----	METKVKAHVAPWKIEEYVNLKLELKNQIIVAVDMMVPPAROLQEIIRDKIR	DKVVLKMSRNTLIERAIEVAEELN	PNKLA	81	
RLA0_PYRAB	-----	MAHVAEWKKKEVEELANLKSFPVIALVDVYSSMPAYPLSQMRRLI	RENGGLLRVSRNTLIERAIEVAEELN	PKPELE	77	
RLA0_PYRHO	-----	MAHVAEWKKKEVEELANLKSFPVIALVDVYSSMPAYPLSQMRRLI	RENGGLLRVSRNTLIERAIEVAEELN	PKPELE	77	
RLA0_PYRFU	-----	MAHVAEWKKKEVEELANLKSFPVIALVDVYSSMPAYPLSQMRRLI	RENGGLLRVSRNTLIERAIEVAEELN	PKPELE	77	
RLA0_PYRKO	-----	MAHVAEWKKKEVEELANLKSFPVIALVDVYSSMPAYPLSQMRRLI	RENGGLLRVSRNTLIERAIEVAEELN	PKPELE	76	
RLA0_HALMA	-----	MSAESEKRTETIPQWQEEVDIVMIESYESVGVVNTIACIPSRLODMRRDLHGT	AEIIRVSRNTLIERALDDVD	DGLE	79	
RLA0_HALVO	-----	MSESEVRQTEVIPQWQEEVDIVMIESYESVGVVNTIACIPSRLODMRRDLHGT	AAVIRVSRNTLIERALDDVD	DGFE	79	
RLA0_HALSA	-----	MSAEQRTTEVIPQWQEEVDIVMIESYESVGVVNTIACIPSRLODMRRDLHGT	AALIRVSRNTLIERALDDVD	DGLD	79	
RLA0_THEAC	-----	MKEYSQKKELVNEITRIKASRSVAIVDLAGIRROIDDIRGNRKG	INLVYIKTLLFKALENLGD	EKLS	72	
RLA0_THEVO	-----	MRKINPKKKEIYSELAQITTKSKVAIVDLAGIRROIDDIRGNRKG	VKIKVYKTLFLKALDSDND	EKLT	72	
RLA0_PICTO	-----	MTEPQWKIDFVKNLENEINSRKAIVSISKGRNNEFQKIRNSIRDK	ARIKYSRRLRLRLAIE	NG	72	
ruler	1.....10.....20.....30.....40.....50.....60.....70.....80.....90					

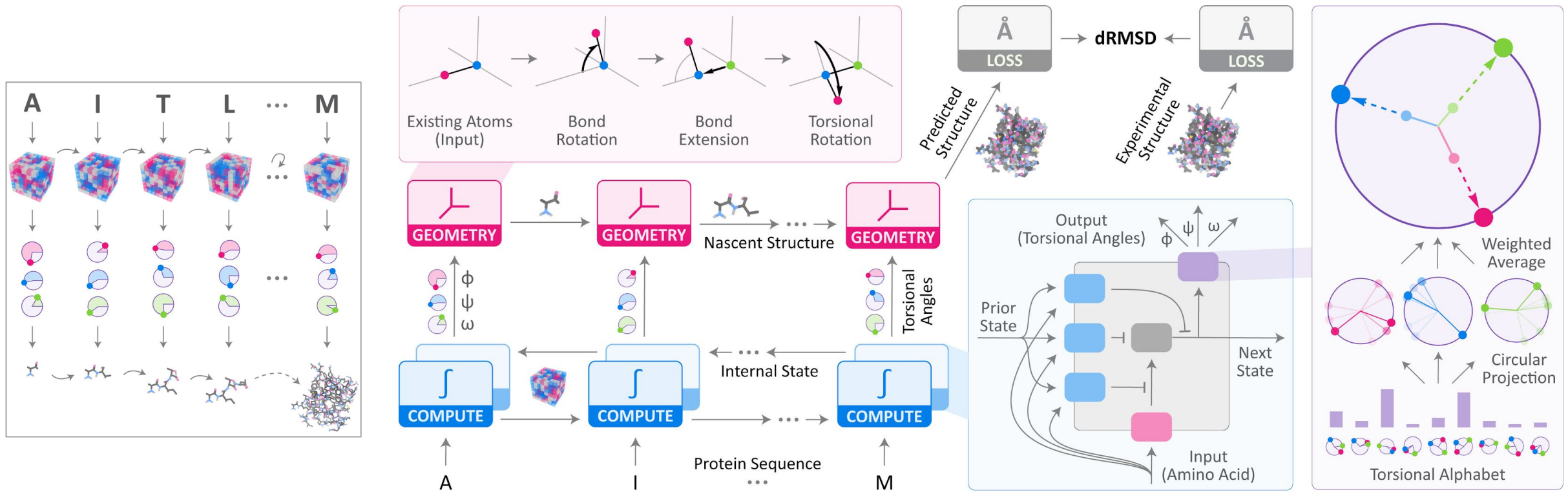
# Protein Folding: Conventional Pipeline

*Example of Feature processing pipeline*



# Another RNN Method for Protein Folding

- [End-to-End Differentiable Learning of Protein Structure](#), by Cell Systems



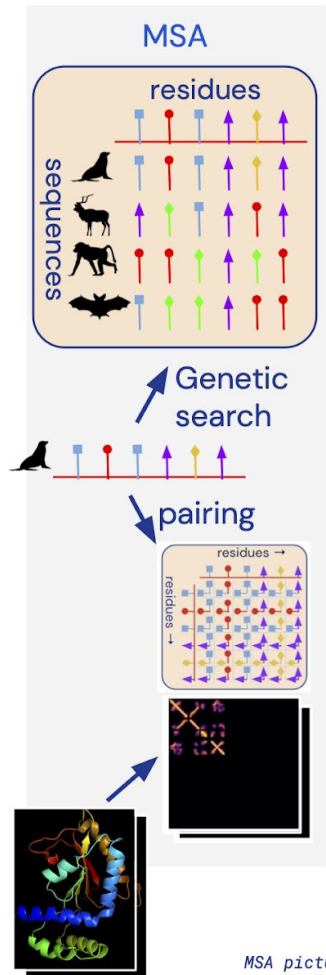
# AlphaFold v2: Highly accurate protein structure prediction with AlphaFold

*Where attention mechanism replace CNN and produce a breakthrough on the folding prediction*

# AlphaFold v2: Glance View

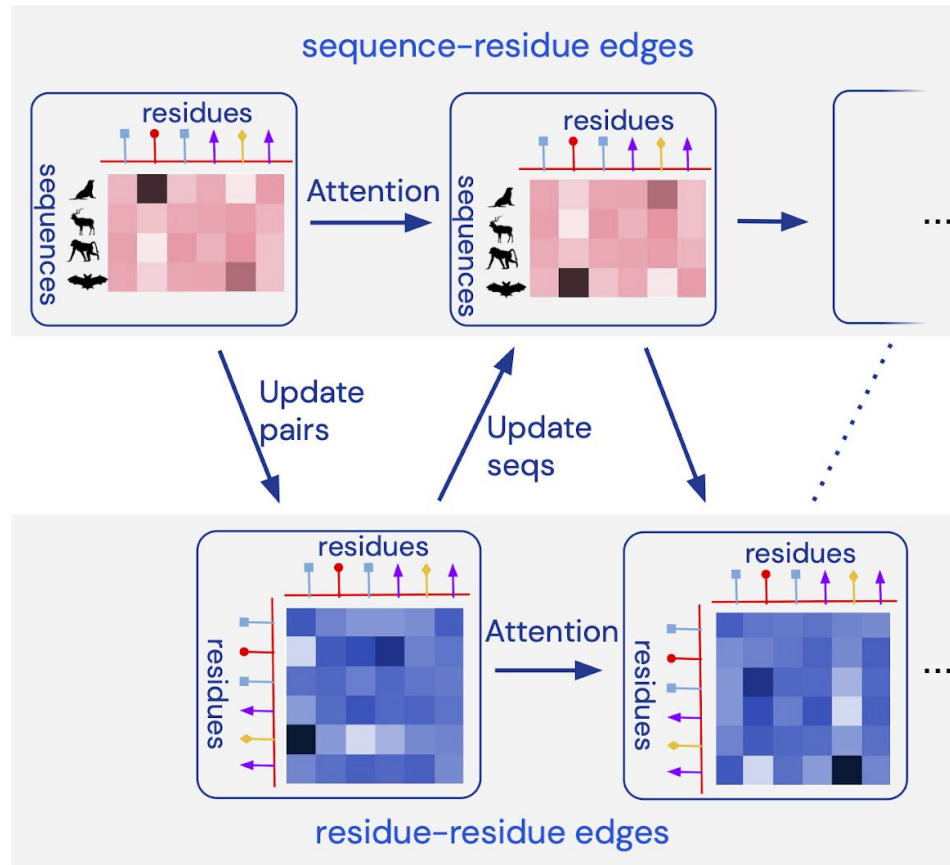


## Embedding



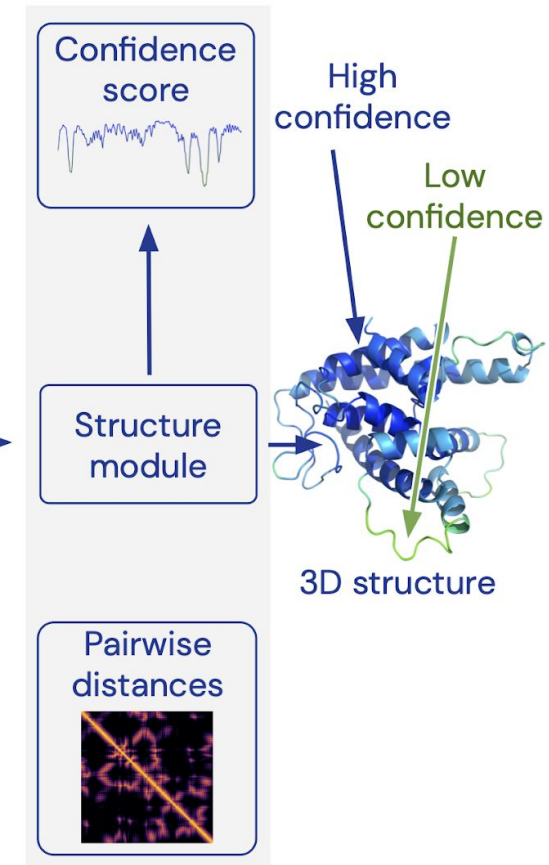
templates

## Trunk



## Heads

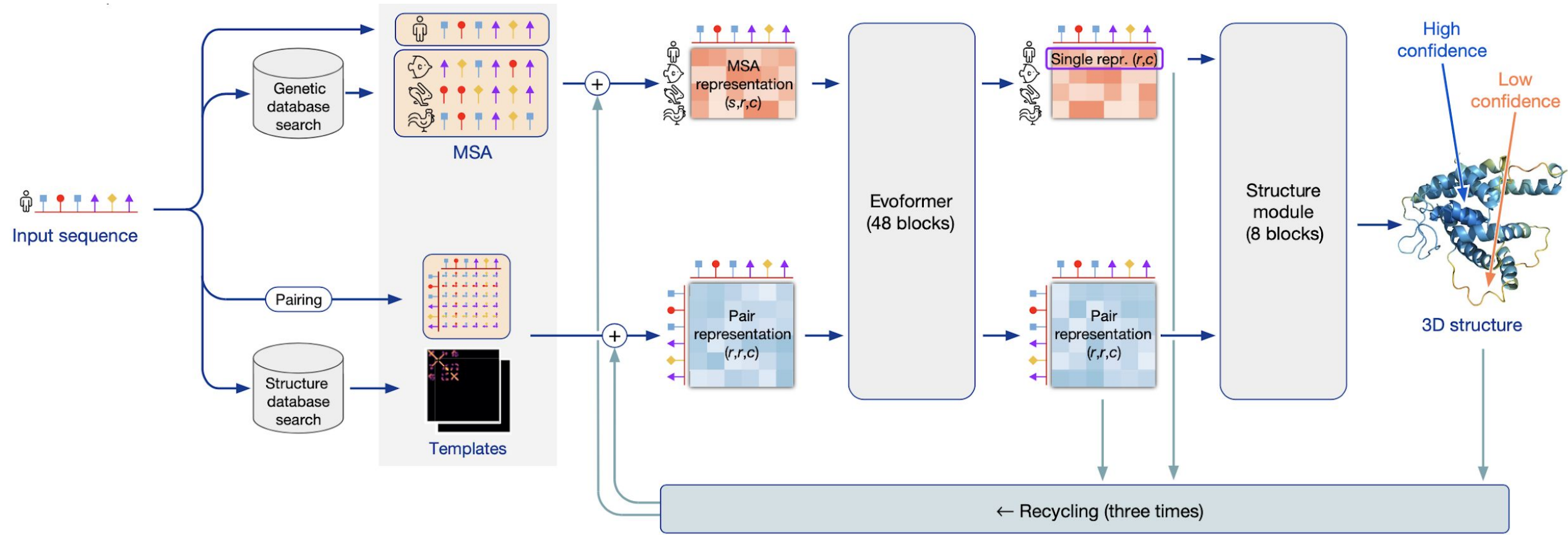
© 2020 DeepMind Technologies Limited



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

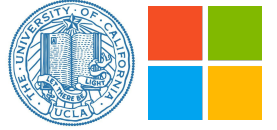


# AlphaFold v2: Glance View (Clearer Version)

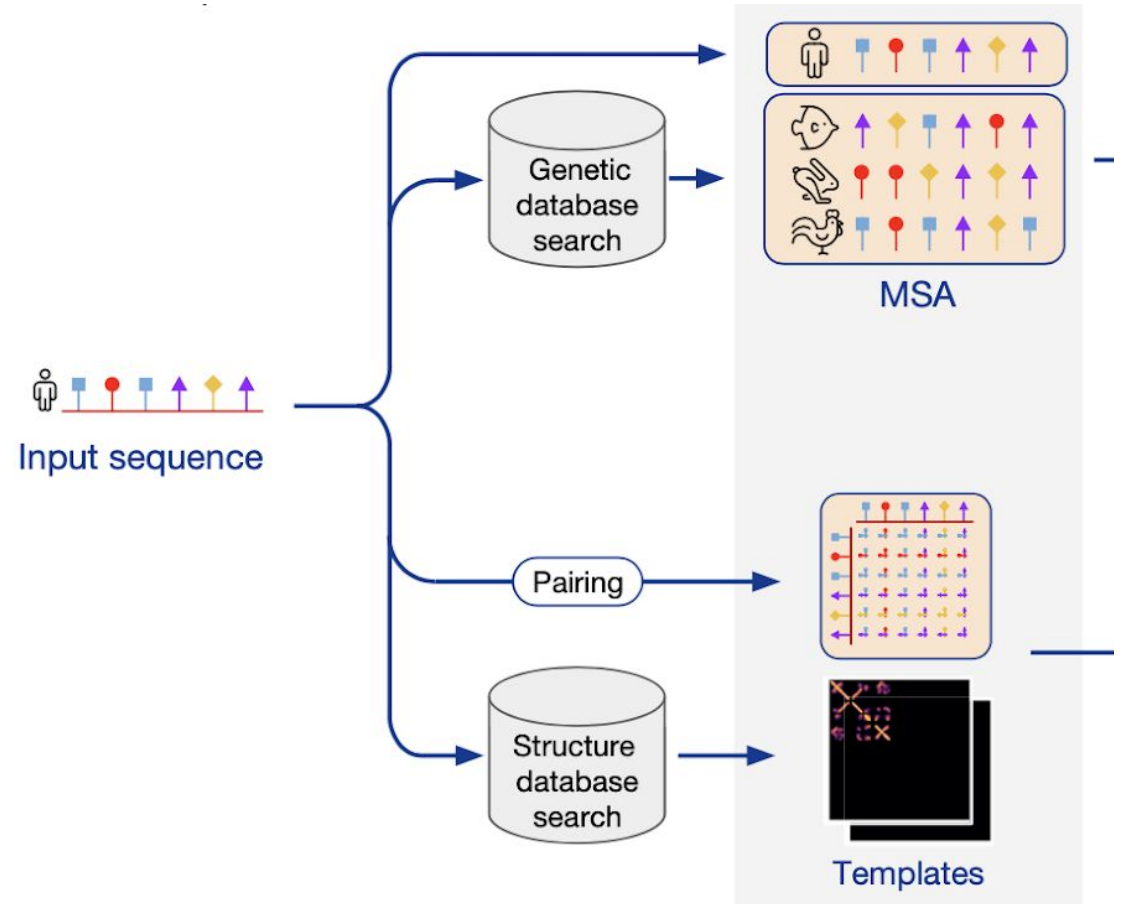




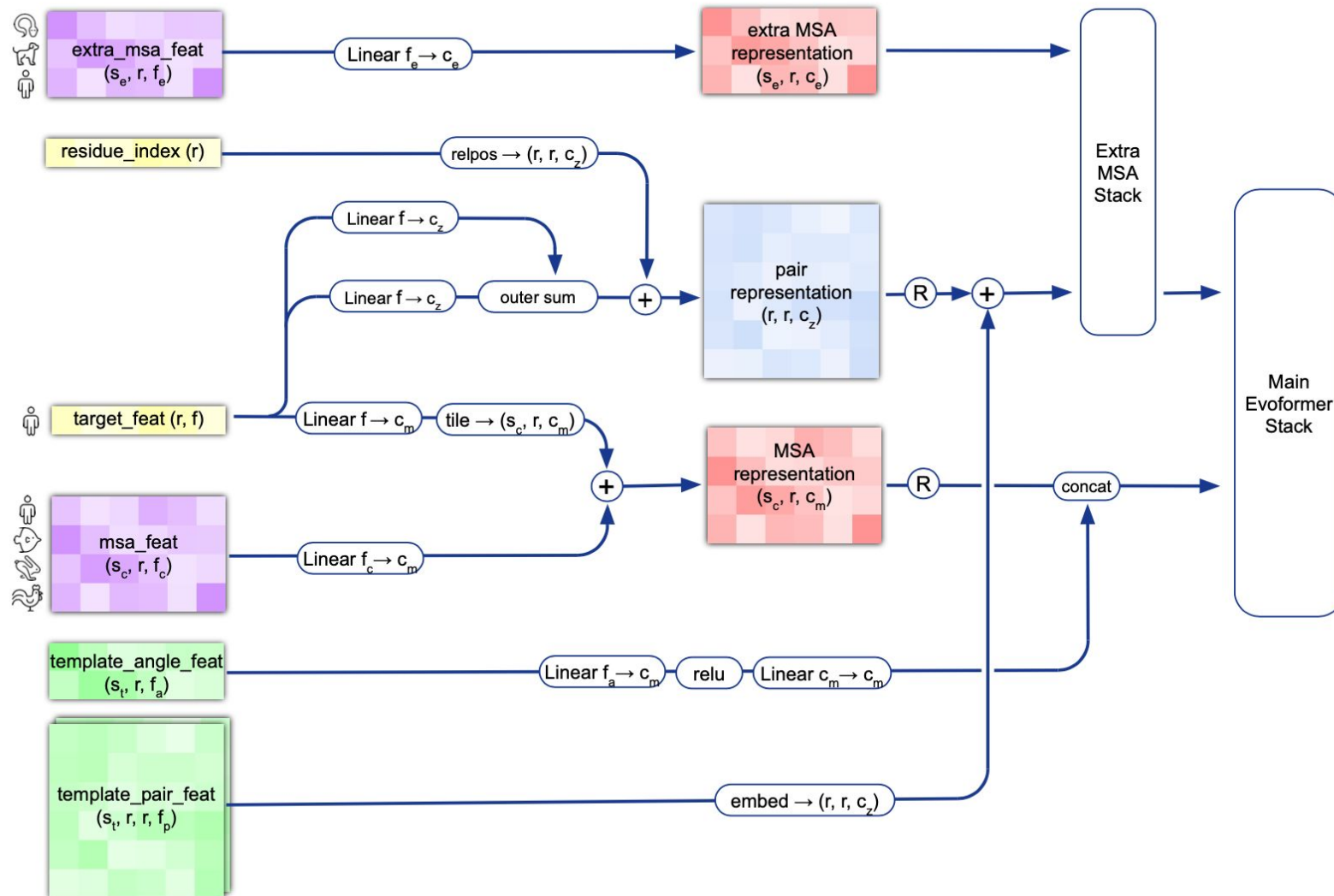
# AlphaFold2: Input



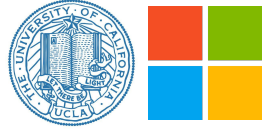
- Not significantly different from AlphaFold v1, or even other models
- Input sequence, and leveraging some known knowledge
- MSA (**sequence-residue** from genetic database), in the shape of  $(s,r,c)$
- Templates (**residue-residue**, structure database from known proteins), in the shape of  $(r,r,c)$



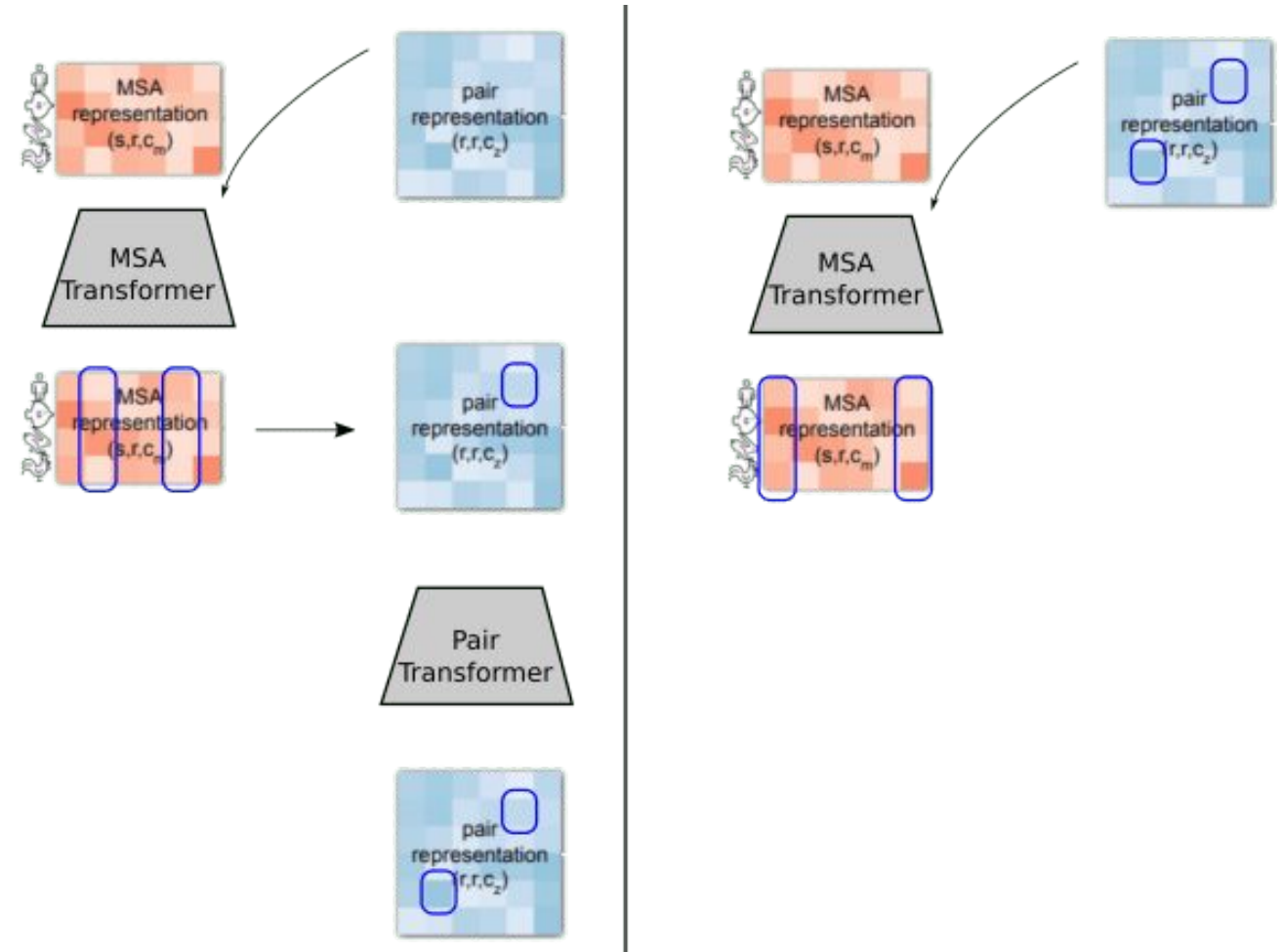
# AlphaFold2: Input (Complete Version)



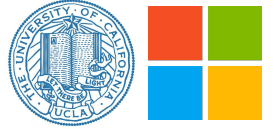
# Evoformer: Evolutionary Transformer?



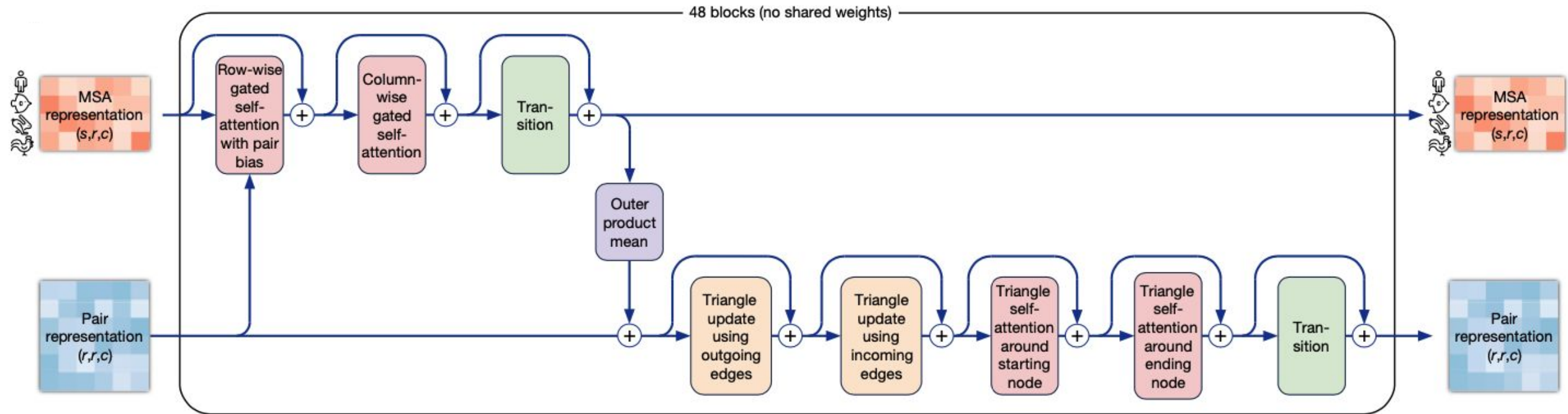
- Central idea: AlphaFold2 leverages the current structural hypothesis to improve the assessment of the multiple sequence alignment, which in turns leads to a new structural hypothesis, back and forth at every cycle.
- two transformers (a “two-tower architecture”), with one clear communication channel.



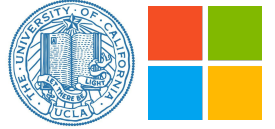
# Evoformer Block



- Information flow in one Evoformer Block. A total of 48 Evo blocks are used.



# Evoformer Stack: Algorithm Workflow



---

## Algorithm 6 Evoformer stack

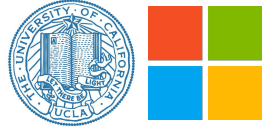
---

```
def EvoformerStack( $\{\mathbf{m}_{si}\}, \{\mathbf{z}_{ij}\}, N_{\text{block}} = 48, c_s = 384$ ) :  
  1: for all  $l \in [1, \dots, N_{\text{block}}]$  do  
    # MSA stack  
    2:  $\{\mathbf{m}_{si}\} += \text{DropoutRowwise}_{0.15}(\text{MSARowAttentionWithPairBias}(\{\mathbf{m}_{si}\}, \{\mathbf{z}_{ij}\}))$   
    3:  $\{\mathbf{m}_{si}\} += \text{MSAColumnAttention}(\{\mathbf{m}_{si}\})$   
    4:  $\{\mathbf{m}_{si}\} += \text{MSATransition}(\{\mathbf{m}_{si}\})$   
    # Communication  
    5:  $\{\mathbf{z}_{ij}\} += \text{OuterProductMean}(\{\mathbf{m}_{si}\})$   
    # Pair stack  
    6:  $\{\mathbf{z}_{ij}\} += \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationOutgoing}(\{\mathbf{z}_{ij}\}))$   
    7:  $\{\mathbf{z}_{ij}\} += \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationIncoming}(\{\mathbf{z}_{ij}\}))$   
    8:  $\{\mathbf{z}_{ij}\} += \text{DropoutRowwise}_{0.25}(\text{TriangleAttentionStartingNode}(\{\mathbf{z}_{ij}\}))$   
    9:  $\{\mathbf{z}_{ij}\} += \text{DropoutColumnwise}_{0.25}(\text{TriangleAttentionEndingNode}(\{\mathbf{z}_{ij}\}))$   
    10:  $\{\mathbf{z}_{ij}\} += \text{PairTransition}(\{\mathbf{z}_{ij}\})$   
  11: end for  
  # Extract the single representation  
  12:  $\mathbf{s}_i = \text{Linear}(\mathbf{m}_{1i})$   $\mathbf{s}_i \in \mathbb{R}^{c_s}$   
  13: return  $\{\mathbf{m}_{si}\}, \{\mathbf{z}_{ij}\}, \{\mathbf{s}_i\}$ 
```

---

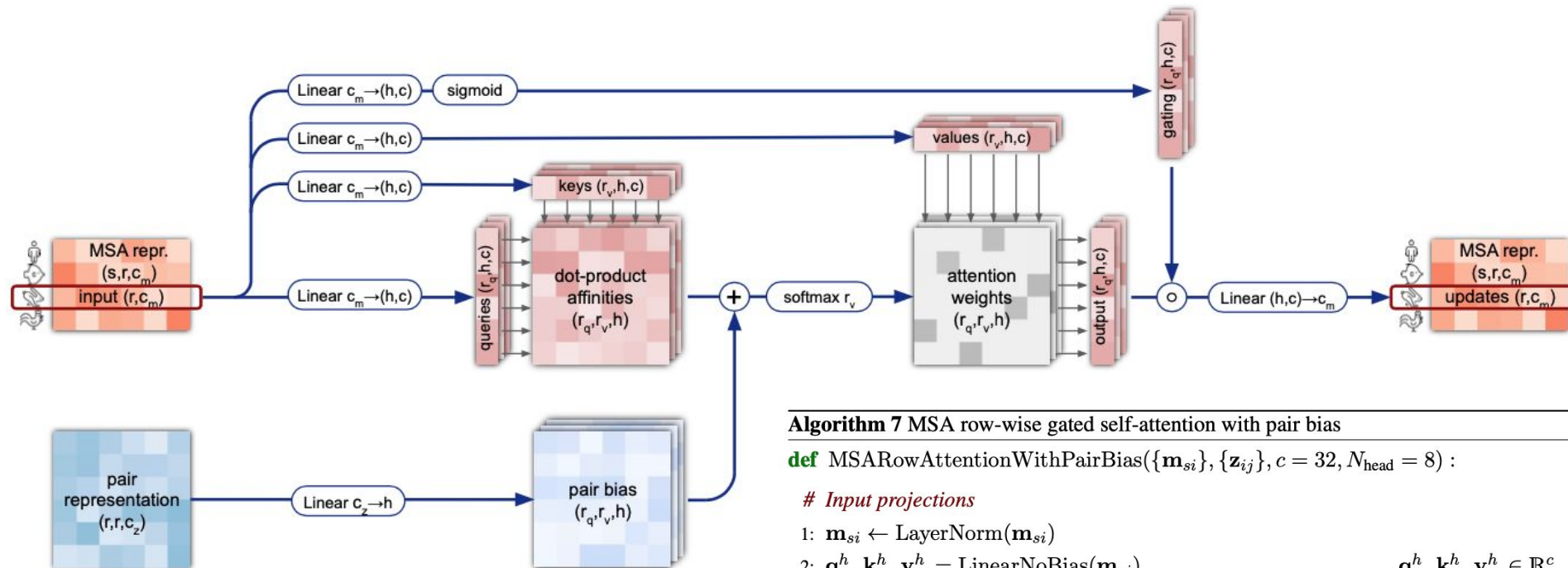
# AlphaFold2's MSA Transformer

---



- The attention is “factorized” in “row-wise” and “column-wise” components.
- MSA Transformer first computes attention in the horizontal direction, allowing the network to identify which pairs of amino acids are more related; and then in the vertical direction, determining which sequences are more informative.
- MSA Transformer’s row-wise (horizontal) attention mechanism incorporates information from the “pair representation”.
- Gated attention applied.

# AlphaFold2's MSA Row-wise Gated Attention



**Algorithm 7** MSA row-wise gated self-attention with pair bias

**def** MSARowAttentionWithPairBias( $\{\mathbf{m}_{si}\}, \{\mathbf{z}_{ij}\}, c = 32, N_{\text{head}} = 8$ ) :

*# Input projections*

- 1:  $\mathbf{m}_{si} \leftarrow \text{LayerNorm}(\mathbf{m}_{si})$
- 2:  $\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h = \text{LinearNoBias}(\mathbf{m}_{si})$
- 3:  $b_{ij}^h = \text{LinearNoBias}(\text{LayerNorm}(\mathbf{z}_{ij}))$
- 4:  $\mathbf{g}_{si}^h = \text{sigmoid}(\text{Linear}(\mathbf{m}_{si}))$

$$\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h \in \mathbb{R}^c, h \in \{1, \dots, N_{\text{head}}\}$$

$$\mathbf{g}_{si}^h \in \mathbb{R}^c$$

*# Attention*

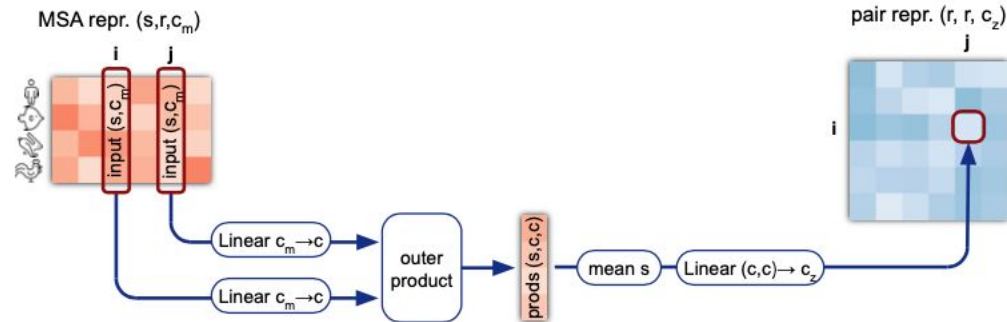
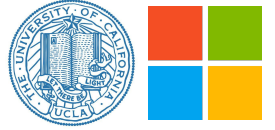
- 5:  $a_{sij}^h = \text{softmax}_j \left( \frac{1}{\sqrt{c}} \mathbf{q}_{si}^{h\top} \mathbf{k}_{sj}^h + b_{ij}^h \right)$
- 6:  $\mathbf{o}_{si}^h = \mathbf{g}_{si}^h \odot \sum_j a_{sij}^h \mathbf{v}_{sj}^h$

*# Output projection*

- 7:  $\tilde{\mathbf{m}}_{si} = \text{Linear}(\text{concat}_h(\mathbf{o}_{si}^h))$
- 8: **return**  $\{\tilde{\mathbf{m}}_{si}\}$

$$\tilde{\mathbf{m}}_{si} \in \mathbb{R}^{c_m}$$

# Evoformer: MSA Stack to Pair Stack



**Supplementary Figure 5** | Outer product mean. Dimensions: s: sequences, r: residues, c: channels.

---

## Algorithm 10 Outer product mean

---

**def** OuterProductMean( $\{\mathbf{m}_{si}\}, c = 32$ ) :

1:  $\mathbf{m}_{si} \leftarrow \text{LayerNorm}(\mathbf{m}_{si})$

2:  $\mathbf{a}_{si}, \mathbf{b}_{si} = \text{Linear}(\mathbf{m}_{si})$

3:  $\mathbf{o}_{ij} = \text{flatten}(\text{mean}_s(\mathbf{a}_{si} \otimes \mathbf{b}_{sj}))$

4:  $\mathbf{z}_{ij} = \text{Linear}(\mathbf{o}_{ij})$

5: **return**  $\{\mathbf{z}_{ij}\}$

$$\mathbf{a}_{si}, \mathbf{b}_{si} \in \mathbb{R}^c$$

$$\mathbf{o}_{ij} \in \mathbb{R}^{c \cdot c}$$

$$\mathbf{z}_{ij} \in \mathbb{R}^{c_z}$$

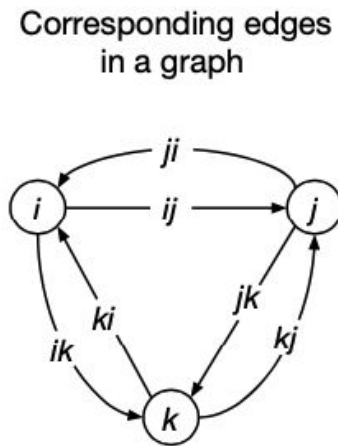
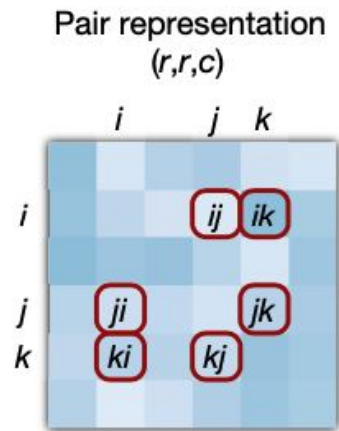

---



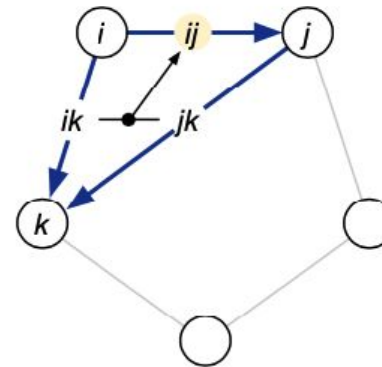
# AlphaFold2's Pair Transformer



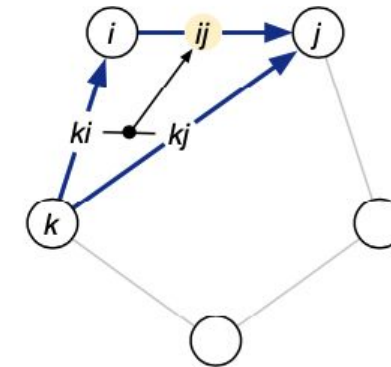
- Attention is arranged in terms of triangles of residues. □ Triangular attention



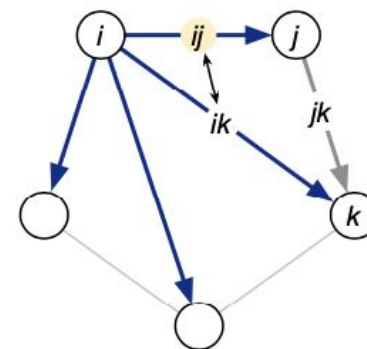
Triangle multiplicative update  
using 'outgoing' edges



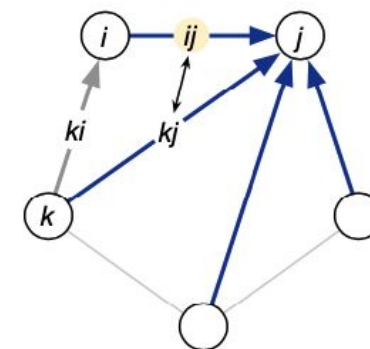
Triangle multiplicative update  
using 'incoming' edges



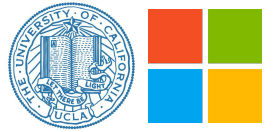
Triangle self-attention around  
starting node



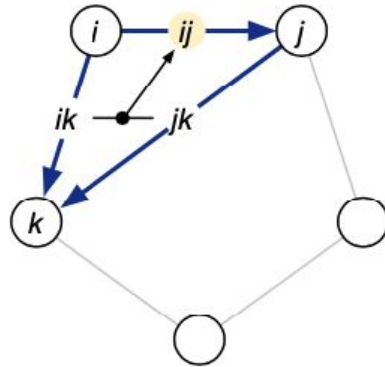
Triangle self-attention around  
ending node



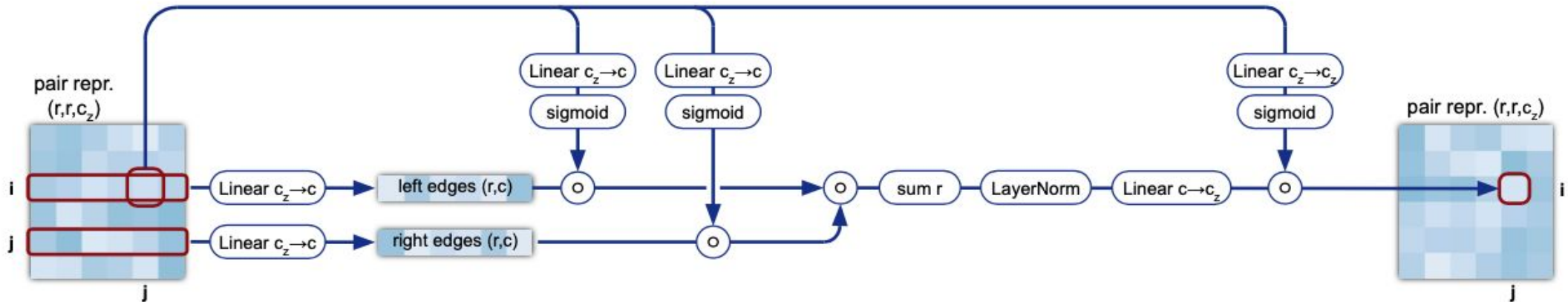
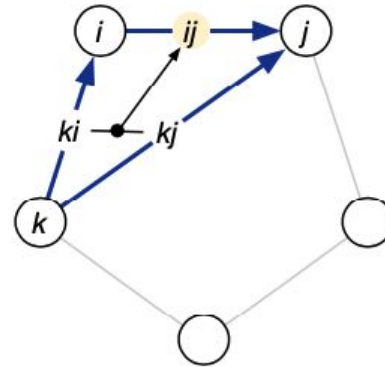
# Pair Transformer: Triangular Multiplicate Update



Triangle multiplicative update using 'outgoing' edges

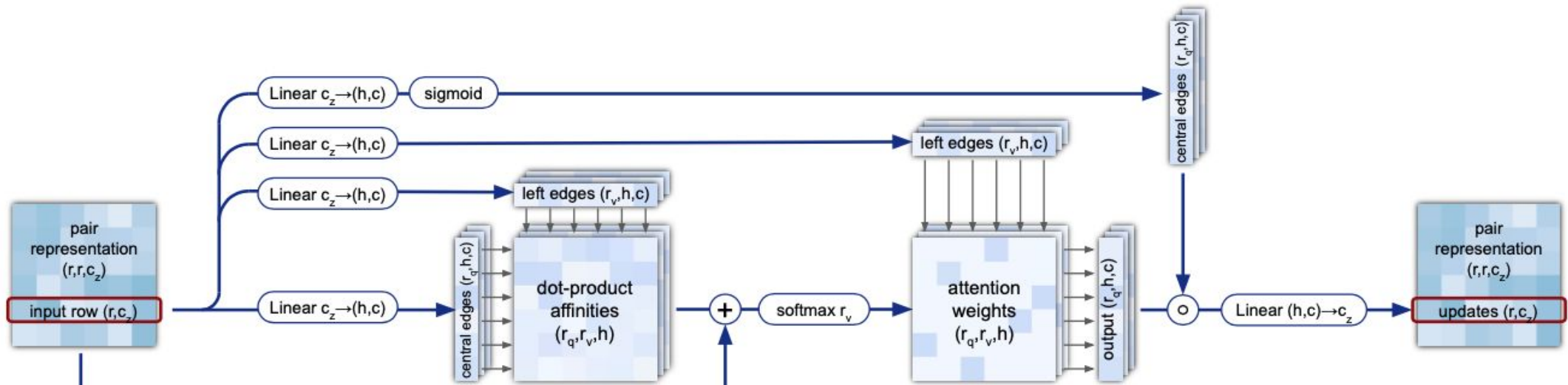
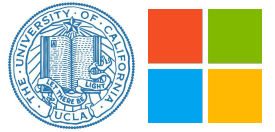


Triangle multiplicative update using 'incoming' edges

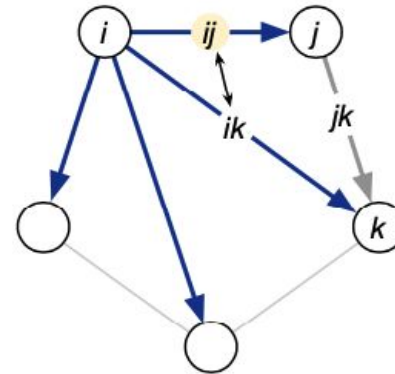


**Supplementary Figure 6** | Triangular multiplicative update using “outgoing” edges. Dimensions:  $r$ : residues,  $c$ : channels.

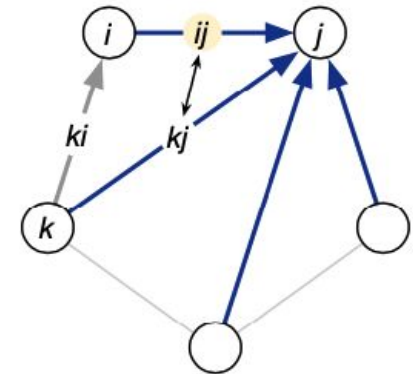
# Pair Transformer: Triangular Self-Attention



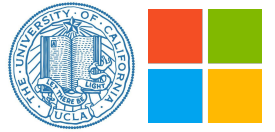
Triangle self-attention around starting node



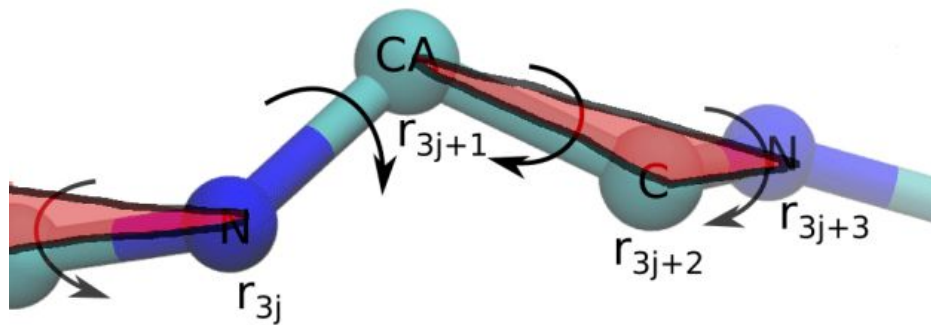
Triangle self-attention around ending node j



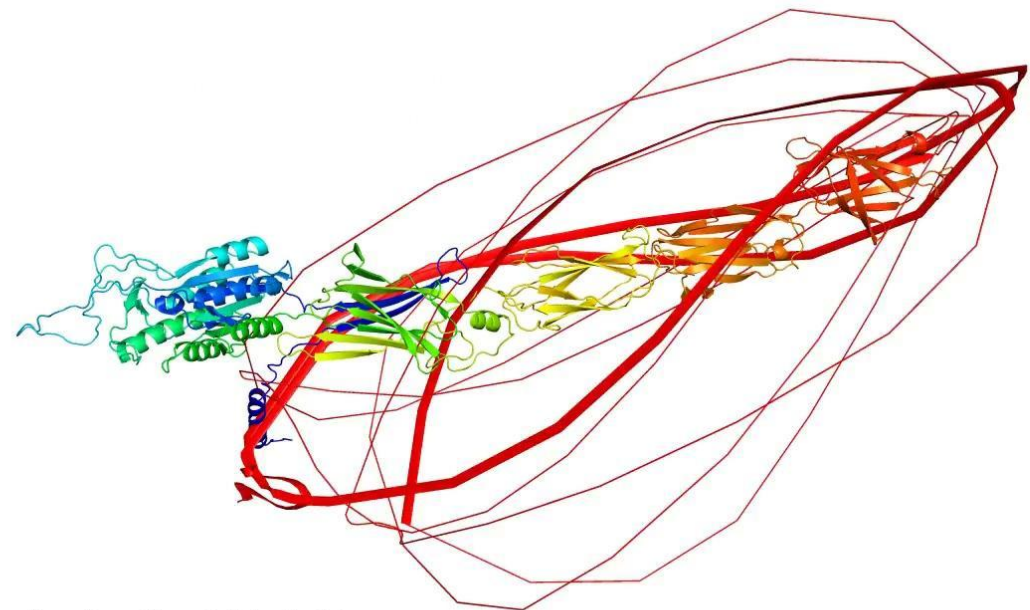
# AlphaFold2: Structure module



- The structure module considers the protein as a “residue gas”, a floating backbone.
- Every amino acid is modelled as a triangle, representing the three atoms of the backbone.
- These triangles float around in space and are moved by the network to form the structure.
- These transformations are parametrized as “affine matrices”.
- At every step of the iterative process, AlphaFold 2 produces a set of affine matrices that displace and rotate the residues in space. □ *There are potential structural violations in stereochemistry.*



$$\mathbf{M} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

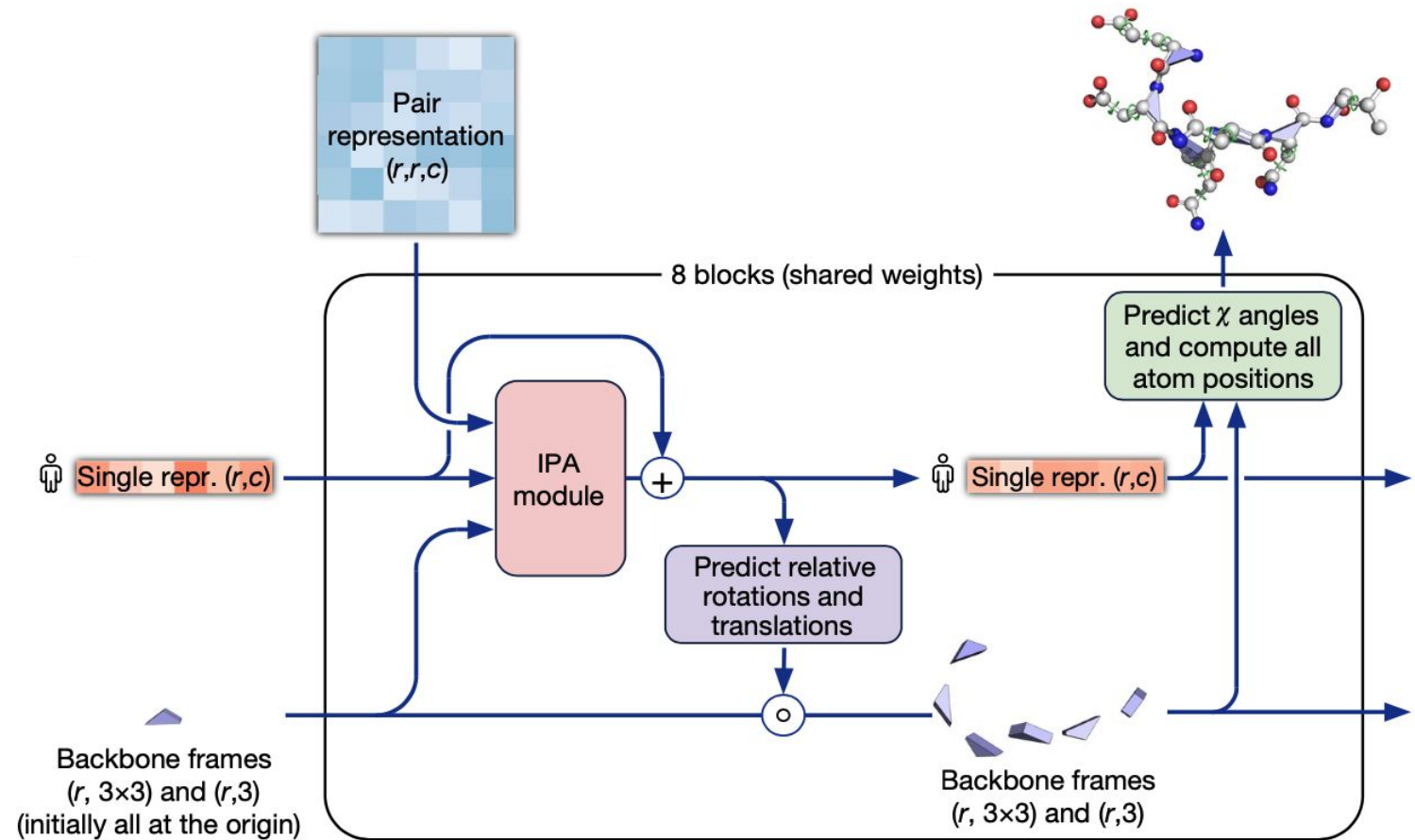
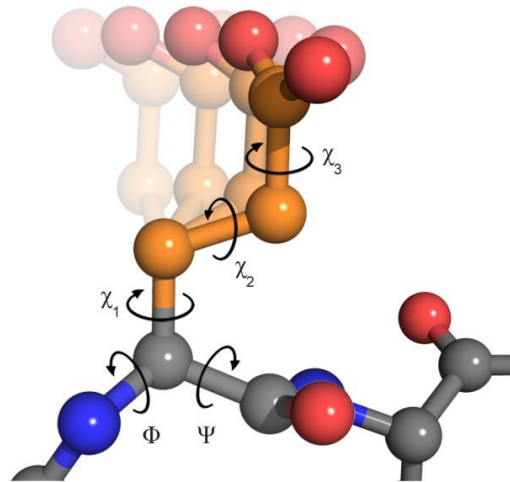


Recycling iteration 2, block 37  
Secondary structure assigned from the final prediction

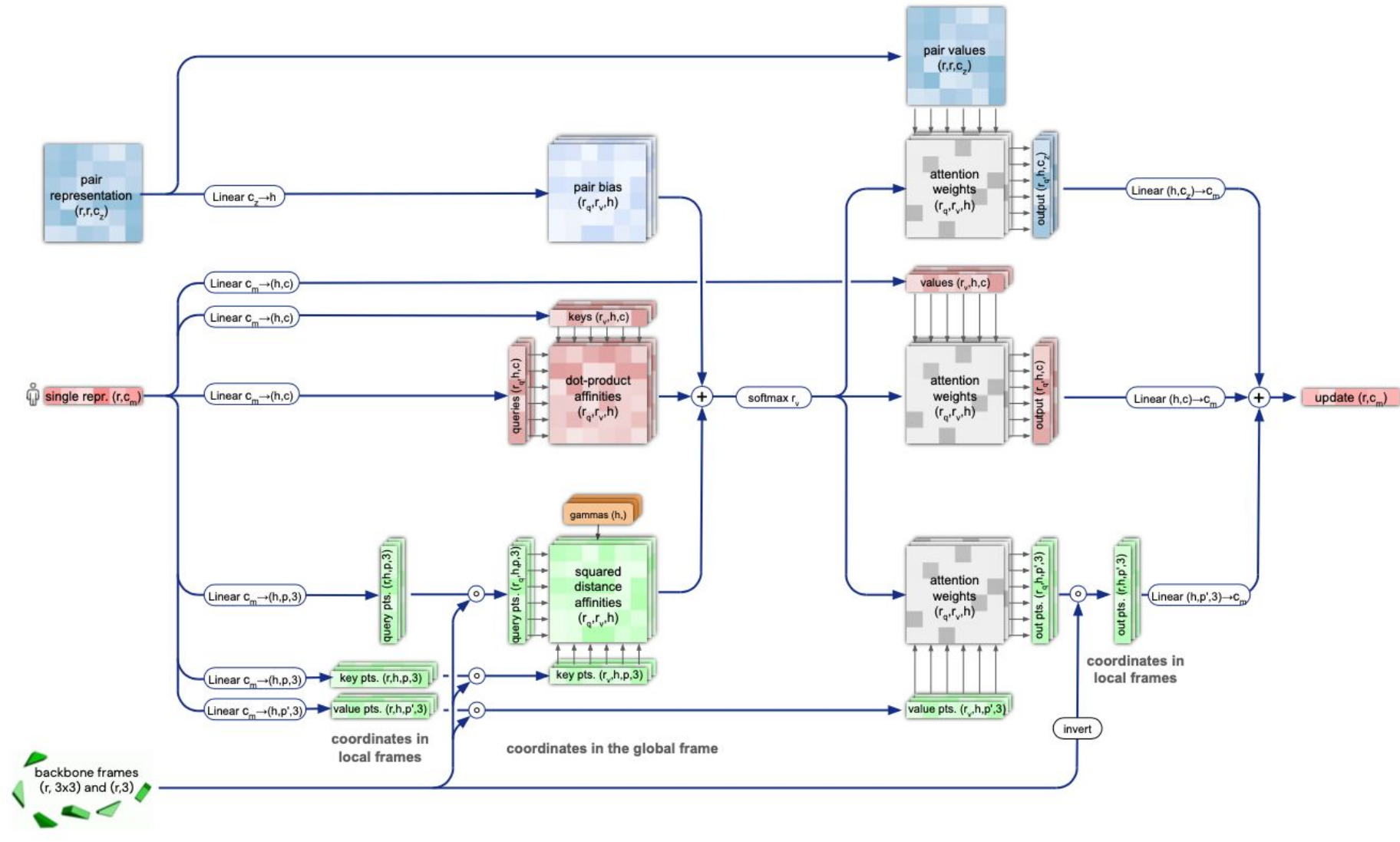
# AlphaFold2: Structure module



- Contains one module named Invariant point attention (IPA)

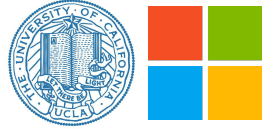


# AlphaFold2: IPA Module



# AlphaFold2: Quick Fact of Training Losses

---

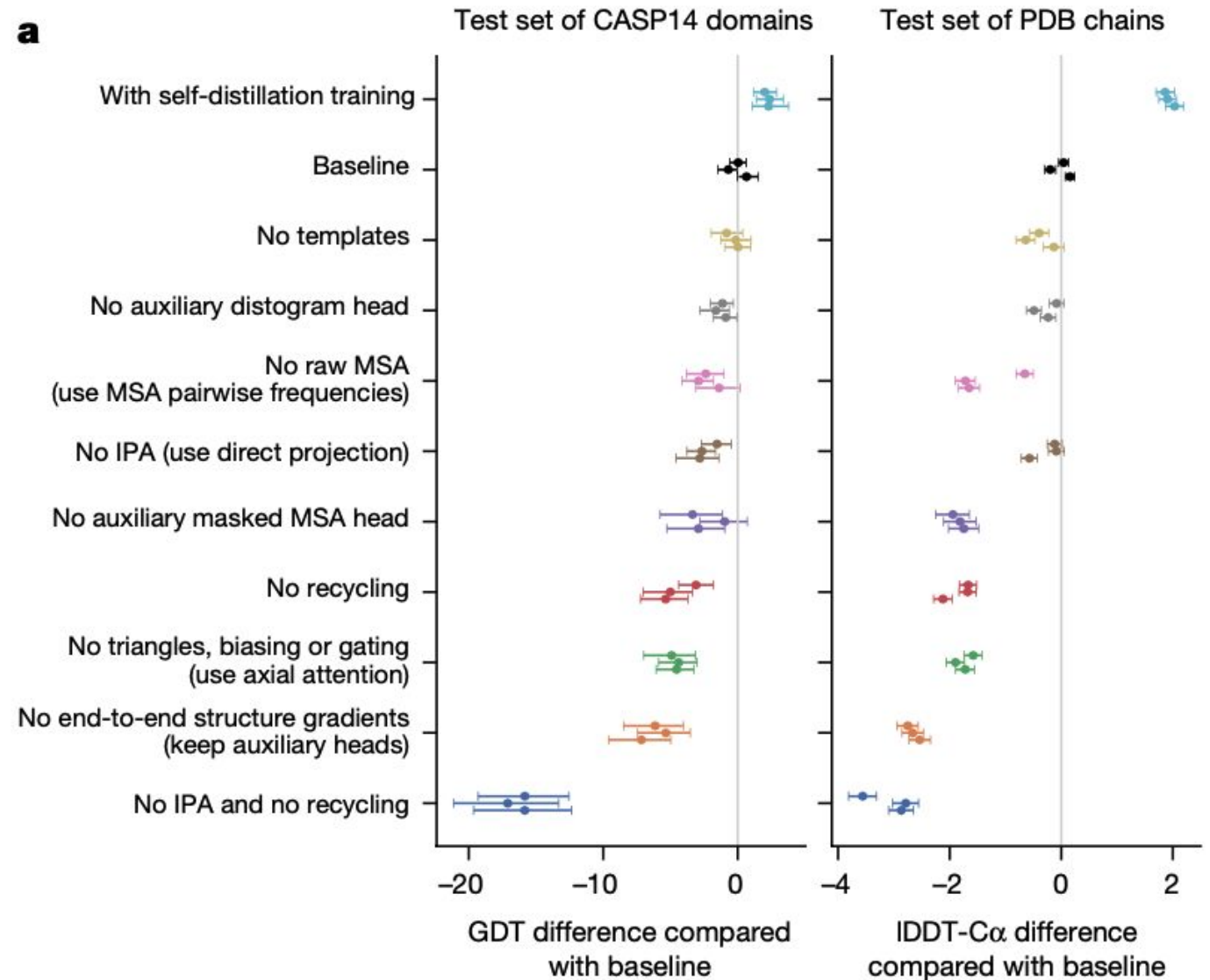


- Specific structural loss which is called FAPE (Frame Aligned Point Error)
- Auxiliary loss: MSA Masking
  - The model is given a multiple sequence alignment with some symbols “masked out” and asked to predict these symbols. → Self-supervision
- **Self-distillation**
  - In this approach, they took a model trained exclusively on the *PDB (full structure details available) and predicted the structures of ~300k diverse protein sequences obtained from Uniclust (no structure available).*
  - They then retrained the full model, incorporating a small random sample of these structures (a high-confidence subset) at every training cycle.
  - They claim this allows the model to leverage the large amount of unlabeled data available in protein sequence repositories.
- Other tricks...

# AlphaFold2: Tons of Engineering and Design



*Ablation study of multiple variants*





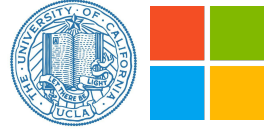
# Summary

- AI can contribute to basic scientific discovery, with the hope of making real-world impact, such as AlphaFold(2) in the realm of protein biology.
- A tool like AlphaFold might help rare disease researchers predict the shape of a protein of interest rapidly and economically.
- Physical insights are built into the network structure, instead of just data preprocessing or feature selection and curation.
- However, AlphaFold(2), similar to many computational biology models, are not verified nor experimented in “wet lab” and still skeptical to many biologists and pharmaceutical industry.

# AlphaFold v2: Protein Structure Database, Source Code and Demo

*Run AlphaFold2 on Google Colab*

# AlphaFold2 Protein Database



Demo (ACE2-HUMAN): <https://alphafold.ebi.ac.uk/entry/Q9BYF1>

Protein	Angiotensin-converting enzyme 2
Gene	ACE2
Source organism	Homo sapiens <a href="#">go to search</a>
UniProt	Q9BYF1 <a href="#">go to UniProt</a>
Experimental structures	63 structures in PDB for Q9BYF1 <a href="#">go to PDBe-KB</a>
Biological function	(Microbial infection) Non-functional as a receptor for human coronavirus SARS-CoV-2. <a href="#">go to UniProt</a>

## 3D viewer

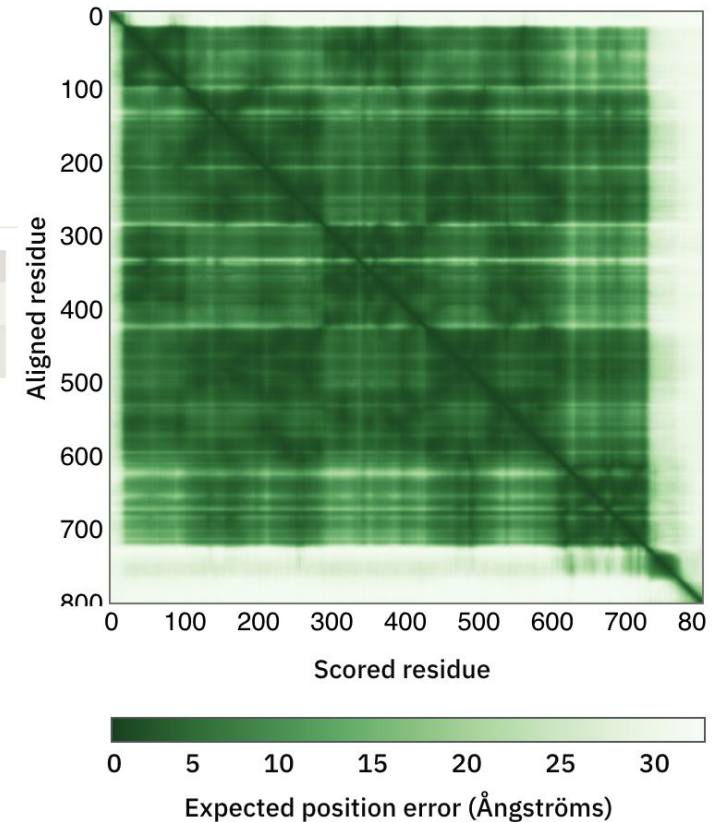
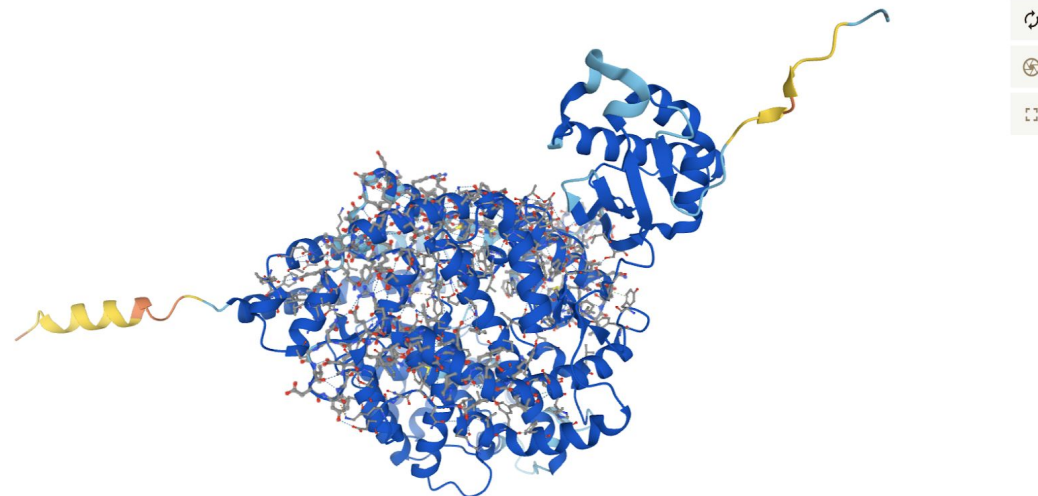
### Model Confidence:

- Very high (pLDDT > 90)
- Confident (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

AlphaFold produces a per-residue confidence score (pLDDT) between 0 and 100. Some regions below 50 pLDDT may be unstructured in isolation.

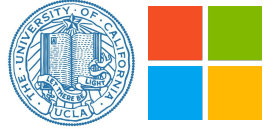
Sequence of AF-Q9BYF1-F1 1: Angiotensi... A

```
MSSSSWLLLSIVAVTAAQSTIEEQAKTFLDKFNHEAEDLFYQSSLASWNYNTNITEENVQNMNAGDKWSAFLKEQSTLAQMYPLQEIQNLTVKQLQALQQNGSSVLSSEDKSKRLNTILNTMS
TIYSTGKVCNPDNPQECLELLEPGLNEIMANSLDYNERLWAWESWRSEVGKQLRPLYEEYVVLKNEMARANHVEDYDGYWRGDYEVNGVDGYDYSRGQLIEDVEHTFEEIKPLYEHLHAYVRAKL
MNAYPSYISPIGCLPAHLLGDMWGRFWTNLYSLTVPFGQKPNIDVTAMVDQAWDAQRIFKEAEKFFVSVGLFNMTQGFWENSMLTDPGNVQKAVCHPTAWDLGKGFRIILMCTKVTMDDFLTA
```



# Source Code and AlphaFold on Google Colab

---



Source Code: <https://github.com/deepmind/alphafold/>

Original AlphaFold Colab:

<https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb>

AlphaFold2 and advanced version (\*not\* authored by Google/DeepMind):

<https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/AlphaFold2.ipynb>

[https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/beta/AlphaFold2\\_advanced.ipynb](https://colab.research.google.com/github/sokrypton/ColabFold/blob/main/beta/AlphaFold2_advanced.ipynb)

More Colab notebooks: <https://github.com/sokrypton/ColabFold/>

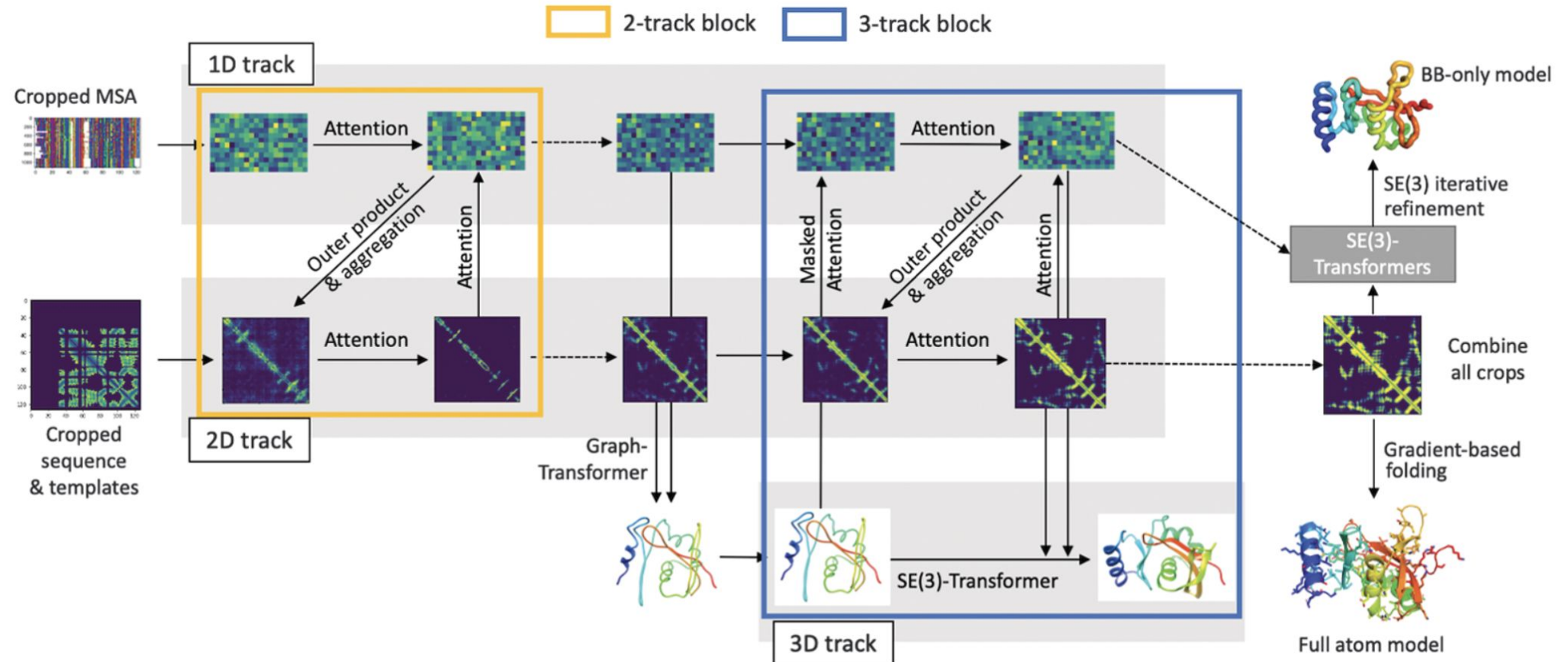
OpenFold2:

*Run AlphaFold2 on Google Colab*

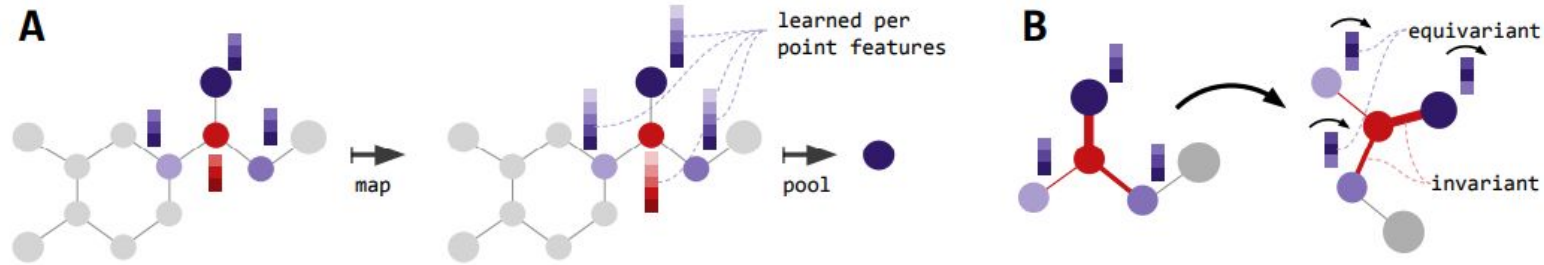
# New Paper on Science: RoseTTAFold



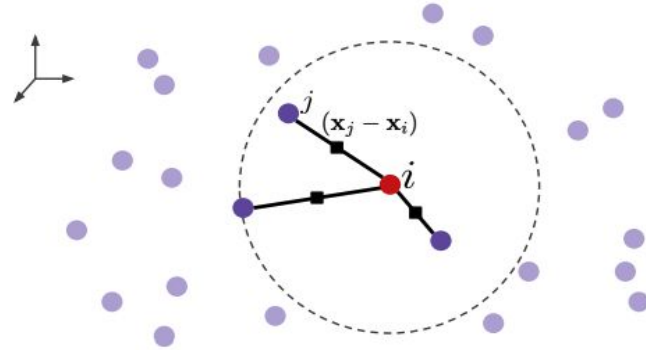
- [Accurate prediction of protein structures and interactions using a three-track neural network](#)
- Accuracy approaching closely on DeepMind's
- Claimed the model enables rapid generation of accurate protein-protein complex models



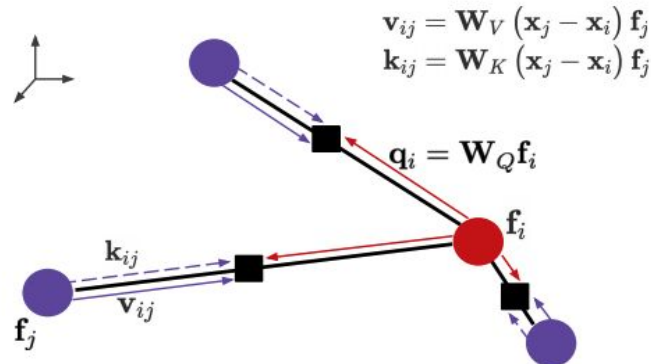
# SE(3)-Transformers [\[Paper\]](#)



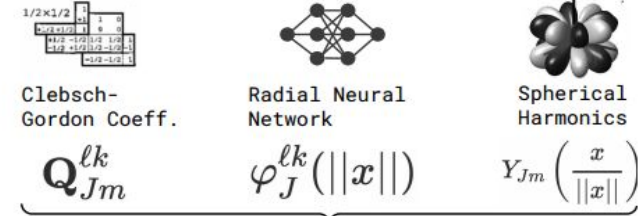
**Step 1: Get nearest neighbours and relative positions**



**Step 3: Propagate queries, keys, and values to edges**



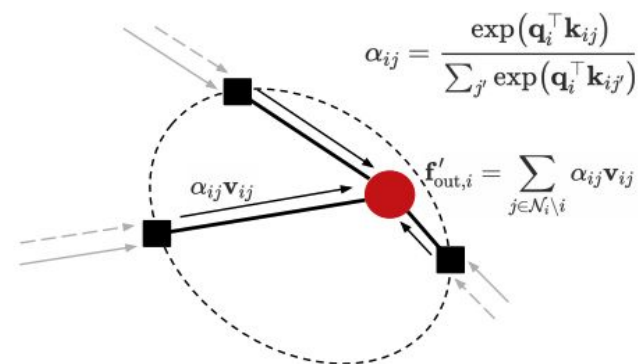
**Step 2: Get SO(3)-equivariant weight matrices**



Matrix  $W$  consists of blocks mapping between degrees

$$W(x) = W \left( \left\{ Q_{Jm}^{\ell k}, \varphi_J^{\ell k}(\|x\|), Y_{Jm} \left( \frac{x}{\|x\|} \right) \right\}_{J,m,\ell,k} \right)$$

**Step 4: Compute attention and aggregate**



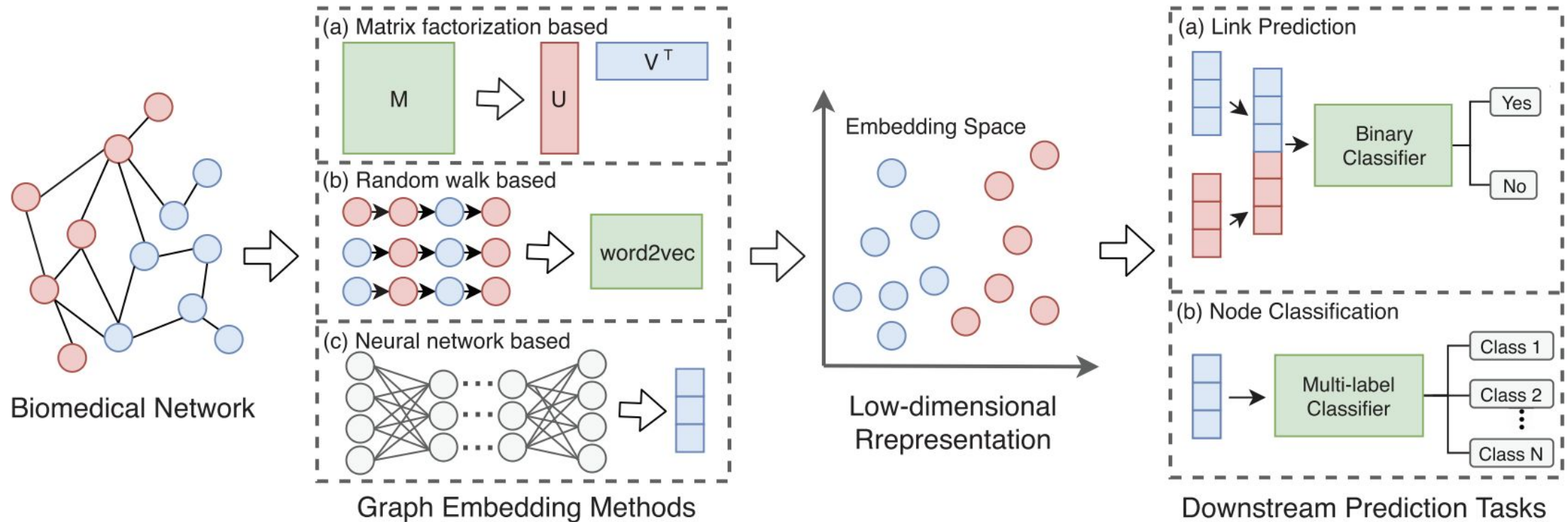
Discussion

## Network Science and Graph in Bioinformatics

*At the boundary between different fields, new  
“mountains” rise up.*

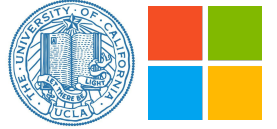


# Learn Embeddings on Biological Networks

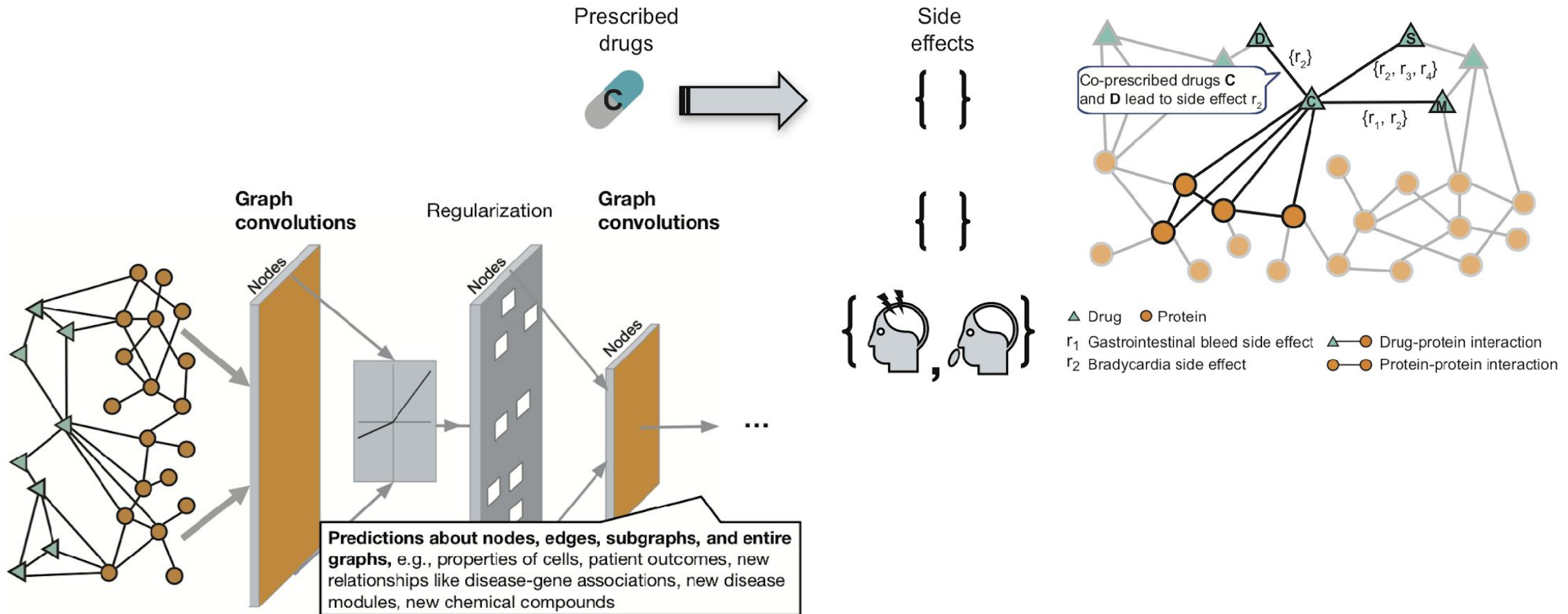


Reference: Yue, Xiang, et al. "Graph embedding on biomedical networks: methods, applications and evaluations." *Bioinformatics* 36.4 (2020): 1241-1251.

# Example: Drug-Protein Network, Side Effects

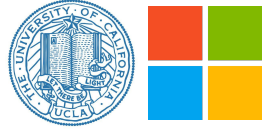


Credit: <https://zitniklab.hms.harvard.edu/research/>



# More Applications

---



- Molecular biology, compound structures, pathways
- Pandemic prediction, disease spreading
- Healthcare knowledge graphs, biomedical ontologies
- Clinical report analysis and personal health record

# DeepMind's AlphaFold Team & Posts

- <https://deepmind.com/blog/article/AlphaFold-Using-AI-for-scientific-discovery> (AlphaFold v1, Jan 2020)
- <https://deepmind.com/blog/article/alphafold-a-solution-to-a-50-year-old-grand-challenge-in-biology> (AlphaFold v2, Dec 2020)
- <https://deepmind.com/blog/article/putting-the-power-of-alphafold-into-the-worlds-hands> (AlphaFold v2 release, Jul 2021)

Resource List:

# AlphaFold and New Frontier of Protein Folding

[Tutorials, Blogs, and Related resources of AlphaFold and AlphaFold2](#) (collected by Junheng)

**UCLA**

**Samueli**  
Computer Science

Thank you!

Contact: [jhao@cs.ucla.edu](mailto:jhao@cs.ucla.edu)

Website: <http://www.haojunheng.com/>

Appendix

## Related Topics and Tutorials

*More about MSA, Protein structure and spatial representation, etc.*