# **AlphaFold:** Algorithmic Building Blocks for Protein Structure Prediction

**BIOE 145/245: Introduction to Machine Learning for Computational Biology**
Spring 2024

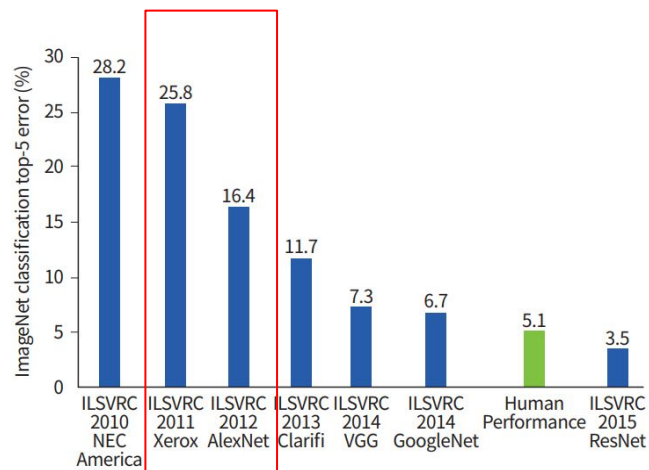**Guest Lecturer: Amy X. Lu**
PhD Student
Department of Electrical Engineering and Computer Science
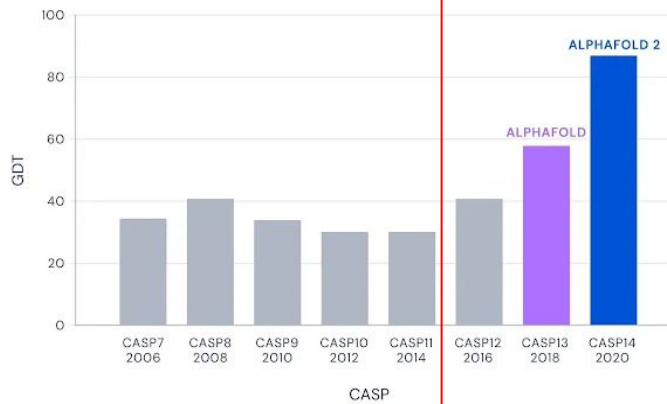Berkeley Artificial Intelligence Lab

# Agenda

1. The protein folding problem
2. AlphaFold2 and its internals
   a. Training data
   b. Evoformer
   c. Structure module
   d. Loss functions
3. Using AlphaFold
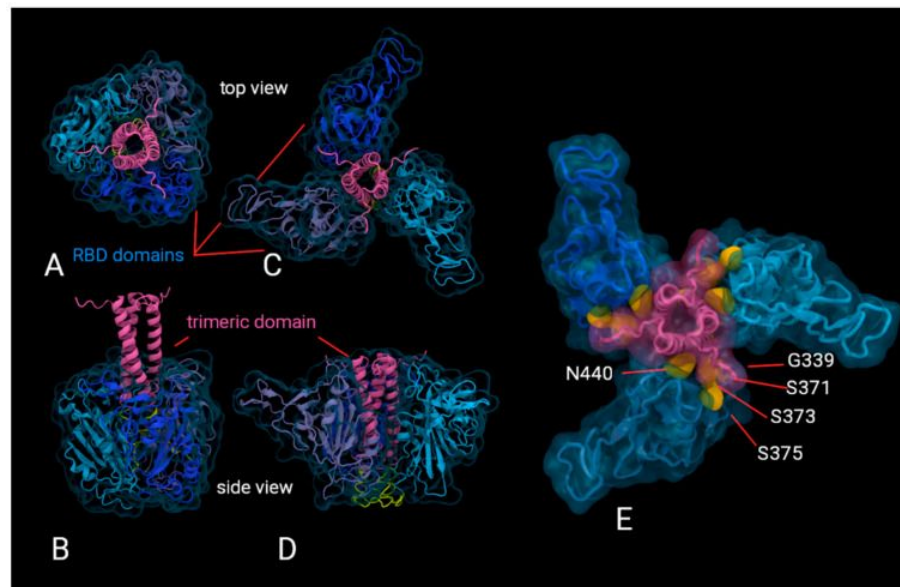4. Limitations and open problems

# AlphaFold: Biology's AlexNet moment?



**How can the key ingredients for AlphaFold's success be used in other AI for biology tasks?**

# Why study AlphaFold?

- As a **tool** for forming biological hypotheses
  - Democratization of estimated structure unlocks new opportunities in understanding biological mechanisms
- As a **repertoire of techniques**
  - Transferable lessons for capturing inductive biases in scientific data
  - Bridges the adoption of machine learning methods for science



Example: AlphaFold-predicted structures help us hypothesize how sequence-level mutations in the SARS-CoV2 Omicron variant impacts its mechanism.
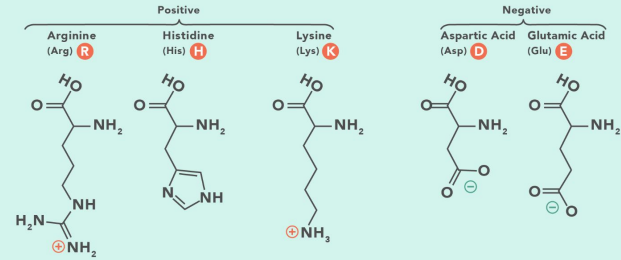(Source: van Vuren et al., 2022)

*Background:*
# The Protein Folding Problem

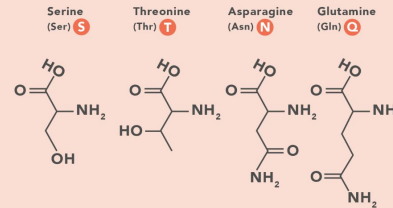*tl;dr: Predict 3D coordinates from 1D string*

# What is a protein?

- Workhorse of biological functions
- Built as an assortment of 20 types of amino acids
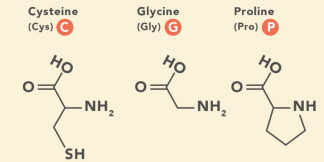  - Each amino acid consists of a distinct assortment of atoms
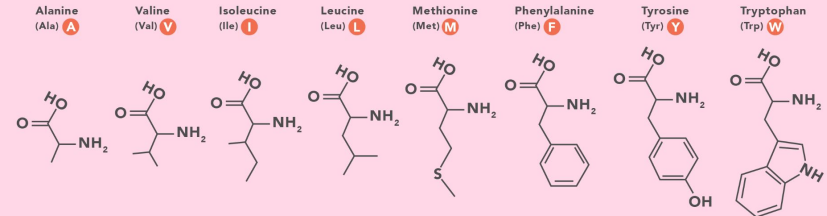


A. Amino Acids with Electrically Charged Side Chains
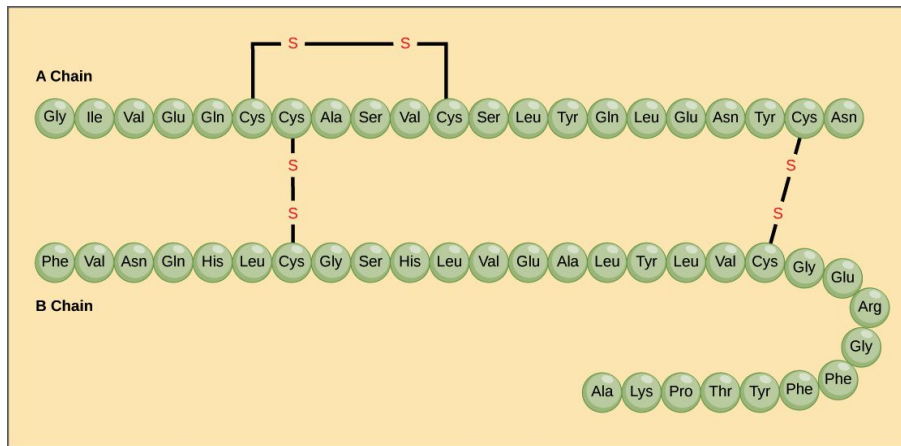B. Amino Acids with Polar Uncharged Side Chains
C. Special Cases
D. Amino Acids with Hydrophobic Side Chains

# What is a protein *sequence*?

- Since there are only 20 amino acids, we can give each amino acid an letter as an abstraction
  - Allows for representation as a string, which is more amenable to computation



```
>UniRef90_P02057 Hemoglobin subunit beta-1/2 n=6
Tax=Euarchontoglires TaxID=314146 RepID=HBB_RABIT
MVHLSSEEKSAVTALWGKVNVEEVGGEALGRLLVVYPWTQRFFESFGDLSSANAVMNNPK
VKAHGKKVLAAFSEGLSHLDNLKGTFAKLSELHCDKLHVDPENFRLLGNVLVIVLSHHFG
KEFTPQVQAAYQKVVAGVANALAHKYH
```

# What is a protein *structure*?

- **Secondary structure:** Interactions between atoms cause the series of amino acids to form regular substructures
- **Tertiary structure:** secondary structures are assembled into folds in 3D space



Primary structure
amino acid sequence

Secondary structure
regular sub-structures

alpha helix

beta sheet

Tertiary structure
three-dimensional structure

P13 protein

Quaternary structure
complex of protein molecules
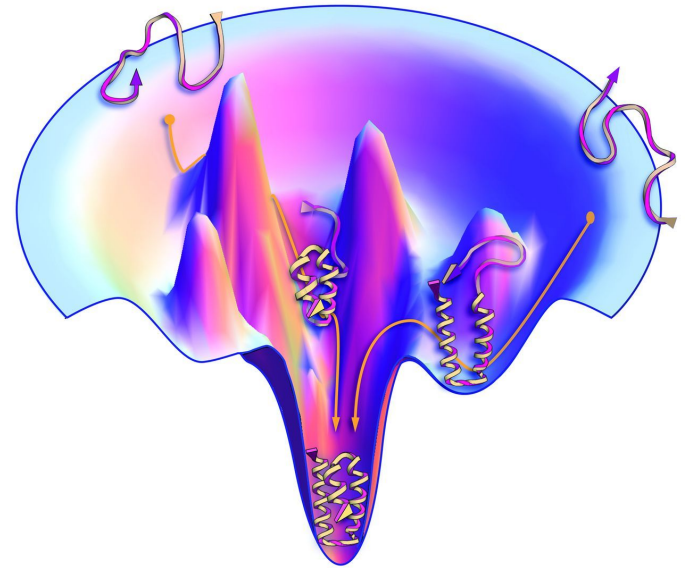
hemoglobin

Image source

# *Backbone* structure vs. *all-atom* structure



- Protein backbone atoms: (N, C$\alpha$, C)
- Side chain groups (denoted generically as R) determine the amino acid identity
  - Analogy: "protrudes" from the backbone like t-shirts on this clothesline, where color of the t-shirts defines the sequence.
- The **all-atom structure** include positions of the backbone and sidechain atoms.
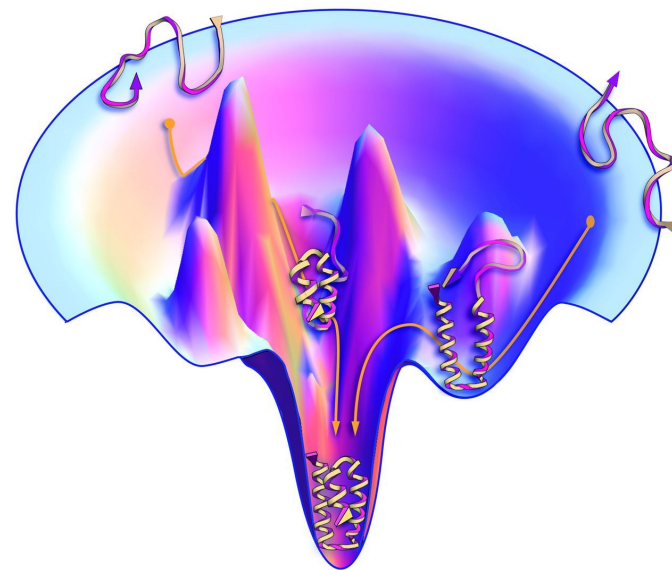
# The protein folding problem

- Proteins seek to "fold" into a structure that minimizes free energy
- **Anfinsen's dogma:** 3D structure of (most) natural proteins is determined only by its amino acid sequence.
- Protein folding problem: predicting 3D structure from 1D sequence.
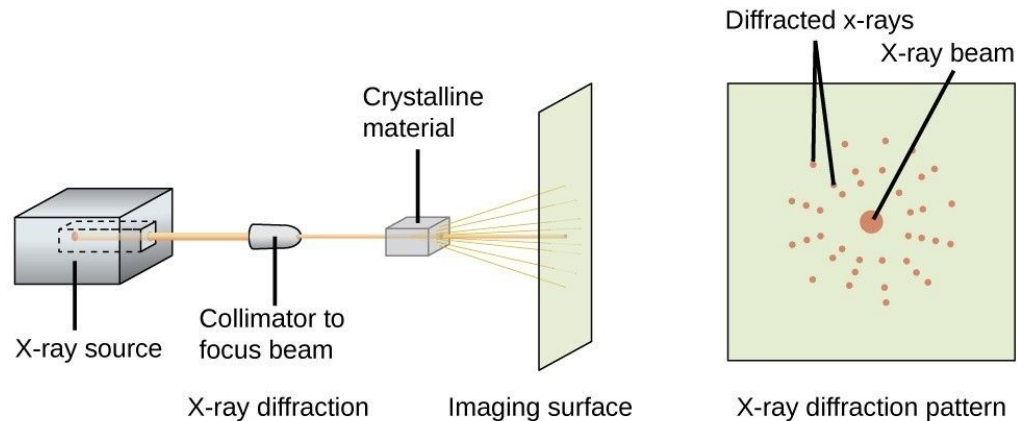
# The protein folding problem

- **Levinthal's paradox:** searching through possible structural configurations for a protein should take longer than the age of the universe, yet proteins can fold in seconds.

Determining structure from sequence alone should theoretically be possible from our understanding of nature, yet very difficult from our current understanding of biophysics.
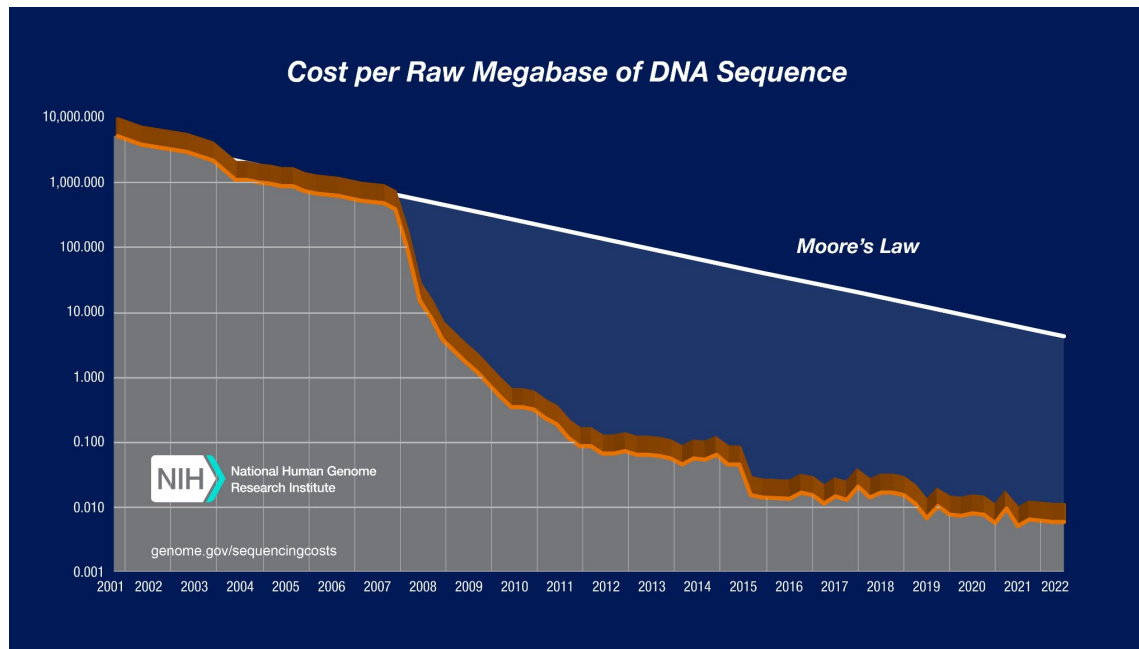
# Experimentally–resolved structure is expensive & slow

- Structures typically **experimentally resolved** by X-ray crystallography, cryo-EM, etc.
- Can take up to months and years to resolve a single structure



Source: https://university.pressbooks.pub/chemistryucf/chapter/10-8-x-ray-crystallography/
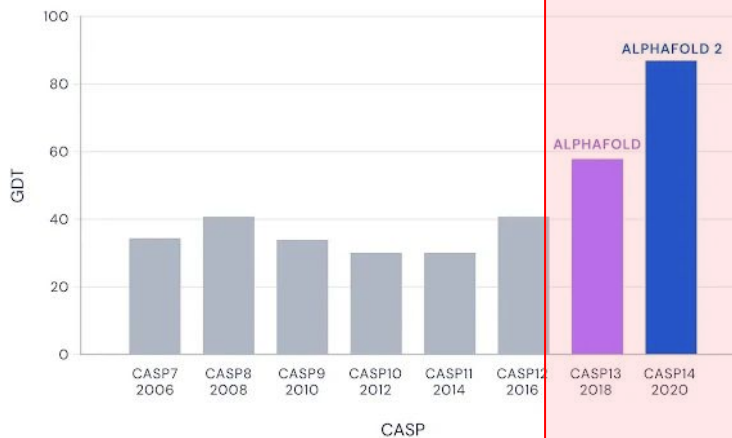
# Sequence data is more abundant than structure



Source: https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data

# Benchmarks drive progress: the CASP competition

- CASP: Critical Assessment of Structure Prediction (CASP)
  - Biennial competition with structural data for the held out set generated as teams prepare their predictions.
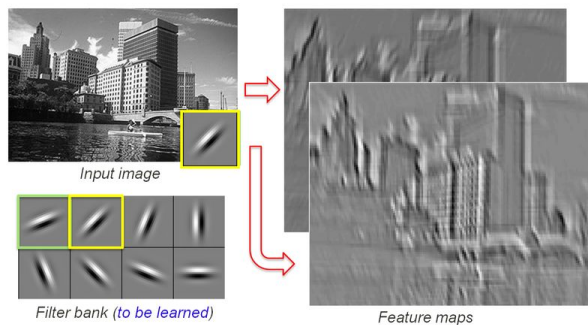- Traditional approaches typically leverage biophysics for predictions

**Median Free-Modelling Accuracy**

*AlphaFold mechanism:*
# Overview

*tl;dr: balancing provision of evolutionary and graph-like priors with effective end-to-end learning*

# Inductive biases in conventional architectures



CNNs: adopts ideas from signal processing to learn composite lower level features

Transformers: attention designed to capture local structure (grammar) and global structure (content coherence)

What are the inductive biases to capture for protein structure prediction?

# Designing an architecture for biological inductive biases



harness additional sequence-based priors

learn structural features from latents

generate structures

# AlphaFold mechanism: Training Data



*tl;dr: Labelled structures, unlabelled sequences, and self-distillation*

# The Protein Data Bank (PDB)

- Publicly available dataset of experimentally-resolved protein structures.

# Incorporating labelled and unlabelled data

- Consider pairs of {sequence, structure} data as {input, label}.
  - As of 2024:
    - # of {sequence, structure} pairs in PDB: **218,196**
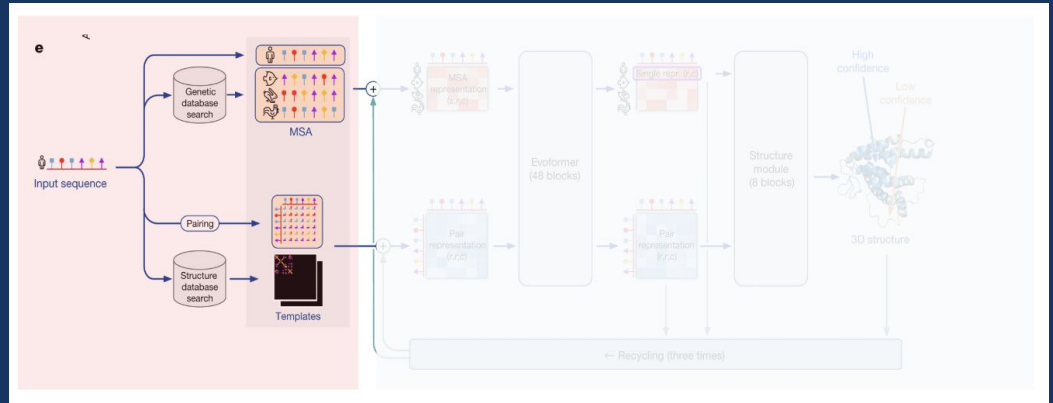    - # of sequence-only data in UniRef50: **63,849,054**
  - Relatively small labelled dataset, but a large number of unlabelled sequence data
- Aim – improve structure prediction using unlabelled sequence data:
  - 1. Capture co-evolutionary patterns from sequence
  - 2. Self-distillation – use labelled sequences to predict labels for unlabelled sequences

# Coevolution



If two positions are in <u>contact</u>, then if one is mutated, the other also loses function.

Thus, statistically observed pairs of residues that "change together" across evolution often denotes contact.

Image credit: Hetu Kamisetty

# Multiple sequence alignment (MSA) input

- Multiple sequence alignment:
  - Statistically line up similar string positions across sequences
  - Used to infer evolutionary relationships
- If a statistical sequence pattern is preserved across species, it's likely to imply function
  - Based on principles of Darwinian evolution

# A natural language analogy...

**Text 1**   **Pad Thai** is stir-fry dish made with <u>rice noodles</u>, <u>shrimp</u>, <u>chicken</u> (or <u>tofu</u>), <u>peanuts</u>, all sautéed together in a wok and tossed in a delicious <u>sauce</u>.

**Text 2**   …a delicious <u>Pad Thai</u> recipe. You'll need <u>rice noodles, shrimp, chicken, tofu, peanuts</u>, scrambled eggs, red peppers, and bean sprouts.

**Text 3**   This delicious thai dish is traditionally made up of <u>rice noodles</u> tossed with a deeply flavored, sweet and sour <u>sauce</u>. The <u>sauce</u> ...

**Text 4**   Our favorite vegan <u>Pad Thai</u> recipe involves stir-frying <u>rice noodles</u>, <u>tofu</u>, and <u>peanut sauce</u> in a wok, with bean sprouts and vegetables added to taste. You'll love…

A pad thai recipe that doesn't include `rice noodles` would likely be removed by the editor, because it would no longer be pad thai! Thus, `rice noodles` is observed to be conserved across recipes.

# Template Input

- For all sequence homologs found during MSA construction: if an existing structure exists in the PDB, use it as input.
  - In practice, few structures have available templates at inference



**b** Pair representation (r,r,c)

Corresponding edges in a graph

# Noisy Student Training

**Self-training with Noisy Student improves ImageNet classification**

Qizhe Xie, Minh-Thang Luong, Eduard Hovy, Quoc V. Le



1. Train "undistilled" model on PDB dataset only
2. Predict structure for all available sequences
3. Retrain network with distillation dataset

During student training:
- 75% from self-distillation dataset
- 25% from PDB

# AlphaFold mechanism: Evoformer



*tl;dr: Criss-cross attention to evolutionary details*
*+ triangulate attention for 3D awareness*

# Breaking down the Evoformer



Axial attention

Triangular self-attention

a

48 blocks (no shared weights)

MSA representation (s,r,c)

Row-wise gated self-attention with pair bias

Column-wise gated self-attention

Tran-sition

Outer product mean
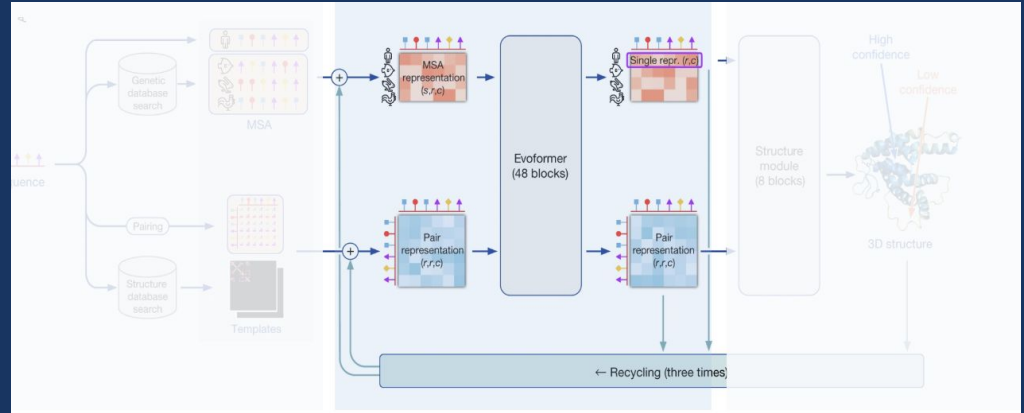
Pair representation (r,r,c)

Triangle update using outgoing edges

Triangle update using incoming edges

Triangle self-attention around starting node

Triangle self-attention around ending node

Tran-sition

MSA representation (s,r,c)

Pair representation (r,r,c)

48 stacked blocks

# Review: transformers and self-attention

Scaled Dot-Product Attention

Multi-Head Attention



high attention

low attention

She is eating a green apple.

$$\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$$

Attention weights learn relationships
between sequence positions
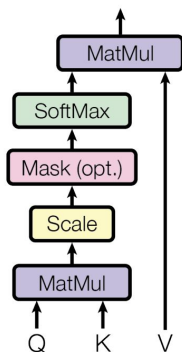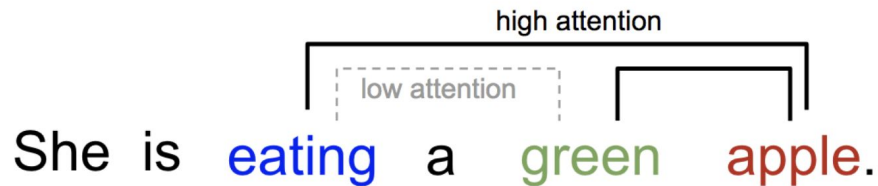
# Axial attention: encoding evolution



[Submitted on 28 Nov 2018 (v1), last revised 9 Jul 2020 (this version, v2)]

**CCNet: Criss-Cross Attention for Semantic Segmentation**

Zilong Huang, Xinggang Wang, Yunchao Wei, Lichao Huang, Humphrey Shi, Wenyu Liu, Thomas S. Huang

**MSA Transformer**

Roshan Rao, Jason Liu, Robert Verkuil, Joshua Meier, John F. Canny, Pieter Abbeel, Tom Sercu, Alexander Rives

Original attention block

Axial transformer block

# Axial attention: encoding evolution



Sequence A

115    120    125

F S T A A F R F G H A V H P L V R R L

Query for similar sequences

Evolutionary relationships

Sequence A
Sequence B
Sequence C
Sequence D
Sequence E
Sequence F

Original attention block

Axial transformer block

# Axial attention: encoding evolution



- **Row-wise attention** captures intra-sequence relationships
- **Column-wise attention** captures evolutionary relationships

# Triangular self-attention

# Triangular self-attention: encoding graph structure

- Pairwise representation captures relationships between <u>residue positions</u>
    - I.e. captures positions in close contact

# Triangular self-attention: encoding graph structure

- **Pairwise representation captures relationships between <u>residue positions</u>**
  - I.e. captures positions in close contact
- **Graph structures represent these edge relationships well**
  - Aim: update transformer learning to implicitly capture graph relationships



**b** Pair representation (r,r,c)

Corresponding edges in a graph

# Triangular self-attention: enforcing 3D common sense

- If pairwise edges are to represent distances in 3D space, they should obey the **triangle inequality**
  - Note: the inequality refers to <u>distances between coordinates</u> rather than the 3D coordinate positions



$$z = x+y$$

$$\|z\| = \|x+y\| < \|x\|+\|y\|$$

# Triangular self-attention: enforcing 3D common sense

- If pairwise edges are to represent distances in 3D space, they should obey the triangle inequality
- Output values by the naive attention mechanism may violate this
  - need to implement a soft constraint

3 + 6 > 8

3 + 6 < 10

# Triangular self-attention: multiplicative update



- To update pairwise representation at {i, j}, take the outer product of {i, k} and {j, k} for all k, and update with new outer product



$$\mathbf{g}_{ij} = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right)$$

$$\tilde{\mathbf{z}}_{ij} = \mathbf{g}_{ij} \odot \text{Linear}(\text{LayerNorm}(\textstyle\sum_k \mathbf{a}_{ik} \odot \mathbf{b}_{jk}))$$

# Triangular self-attention: attention update

- To update attention value of pairwise representation at {i, j} ("starting node"), update by the attention value at {i,k} for all k.
  - Additionally: bias the query-key multiplication {i, j} by value at {i, k}



Triangle self-attention around starting node

# Triangular self-attention: attention update

- To update attention value of pairwise representation at {i, j} ("starting node"), update by the attention value at {i,k} for all k.
  - Additionally: bias the query-key multiplication {i, j} by value at {i, k}

**Algorithm 13** Triangular gated self-attention around starting node

$\textbf{def}$ TriangleAttentionStartingNode($\{\mathbf{z}_{ij}\}, c = 32, N_{\text{head}} = 4$) :

*# Input projections*

1: $\mathbf{z}_{ij} \leftarrow \text{LayerNorm}(\mathbf{z}_{ij})$

2: $\mathbf{q}_{ij}^h, \mathbf{k}_{ij}^h, \mathbf{v}_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$

3: $b_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$

4: $\mathbf{g}_{ij}^h = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right)$

*# Attention*

5: $a_{ijk}^h = \text{softmax}_k\left(\frac{1}{\sqrt{c}}\, \mathbf{q}_{ij}^{h\top} \mathbf{k}_{ik}^h + b_{jk}^h\right)$

6: $\mathbf{o}_{ij}^h = \mathbf{g}_{ij}^h \odot \sum_k a_{ijk}^h \mathbf{v}_{ik}^h$

*# Output projection*

7: $\tilde{\mathbf{z}}_{ij} = \text{Linear}\left(\text{concat}_h(\mathbf{o}_{ij}^h)\right)$

8: $\textbf{return}\ \ \{\tilde{\mathbf{z}}_{ij}\}$

Berkeley
UNIVERSITY OF CALIFORNIA

# Triangular self-attention: attention update

- Repeat for "ending nodes":
  - i.e. having just updated attention value at {i, j} by all {i, k} ∀ k, we now update attention value at **{j, i}** by attention values at **{j, k}** ∀ k.



Triangle self-attention around starting node

Triangle self-attention around ending node

# Integrating graph and evolutionary representations



Bias from pairwise representations added to the corresponding {i, j} position in row attention, followed by outer product mean

**Row attention with bias**

$$a_{sij}^h = \text{softmax}_j \left( \frac{1}{\sqrt{c}} \, \mathbf{q}_{si}^{h^\top} \mathbf{k}_{sj}^h + b_{ij}^h \right)$$

**Column attention**

$$a_{sti}^h = \text{softmax}_t \left( \frac{1}{\sqrt{c}} \, \mathbf{q}_{si}^{h^\top} \mathbf{k}_{ti}^h \right)$$

**Outer product mean**

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

# Evoformer: Summary

*AlphaFold mechanism:*
# Structure Module



*Tl;dr: simplify prediction by considering protein structures as rigid groups connected by "joints".*

# How should we represent the 3D structure?

**Idea 1:** As a point cloud of 3D coordinates

# How should we represent the 3D structure?

**Idea 1:** As a point cloud of 3D coordinates

→ issue: lacks 3D "equivariance"

Same object in different orientations have different representations.

We want same object in different orientations to the same representation

# How should we represent the 3D structure?

**Idea 1:** As a point cloud of 3D coordinates
→ issue: lacks 3D "equivariance"

Same object in different orientations have different representations.

We want same object in different orientations to the <u>same representation</u>

For any rotations applied to the input, we want our model to apply that same rotation to the output:



$g$: rotation
$f$: function (e.g. neural net)

# How should we represent the 3D structure?

**Idea 2:** Use equivariant neural network architectures from geometric deep learning.

→ Details are outside the scope of this lecture, but such methods are often complex and computationally expensive.



Source:
https://thegradient.pub/towards-geometric-deep-learning/

# How should we represent the 3D structure?

**Idea 3:** as "rigid blobs" with translation and rotation matrices from a reference frame

→ functions built on this representation is equivariant by construction
→ easily adaptable as input to transformer



Figure by: Nazim Bouatta

# AlphaFold representation: "Residue gas"



Figure by: Nazim Bouatta

- Represent each residue as a triangle with {N, C$\alpha$, C} as vertices.
- Backbone is a series of tiled triangles in 3D.

# Residue gas representation



Figure by: Nazim Bouatta

- "Chunk up" the structure, and consider each residue as "independently floating".
  - Network learns to place the backbone frames sequentially.
- All frames are initialized at the origin
  - Termed "black hole initialization"

# Representing residue positions as transformations

- Each residue is represented as a rigid frame transformed by a tuple of **rotation** ($R_i$) and **translation** ($t_i$) matrices from the origin:
  - Rotation: $R_i \in \mathbb{R}^{3 \times 3}$
  - Translation: $T_i \in \mathbb{R}^{3}$
- Obtained from 3D coordinates via the Gram-Schmidt process



Figure by: Nazim Bouatta

# Structure module: overview

Resolve side chain positions from predicted torsion angles

Backbone positions

**d**

Pair representation $(r,r,c)$

8 blocks (shared weights)

Predict χ angles and compute all atom positions

Single repr. $(r,c)$

IPA module

Single repr. $(r,c)$

Predict relative rotations and translations

Backbone frames $(r, 3×3)$ and $(r,3)$ (initially all at the origin)

Backbone frames $(r, 3×3)$ and $(r,3)$

Berkeley
UNIVERSITY OF CALIFORNIA

# Invariant Point Attention

Workhouse of the structure module: **Invariant Point Attention (IPA).**

→ A modified attention mechanism to act on $T_i = (R_i, t_i)$ tuples.

# Invariant Point Attention

From Evoformer

$$a_{ij}^h = \text{softmax}_j \left( w_L \left( \frac{1}{\sqrt{c}}\, \mathbf{q}_i^{h\top} \mathbf{k}_j^h \; + b_{ij}^h \right.\right.$$

Query/key over single
representation

Bias from
pairwise
representation

# Invariant Point Attention

From Evoformer

Bias by attending to frame transformations for residues $i$ and $j$

$$a_{ij}^h = \text{softmax}_j \left( w_L \left( \frac{1}{\sqrt{c}} \, \mathbf{q}_i^{h\top} \mathbf{k}_j^h + b_{ij}^h \right) - \frac{\gamma^h w_C}{2} \sum_p \left\| T_i \circ \vec{\mathbf{q}}_i^{hp} - T_j \circ \vec{\mathbf{k}}_j^{hp} \right\|^2 \right) \right)$$

Query/key over single representation

Bias from pairwise representation

**Note: L2-norm between Euclidean matrices is invariant by construction,** since the distance between matrices remains constant if the same rotation is applied to both!

# Backbone update

- Predict backbone frames from IPA module output

# Backbone update

- Rotation matrices are stored as quaternions
  - Rather than as a $\mathbb{R}^{3\times3}$ matrix, use a $\mathbb{R}^4$ to denote **rotation axis vector and associated angle**
    - Only two dimensions need to be provided for the axis vector to define a unique rotation
  - Reduces memory usage
  - First component is fixed to 1, and network predicts the remaining 3 components.



Rotation axis

Rotation angle

# Backbone update

**Algorithm 23** Backbone update

**def** BackboneUpdate($\mathbf{s}_i$) :

1: $b_i, c_i, d_i, \vec{\mathbf{t}}_i = \text{Linear}(\mathbf{s}_i)$         → Predict rotation (quaternion components) and translation vector

  # *Convert (non-unit) quaternion to rotation matrix.*

2: $(a_i, b_i, c_i, d_i) \leftarrow (1, b_i, c_i, d_i) / \sqrt{1 + b_i^2 + c_i^2, + d_i^2}$

3: $R_i = \begin{pmatrix} a_i^2 + b_i^2 - c_i^2 - d_i^2 & 2b_i c_i - 2a_i d_i & 2b_i d_i + 2a_i c_i \\ 2b_i c_i + 2a_i d_i & a_i^2 - b_i^2 + c_i^2 - d_i^2 & 2c_i d_i - 2a_i b_i \\ 2b_i d_i - 2a_i c_i & 2c_i d_i + 2a_i b_i & a_i^2 - b_i^2 - c_i^2 + d_i^2 \end{pmatrix}$     → Convert back to $\mathbb{R}^{3 \times 3}$ representation

4: $T_i = (R_i, \vec{\mathbf{t}}_i)$

5: **return** $T_i$

# Placing the side chain atoms

- Predicting frame transformations only places the backbone atoms
- To place the side chains, we also predict the torsion angles.
  - Consider proteins as having "moveable joints", while other components are kept rigid.

# Predicting torsion angles

- Predict torsion angles $\alpha \in \mathbb{R}^2$ with a shallow ResNet
- 7 total types of rotations:
  - **Backbone: {ω,φ,ψ}**
  - **Sidechain: {$X_1$,$X_2$,$X_3$,$X_4$}**

# Predicting torsion angles

| 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 | $C^\gamma$ | $\boxed{O^{\delta 1}, O^{\delta 2}}$ | - | - |
| 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 | $C^\gamma$ | $C^\delta$ | $\boxed{O^{\epsilon 1}, O^{\epsilon 2}}$ | - |
| 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 | $C^\gamma$ | $\boxed{C^{\delta 1}, C^{\delta 2}, C^{\epsilon 1}, C^{\epsilon 2}, C^\zeta}$ | - | - |
| 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 | $C^\gamma$ | $\boxed{C^{\delta 1}, C^{\delta 2}, C^{\epsilon 1}, C^{\epsilon 2}, O^\eta, C^\zeta}$ | - | - |
| VAL | $N, C^\alpha, C, C^\beta$ | O | $C^{\gamma 1}, C^{\gamma 2}$ | - | - | - |

<u>Note</u>: some residues are symmetric and has two "correct" torsion angle values

- Each atom is grouped into a "rigid group" based on their dependence on a specific torsion angle
- Number of χ angles present is dependent on the amino acid identity.

# From frames to coordinates

- Map rigid body coordinates onto the backbone frames with the predicted torsion angles to obtain final atom positions
  - Coordinates of rigid bodies in reference frame are defined as constants
- Molecular dynamics (non-deep learning) used on final coordinates to relax explicit clashes



Torsion angles

Backbone frames

# Predicted LDDT (pLDDT)

- lDDT: local distance difference test
  - A per-atom measure of local accuracy
- Train model to predict lDDT-Cα using ground truth structure
  - I.e. pLDDT is a measure of predicted confidence
  - One value per residue
- Experimentally shown to correspond well to experimental accuracy

**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 with low pLDDT may be unstructured in isolation.

# Other outputs

- Distogram
  - "Distance" + "histogram"
  - Bin pairwise contacts into 64 discrete bins
- Experimentally resolved prediction
  - Per-atom prediction with values between (0, 1)

# Structure module: summary

**Predict torsion angles** with separate ResNet to place side chains

Process coordinates into **backbone frames** of rotation & translation from origin

$$T_i := (R_i, \vec{t_i})$$

**Final output** of **atom positions** (placed with idealized bond lengths & predicted angles) and auxiliary outputs (e.g. **pLDDT**)

**Modified attention mechanism** that attends to relationships between frame representations.



d

Pair representation $(r,r,c)$

8 blocks (shared weights)

Single repr. $(r,c)$

IPA module

Predict $\chi$ angles and compute all atom positions

Single repr. $(r,c)$

Predict relative rotations and translations

Backbone frames $(r, 3\times3)$ and $(r,3)$ (initially all at the origin)

Backbone frames $(r, 3\times3)$ and $(r,3)$

Berkeley
UNIVERSITY OF CALIFORNIA

*AlphaFold mechanism:*
# Loss Functions

*Tl;dr: L2 distance on frames + losses to enforce contributions from all network components*

# Assessing structure prediction quality

- Desiderata for a loss function:
  - Differentiable
  - Computationally fast to calculate
- Idea 1: distance root-mean squared distance (dRMSD)
  - Take pairwise distances and calculate root-mean-squared distance
- However, it does not account for chirality
  - Recall: mirror image of a protein might have different biological functions



Bimodal errors for AlphaFold trained with the chirality-unaware dRMSD loss

# Frame Adjusted Point Error (FAPE)

- Leverage the transformations associated with predicted coordinates to project back to the reference frame:

$$\vec{\mathbf{x}}_{ij} = T_i^{-1} \circ \vec{\mathbf{x}}_j$$

$$\vec{\mathbf{x}}_{ij}^{true} = T_i^{true-1} \circ \vec{\mathbf{x}}_j^{true}$$

"Undo" the predicted from-reference-frame transformation

$$d_{ij} = \sqrt{\|\vec{\mathbf{x}}_{ij} - \vec{\mathbf{x}}_{ij}^{true}\|^2 + \epsilon}$$

Calculate L2 distance for coordinates in reference frame orientation

Berkeley
UNIVERSITY OF CALIFORNIA

# Training loss functions

$$\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$$

Training losses are designed to emphasize
contributions from different architectural parts:

- $L_{\text{aux}}$: averaged FAPE and torsion angle losses
  from intermediate structure module layers
- $L_{\text{dist}}$: distogram prediction (cross-entropy)
- $L_{\text{msa}}$: masked-MSA prediction (cross-entropy)

# Training loss functions

$$\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$$

Training losses are designed to emphasize contributions from different architectural parts:

- $L_{\text{aux}}$: averaged **FAPE and torsion angle losses** from intermediate structure module layers
- $L_{\text{dist}}$: distogram prediction (cross-entropy)
- $L_{\text{msa}}$: masked-MSA prediction (cross-entropy)

- enforces contributions from structure module
- Torsion angles calculated as unit-circle normalized L2-loss

Berkeley
UNIVERSITY OF CALIFORNIA

# Training loss functions

$$\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$$

Training losses are designed to emphasize contributions from different architectural parts:

- $L_{\text{aux}}$: averaged **FAPE and torsion angle losses** from intermediate structure module layers
- $L_{\text{dist}}$: distogram prediction (cross-entropy)
- $L_{\text{msa}}$: masked-MSA prediction (cross-entropy)

- Calculates cross-entropy between real and predicted distogram
- Enforces accuracy of pairwise relationships

Berkeley
UNIVERSITY OF CALIFORNIA

# Training loss functions

$$\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$$

Training losses are designed to emphasize contributions from different architectural parts:
- $L_{\text{aux}}$: averaged **FAPE and torsion angle losses** from intermediate structure module layers
- $L_{\text{dist}}$: distogram prediction (cross-entropy)
- $L_{\text{msa}}$: masked-MSA prediction (cross-entropy)

- BERT-like masked language modelling loss
- Enforces contributions from the MSA trunk
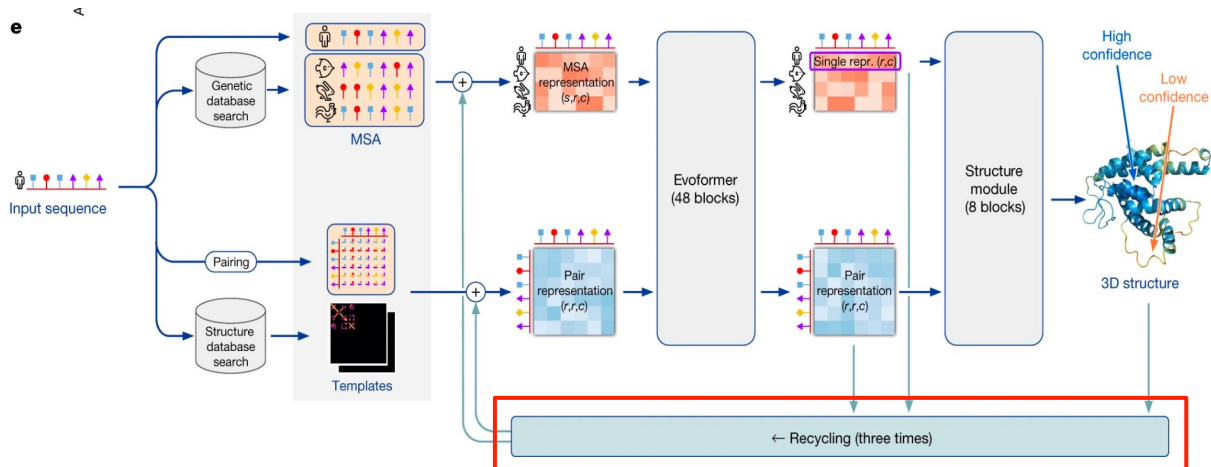
Berkeley
UNIVERSITY OF CALIFORNIA

# Finetuning loss functions

$$\mathcal{L} = \begin{cases} 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} & \text{training} \\ 0.5\mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 2.0\mathcal{L}_{\text{msa}} + 0.01\mathcal{L}_{\text{conf}} + 0.01\mathcal{L}_{\text{exp resolved}} + 1.0\mathcal{L}_{\text{viol}} & \text{fine-tuning} \end{cases}$$
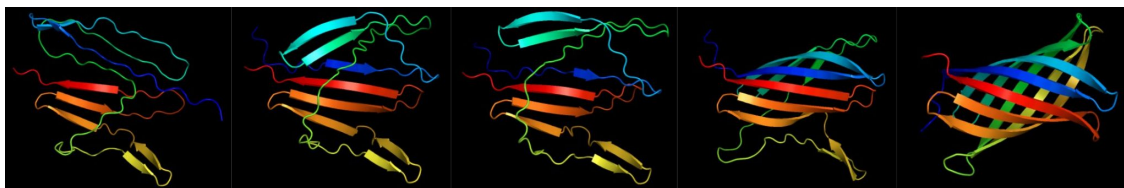
Some are only used during finetuning:

- $L_{\text{exp resolved}}$: finetunes the experimentally-resolved output with explicit labels from high-resolution structures
- $L_{\text{viol}}$: distogram prediction (cross-entropy)

# Recycling



- Output single and pair representations are reused as inputs to the Evoformer!
- Unlike stacked blocks, **weights are shared** across iterations
  - Uses stop-gradient after each iteration.

For some proteins, insufficient recycling cycles can greatly impact performance.



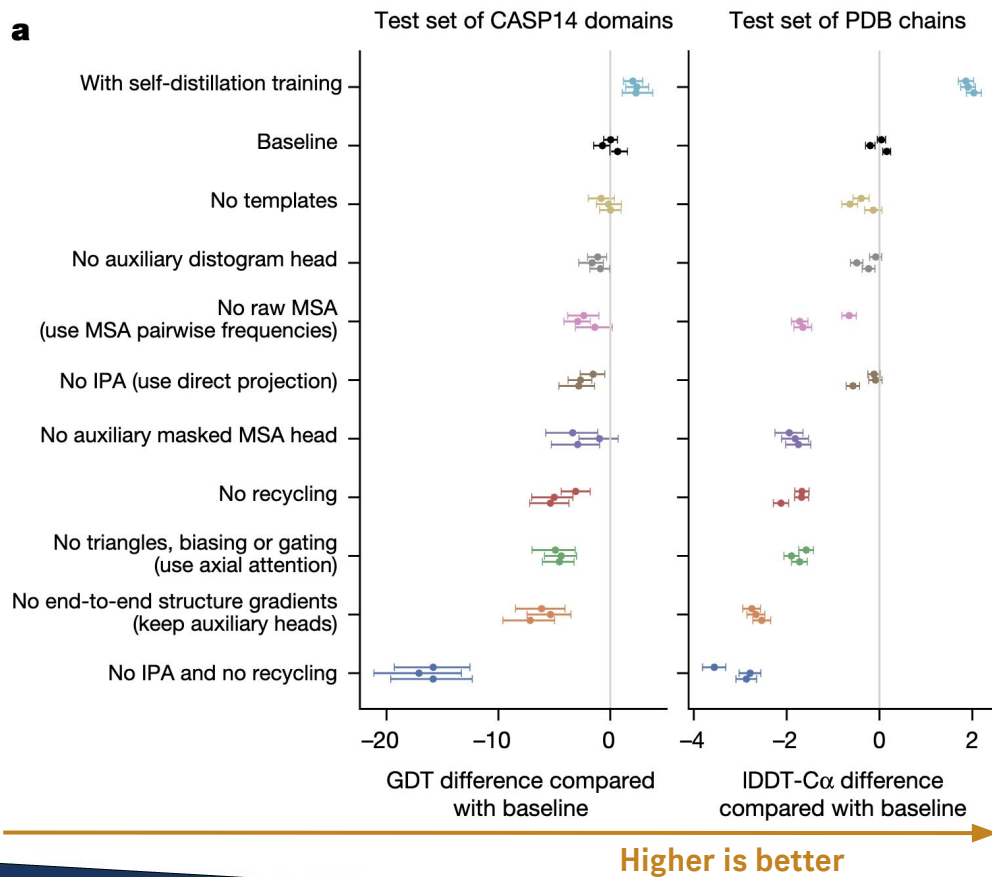recycle 1    recycle 2    recycle 3    recycle 4    recycle 5

Image source: Sergey Ovchinnikov

# Ablations

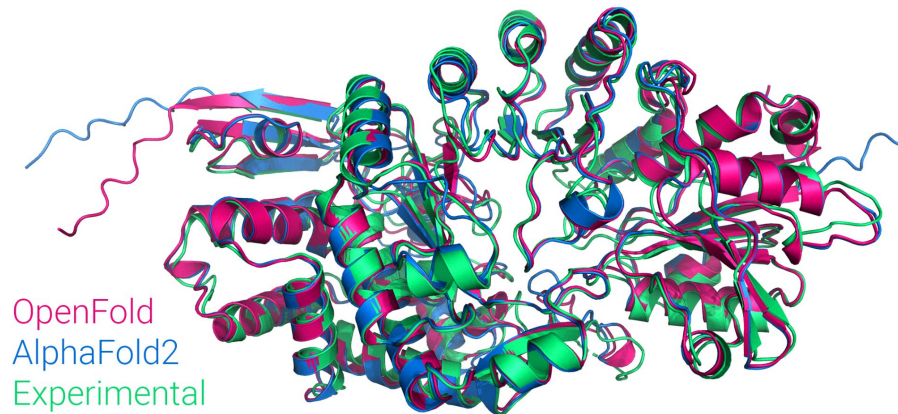- "Knockdown study" to see which network components are most important

# Limitations and Open Problems

*tl;dr: generalization failures; moving towards dynamical & heterogeneous systems.*
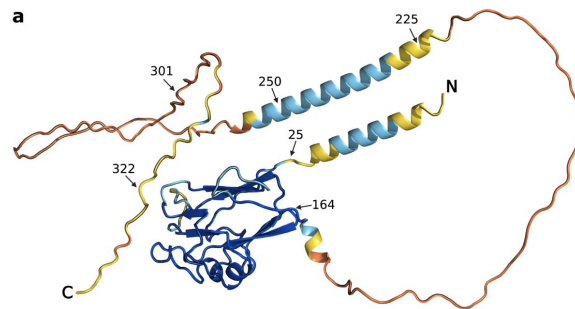
# Alternative methods

- RoseTTAFold
  - Methodologically similar to AlphaFold2, released after CASP14
- ESMFold and OmegaFold:
  - Replaces MSA input with a protein language model
- OpenFold
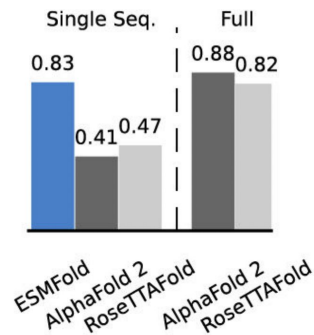  - PyTorch and open-sourced reimplementation of AlphaFold



OpenFold
AlphaFold2
Experimental

Source: https://github.com/aqlaboratory/openfold/

# Limitations: Generalization failures

- PDB biases towards crystallizable proteins with stable structure
  - Performance degrades on antibody loops and intrinsically disordered regions
- MSA trunk relies on the existence of homologs
  - Performance degrades for proteins without homologs
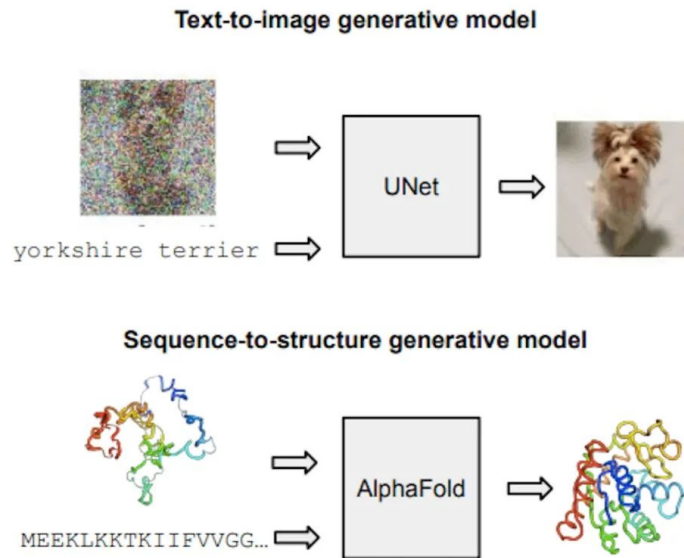  - Non-MSA methods such as OmegaFold and ESMFold performs better for these use cases.



Source: Piovesan et al.



Source: Lin et al.

# Limitations: Conformationally static

- Proteins are dynamic!
- Because AlphaFold is trained on crystallized structures, it is only capable of predicting static structures
- Recent works finetune AlphaFold as a generative model to predict a distribution over structures instead
  - Conceptually recapitulates the Boltzmann distribution



Source: Jing et al.

# Limitations: Modelling Complexes

- Multimers
  - AlphaFold2 at CASP14 was unable to model multi-chain protein assemblies (i.e. quaternary structure)
  - Community hacks arounds include adding a glycine linker to "glue" multiple chains into one sequence
  - AlphaFold-Multimer model released in 2022.
- Nucleic acids
  - Many important protein functions occur in tandem with nucleic acids, for e.g. protein-RNA interactions, gene editing enzymes, etc.
  - Recent work such as RoseTTAFold All-Atom (Krishna, 2024) make progress towards predicting structure for these complexes