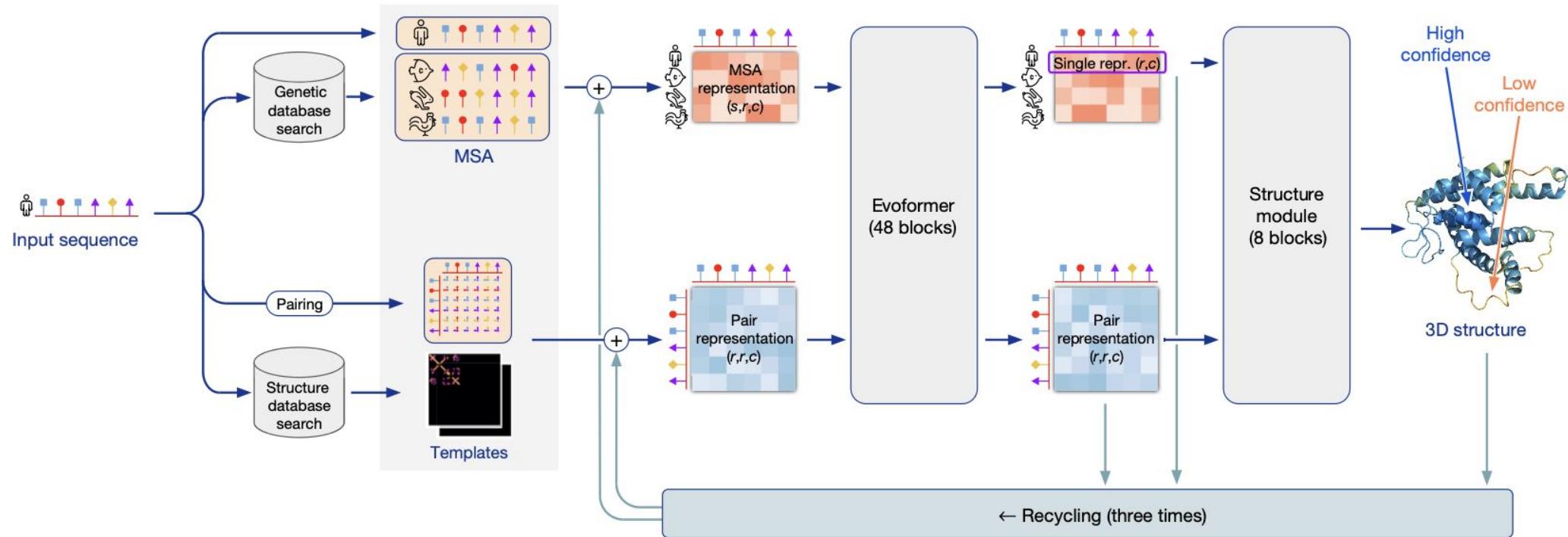


# Protein Structure Prediction with Language Models

Shentong Mo

Dec 8, 2022

# Recap AlphaFold2



- Multiple Sequence Alignment (MSA) + Templates of Similar Protein Structures
- Evoformer
- Structure Module

# Follow-up Papers

- **ESMFold** (Meta AI, Alexander Rives)
  - *Lin et al*, Language models of protein sequences at the scale of evolution enable accurate structure prediction
- **OmegaFold** (Helixon, Jian Peng)
  - *Wu et al*, High-resolution *de novo* structure prediction from primary sequence
- **HelixFold** (Baidu & Biomap, Le Song)
  - *Fang et al*, HelixFold-Single: MSA-free Protein Structure Prediction by Using Protein Language Model as an Alternative

# Presentation Line

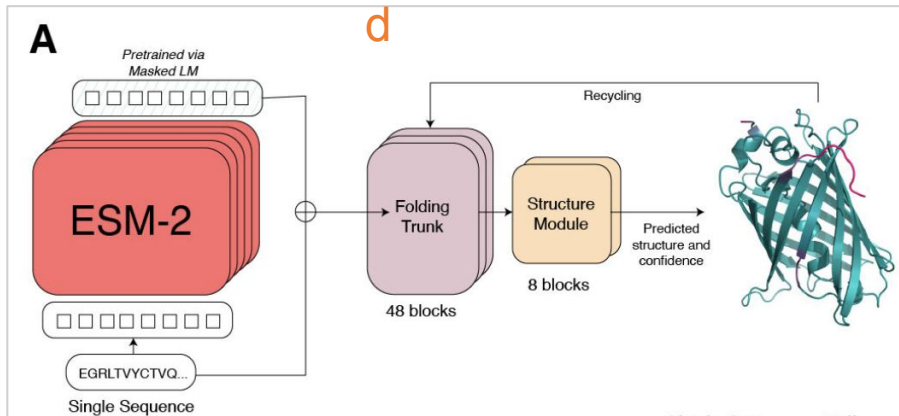
- PLM architecture
- PLM training methods
- PLM training data
- Model size
- Evoformer counterpart
- Structure finetuning data

# Comparisons

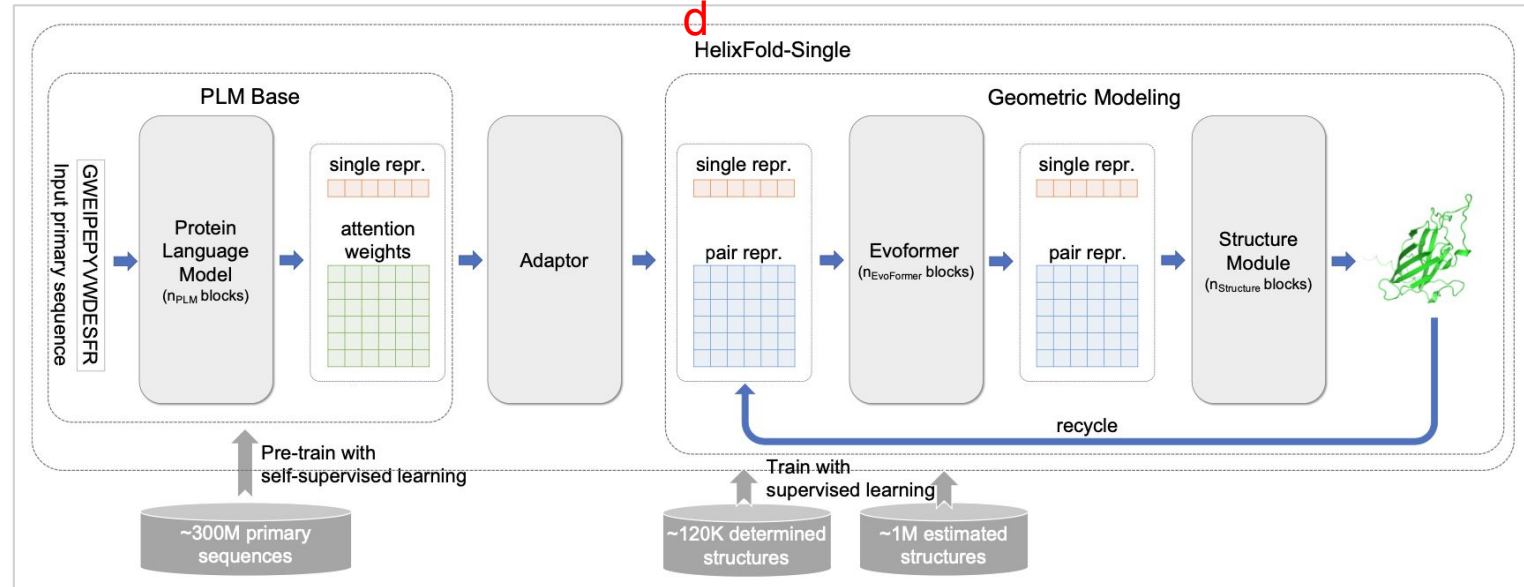
Method	PLM architectures	PLM training data	masked pretraining objective	Evoformer Counterpart	Structure fine-tuning data	Structure validation data	Freeze PLM or not
ESMFold	ESM-2 (15B)	UniRef50 (60M)	residue	ESM Folding Block	PDB + Alphafold2 Structure	CAMEO CASP14	Freeze
OmegaFold	OmegaPLM (670M)	UniRef50 (60M)	residue, motif, subsequence	Geoformer (Geometric Modeling)	PDB	CAMEO CASP13, 14	Freeze
HelixFold	PLM (1B) + Adapter	UniRef30 (260M)	residue	Revised Evoformer (Geometric Modeling)	PDB+ Uniclust30+ Alphafold2 Structure	CAMEO CASP14 MSA Depth Test	Not freeze

# Overall Architecture

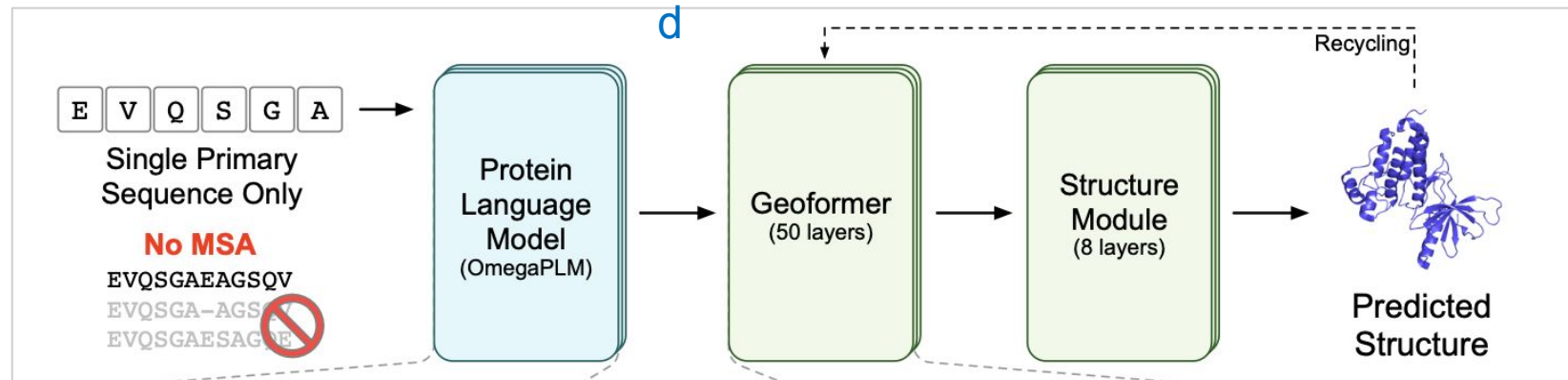
## ESMFol<sup>d</sup>



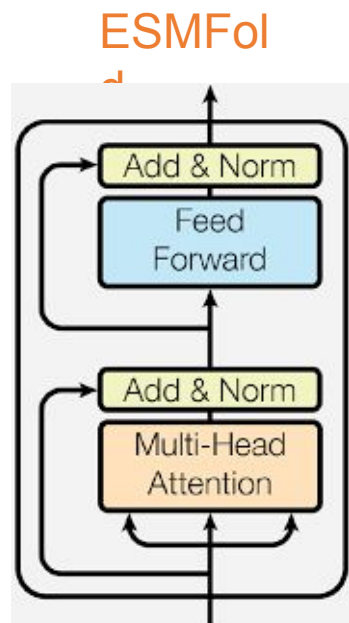
## HelixFol<sup>d</sup>



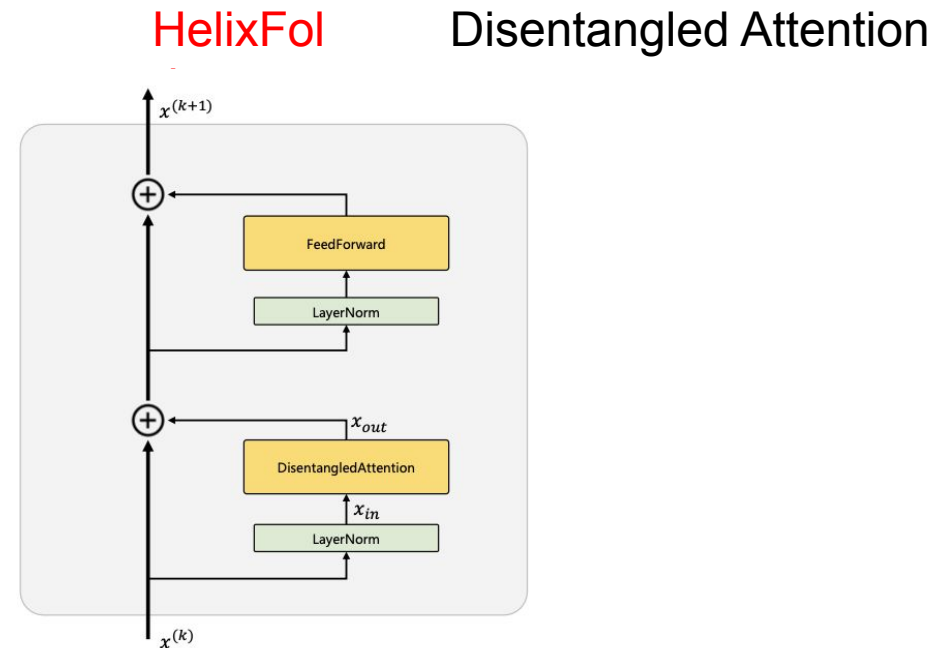
## OmegaFol<sup>d</sup>



# PLM Architecture



Standard Attention



OmegaFol Gated Attention Unit

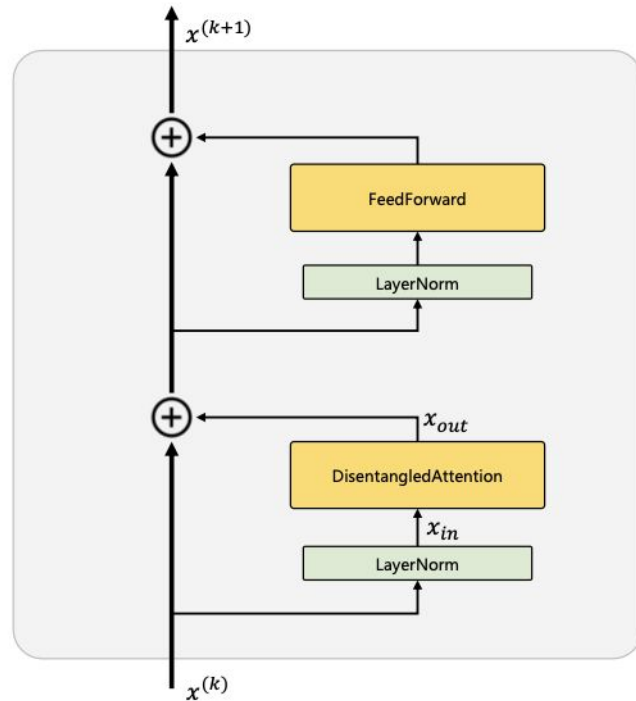
**Algorithm 1:** Protein language models based on the Gated Attention Module (GAU)

```

1 def OmegaPLM (  $\{n_i\}$ ,  $d_k=256$ ,  $d=1,280$ ,  $N_{stack}=66$ ,  $d_v=2,560$  ) :
2   for  $l \in [1, \dots, N_{stack}]$  do
3      $r_i = \text{LayerNorm}(n_i)$ 
4      $u_i, v_i, g_i = \text{SiLU}(\text{Linear}(r_i))$ 
5      $\{q_i\} = \text{RoPE}(\{w_q \odot u_i + b_q\})$ 
6      $\{k_i\} = \text{RoPE}(\{w_k \odot u_i + b_k\})$ 
7      $\alpha_{ij} = \text{softmax}_j \left( \frac{\log n}{\sqrt{d_k}} (q_i^T k_j) + b_{i-j} \right)$ 
8      $o_i = g_i \odot \sum_j \alpha_{ij} v_j$ 
9      $n_i += \text{Linear}(o_i)$ 
10  end
11 return  $\{n_i\}$ 

```

# Disentangled Attention Transformer



$$q = x_{in} W_q, \quad k = x_{in} W_k, \quad v = x_{in} W_v, \quad p = e_p W_p$$

$$A_{i,j} = \underbrace{q_i k_j^T}_{\text{(a) residue-to-residue}} + \underbrace{q_i p_{\delta(i,j)}^T}_{\text{(b) residue-to-position}}$$

$$x_{out} = \text{softmax}\left(\frac{A}{\sqrt{2d_{PLM}}}\right)v$$



# Standard vs Disentangled Attention

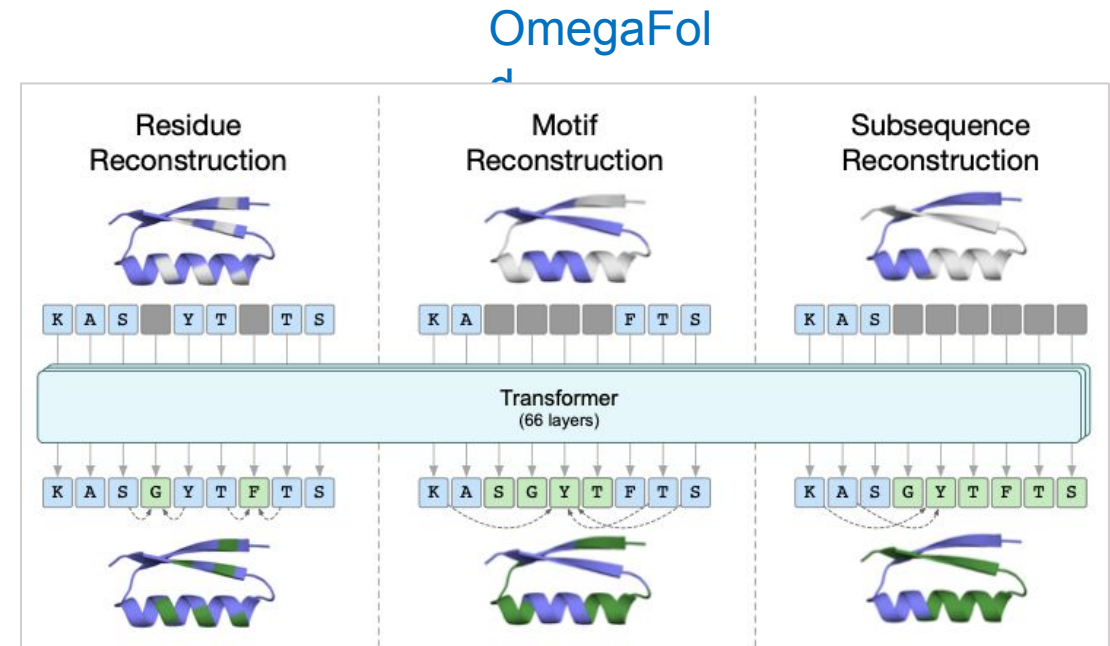
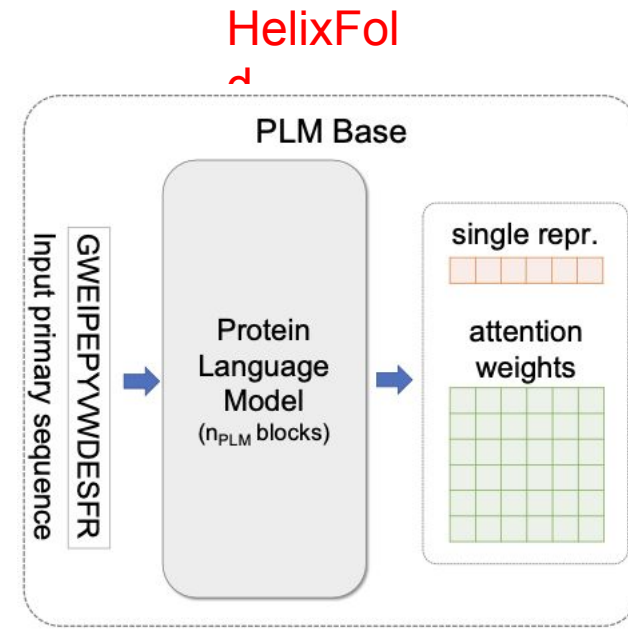
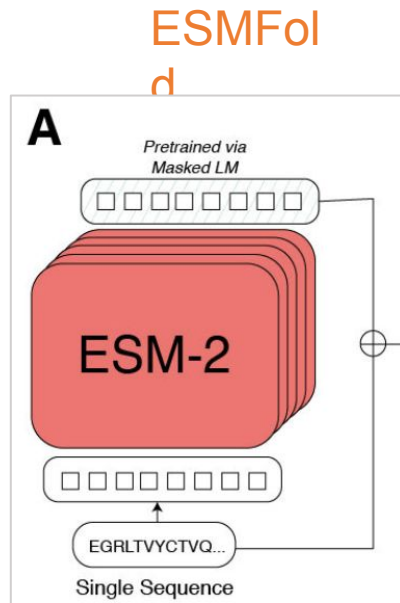
Standard  
Attention

$$Q = HW_q, K = HW_k, V = HW_v, A = \frac{QK^\top}{\sqrt{d}}$$
$$H_o = \text{softmax}(A)V$$

Disentangled  
Attention

$$Q_c = HW_{q,c}, K_c = HW_{k,c}, V_c = HW_{v,c}, Q_r = PW_{q,r}, K_r = PW_{k,r}$$
$$\tilde{A}_{i,j} = \underbrace{Q_i^c K_j^{c\top}}_{\text{(a) content-to-content}} + \underbrace{Q_i^c K_{\delta(i,j)}^r}^{\top}_{\text{(b) content-to-position}} + \underbrace{K_j^c Q_{\delta(j,i)}^r}^{\top}_{\text{(c) position-to-content}}$$
$$H_o = \text{softmax}\left(\frac{\tilde{A}}{\sqrt{3d}}\right)V_c$$

# PLM training objectives



# PLM model size

## ESMFol

	8M	35M	150M	650M	3B	15B
Dataset	UR50/D	UR50/D	UR50/D	UR50/D	UR50/D	UR50/D
Number of layers	6	12	30	33	36	48
Embedding dim	320	480	640	1280	2560	5120
Attention heads	20	20	20	20	40	40
Training steps	500K	500K	500K	500K	500K	270K
Learning rate	4e-4	4e-4	4e-4	4e-4	4e-4	1.6e-4
Weight decay	0.01	0.01	0.01	0.01	0.01	0.1
Clip norm	0	0	0	0	1.0	1.0
Distributed backend	DDP	DDP	DDP	DDP	FSDP	FSDP

## OmegaFol

No. Layers	66
$d$	1280
$d_k$	256
$d_v$	2560
No. Attention Head	1
Tying Embeddings for input & prediction head	True
Cosine normalization with learned scale (28, 36) at output	True
Clipping thresholds of relative positional bias	[-64, 64]
Normalization type	LayerNorm (37)
Pre- or Post-Norm	Pre-Norm
Learnable parameters in normalization	False

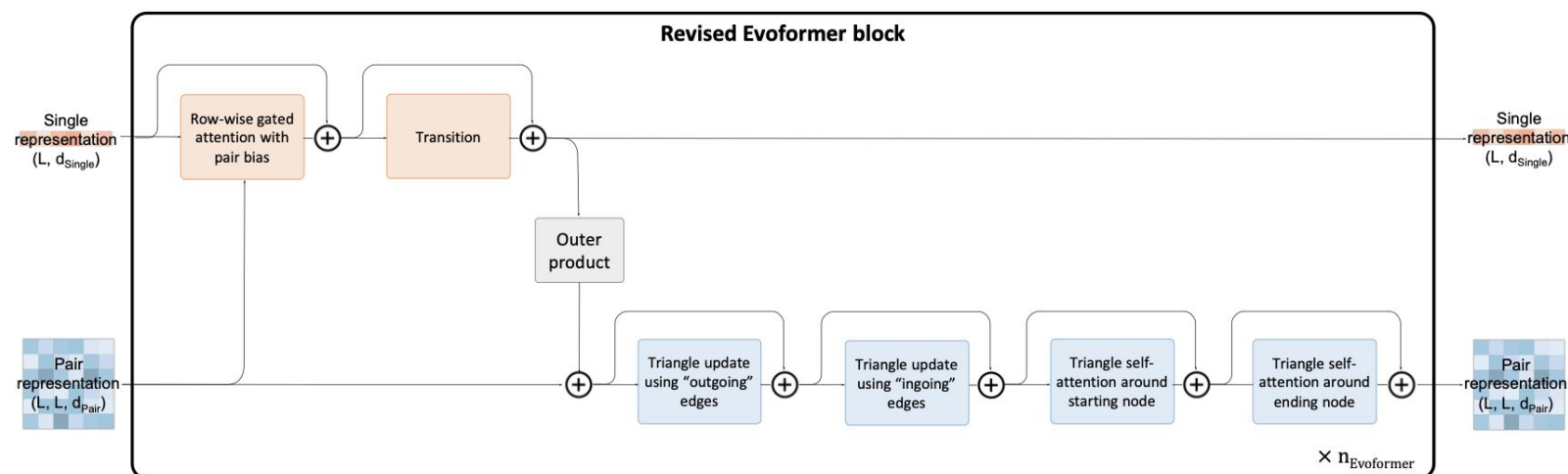
## HelixFol

Components	Model size	Layer num	Hidden size	Intermediate size	Head num
PLM-1B	1.09B	$n_{PLM} = 20$	$d_{PLM} = 2048$	8192	$h_{PLM} = 16$
PLM-100M	100M	$n_{PLM} = 12$	$d_{PLM} = 768$	3072	$h_{PLM} = 12$

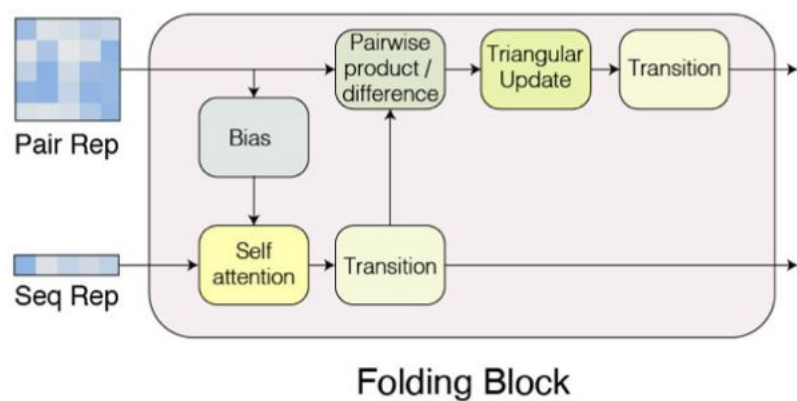
- *Disentangled Attention Transformer*

# Evoformer Counterpart

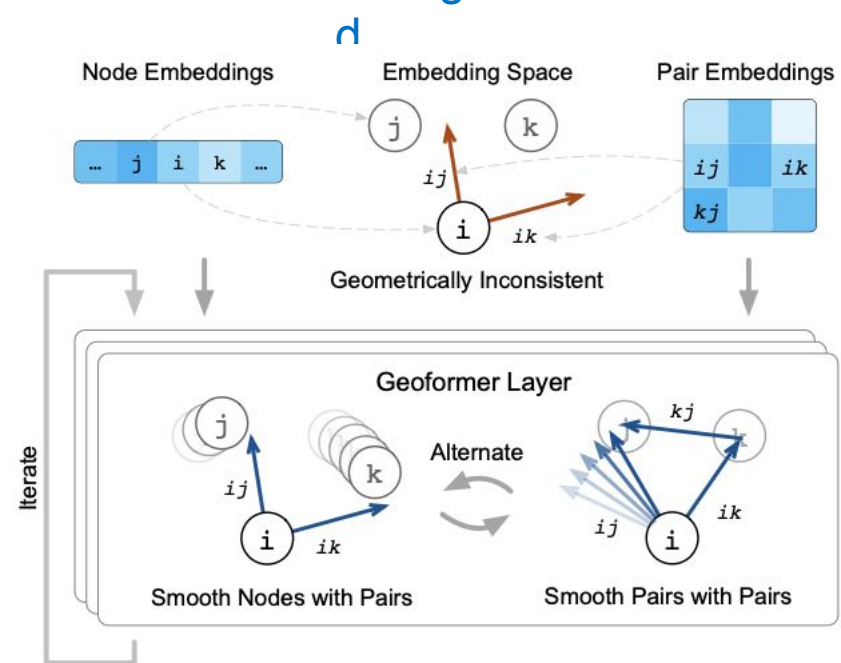
HelixFol



ESMFol



OmegaFol



# Block Algorithm

ESMFol

```
Algorithm 1:  
FoldingBlock(s, z)  
b = Linear(z)  
s = s + MultiHeadSelfAttention(s, bias=b)  
s = s + MLP(s)  
z = z + Linear(Concat([OuterProduct(s), OuterDifference(s)]))  
z = z + TriangularMultiplicativeUpdateOutgoing(z)  
z = z + TriangularMultiplicativeUpdateIncoming(z)  
z = z + TriangularSelfAttentionOutgoing(z)  
z = z + TriangularSelfAttentionIncoming(z)  
z = z + MLP(z)  
return s, z
```

HelixFol

d

No much details about the algorithm,  
but similar to ESMFold block,  
removing

column-wise attention in AlphaFold2

OmegaFol

---

## Algorithm 3: Geometric Transformer

---

```
1 def Geoformer (  $A_{aa(i)}$ ,  $\{\mathbf{n}_i\}$ ,  $\{\mathbf{w}_{ij}\}$ ,  $N_1 = 50$ ,  $N_2 = 8$ ,  $d_n = 256$ ,  $d_w = 128$ ) :  
2   for  $l \in [1, \dots, N_1]$  do  
3      $\{\mathbf{n}_i\} \mathrel{+}= \text{NodeAttention}(\{\mathbf{n}_i\}, \{\mathbf{w}_{ij}\})$   
4      $\{\mathbf{n}_i\} \mathrel{+}= \text{NodeTransition}(\{\mathbf{n}_i\})$   
5      $\{\mathbf{w}_{ij}\} \mathrel{+}= \text{Node2Edge}(\{\mathbf{n}_i\})$   
6     for  $k \in [1, 2]$  do  
7        $\{\mathbf{w}_{ij}\} \mathrel{+}= \text{EdgeAttention}(\{\mathbf{w}_{ij}\})$   
8     end  
9      $\{\mathbf{w}_{ij}\} \mathrel{+}= \text{EdgeTransition}(\{\mathbf{w}_{ij}\})$   
10  end  
11  for  $l \in [1, \dots, N_2]$  do  
12     $\{\mathbf{n}_i\}, \{\vec{x}_i\} \mathrel{+}= \text{StructureModule}(\{\mathbf{n}_i\}, \{\mathbf{w}_{ij}\})$   
13     $\{\mathbf{w}_{ij}\} \mathrel{+}= \text{3Dprojection}(A_{aa(i)}, \{\vec{x}_i\})$   
14     $\{\mathbf{w}_{ij}\} \mathrel{+}= \text{EdgeAttention}(\{\mathbf{w}_{ij}\})$   
15  end  
16 return  $\{\mathbf{n}_i\}, \{\mathbf{w}_{ij}\}$ 
```

---

# Geometric Transformer

## Node Attention

$$\mathbf{q}_i, \mathbf{k}_i, \mathbf{v}_i = \text{Linear}(\mathbf{n}_i^{(\ell-1)})$$

$$\mathbf{b}_{ij} = \text{Linear}(\mathbf{w}_{ij})$$

$$\alpha_{ij} = \text{softmax}_j \left( \frac{1}{\sqrt{c}} \mathbf{q}_i^T \mathbf{k}_j + \mathbf{b}_{ij} \right)$$

$$\mathbf{o}_i = \text{sigmoid}(\text{Linear}(\mathbf{n}_i^{\ell-1})) \odot \sum_j \alpha_{ij} \mathbf{v}_j$$

$$\mathbf{n}_i^\ell = \text{Linear}(\mathbf{o}_i)$$

## Edge Attention

$$\tilde{\mathbf{a}}_{ij}, \tilde{\mathbf{b}}_{ij} = \text{sigmoid}(\text{Linear}(\mathbf{w}_{ij})) \odot \text{Linear}(\mathbf{w}_{ij})$$

$$\mathbf{q}_{ij}, \mathbf{k}_{ij}, \mathbf{v}_{ij}, \mathbf{b}_{ij} = \text{Linear}(\mathbf{w}_{ij})$$

$$\mathbf{g}_{ij} = \text{sigmoid}(\text{Linear}(\mathbf{w}_{ij}))$$

$$\alpha_{itj} = \text{softmax}_t \left( \frac{1}{\sqrt{c}} \mathbf{q}_{ij}^T (\mathbf{k}_{it} + \mathbf{k}_{tj}) + \mathbf{b}_{it} + \mathbf{b}_{tj} \right)$$

$$\mathbf{o}_{ij} = \mathbf{g}_{ij} \odot \left( \sum_k \alpha_{itj} (\mathbf{v}_{it} + \mathbf{v}_{tj} + \tilde{\mathbf{a}}_{ti} \odot \tilde{\mathbf{b}}_{tj}) \right)$$

$$\mathbf{w}_{ij} = \text{Linear}(\mathbf{o}_{ij})$$

# Geometric Transformer (cont.)

---

**Algorithm 4: 3Dprojection**

---

```
1 def 3Dprojection(  $A_{aa(i)} \in \mathbb{R}^L, \{\vec{x}_i\} \in \mathbb{R}^{L \times 14 \times 3}, n = 3$  ) :  
2   for  $i, j \in [1, L]$  do  
3     // Embed all atom pair distances:  
4      $d_{ij} = \|\vec{x}_i - \vec{x}_j\|$   
5     // Embed all atom frames  
6      $f_{aa(i)}^{(n)} = \text{Frame}_n(\vec{x}_i)$   
7     // Embed distances in both rough and fine bins  
8      $\mathbf{a}_{ij} = \text{Linear}(\text{OneHot}(d_{ij}, \text{bins} = [3.25\text{\AA}, 20.75\text{\AA}, 16]))$   
9      $\mathbf{b}_{ij} = \text{Linear}(\text{OneHot}(d_{ij}, \text{bins} = [2.3125\text{\AA}, 21.6875\text{\AA}, 64]))$   
10    // Embed frame-position directed distance in Euclidean  
11    space:  
12     $\mathbf{c}_{ij}^{(n)} = \text{OneHot}(\text{R}(\vec{x}_i, f_{aa(j)}^{(n)}), \text{bins} = [-16\text{\AA}, 16\text{\AA}, 64], \text{space} = \text{Euclidean})$   
13    // Embed frame-position angle distance in spherical space:  
14     $\mathbf{e}_{ij}^{(n)} = \text{OneHot}(\text{R}(\vec{x}_i, f_{aa(j)}^{(n)}), \text{bins}_\phi = [0, 2\pi, 12], \text{bins}_\psi = [0, 2\pi, 12], \text{space} = \text{Spherical})$   
15    // Embed pair amino acid types:  
16     $\mathbf{g}_{ij} = \text{Linear}(\mathbf{a}_{ij} + \mathbf{b}_{ij})$   
17     $\mathbf{h}_{ij} = \text{Linear}(\text{concat}_n(\{\mathbf{c}_{ij}^{(n)}\}) + \text{concat}_n(\{\mathbf{e}_{ij}^{(n)}\}))$   
18     $\mathbf{z}_{ij} = A_{aa(i)} + A_{aa(j)}$   
19     $\mathbf{w}_{ij} = \text{Linear}(\mathbf{z}_{ij} \otimes (\mathbf{g}_{ij} + \mathbf{h}_{ij}))$   
20  end  
21 return  $\{\mathbf{w}_{ij}\}$ 
```

---



# Structure training loss

- AlphaFold2:

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

- ESMFold:

$$\mathcal{L} = \mathcal{L}_{\text{FAPE}} + \mathcal{L}_{\text{dist}}$$

- OmegaFold:

$$\mathcal{L} = \begin{cases} \mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 0.01\mathcal{L}_{\text{conf}} & \text{first stage} \\ \mathcal{L}_{\text{FAPE}} + 0.5\mathcal{L}_{\text{aux}} + 0.3\mathcal{L}_{\text{dist}} + 0.01\mathcal{L}_{\text{conf}} + 1.0\mathcal{L}_{\text{viol}} & \text{second and last stage} \end{cases}$$

- HelixFold: removing msa loss in AlphaFold2



# Structure training data

- ESMFold: (60M)

- **RCSB PDB**: all PDB chains until 2020-05-01 with resolution greater than or equal to 9Å and length greater than 20, cluster resulting in sequences at 40% sequence identity
- **AlphaFold2 Structure Database**: predicted IDDT greater than 70%

- OmegaFold: (number not given)

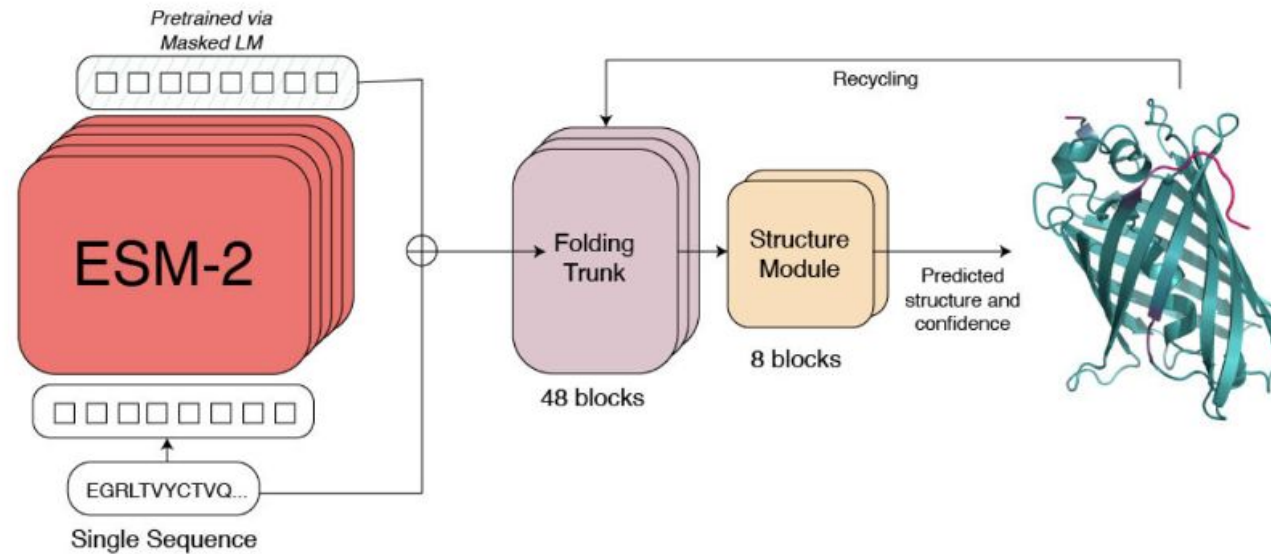
- **RCSB PDB**: all PDB chains, cluster resulting in sequences at 40% sequence identity

- HelixFold: (260M)

- **RCSB PDB**: all PDB chains until 2020-05-14 with resolution greater than or equal to 9Å and length greater than 20, and number of amino acids larger than 10, cluster resulting in sequences at 40% sequence identity
- **Distillation-Uniclust30**: Uniclust30 structure predicted by AlphaFold2
- **Distillation-EBI**: predicted IDDT greater than 50%

Thanks!

# Model Structure



- ESM-2 (Replace MSA and templates in AlphaFold2)
- Folding Block (Change evoformer block in AlphaFold2)
- Structure Module

# ESM-2

- Architecture: BERT-style encoder with transformers
- Modifications:
  - number of layers
  - number of attention heads
  - hidden size
  - feed-forward hidden size
  - sinusoidal positional encoding > learnable positional embedding
  - RoPE: rotary positional embedding (good for small models, bad for large models)
- Training Objective: unsupervised contact prediction with logistic regression

# Unsupervised Contact Prediction

- Let  $c_{ij}$  be a boolean random variable which is true if amino acids  $i, j$  are in contact
- Suppose our transformer has  $L$  layers and  $K$  attention heads per layer.
- Let  $A_{kl}$  be the symmetrized and Average Product Correction (APC)-corrected attention map for the  $k$ -th attention head in the  $l$ -th layer of the transformer,
- and  $\alpha_{ij}^{kl}$  be the value of that attention map at position  $i, j$ .

$$p(c_{ij}) = (1 + \exp(-\beta_0 - \sum_{l=1}^L \sum_{k=1}^K \beta_{kl} \alpha_{ij}^{kl}))^{-1}$$

# Perplexity Estimation

- To measure a language model's uncertainty of a sequence and defined as the exponential of the negative log-likelihood of the sequence
- the perplexity over a large dataset (non-deterministic)

$$\text{Perplexity}(x) = \exp\left\{-\log p(x_{i \in M} | x_{j \notin M} \cup \hat{x}_{i \in M})\right\}$$

where the mask  $M$  be a random variable denoting a set of tokens from input sequence  $x$

- the pseudo-perplexity over a single sequence (deterministic)

$$\text{PseudoPerplexity}(x) = \exp\left\{-\frac{1}{L} \sum_{i=1}^L \log p(x_i | x_{j \neq i})\right\}$$

where  $L$  is the length of the sequence

# ESM-2 Parameters

	<b>8M</b>	<b>35M</b>	<b>150M</b>	<b>650M</b>	<b>3B</b>	<b>15B</b>
Dataset	UR50/D	UR50/D	UR50/D	UR50/D	UR50/D	UR50/D
Number of layers	6	12	30	33	36	48
Embedding dim	320	480	640	1280	2560	5120
Attention heads	20	20	20	20	40	40
Training steps	500K	500K	500K	500K	500K	270K
Learning rate	4e-4	4e-4	4e-4	4e-4	4e-4	1.6e-4
Weight decay	0.01	0.01	0.01	0.01	0.01	0.1
Clip norm	0	0	0	0	1.0	1.0
Distributed backend	DDP	DDP	DDP	DDP	FSDP	FSDP

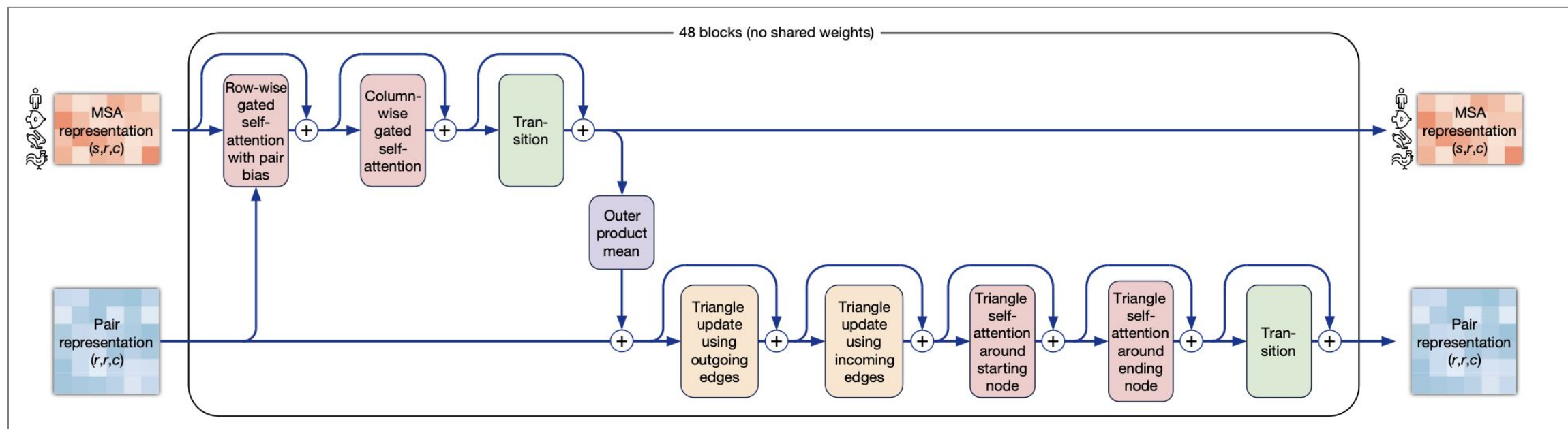
**Table S1: ESM-2 model parameters at different scales**

# Training data

- **Data source:** UniRef50 & UniRef90 (60M protein sequences for training, 250K for validation)
  - MMseqs search to remove all train sequences matching a validation sequence with 50% identity.
- **Filtering de-novo designed proteins:**
  - remove any sequence in UniRef50 and UniRef90 that was annotated as “artificial sequence” by a taxonomy search on the UniProt website
  - use jackhmmer to remove all hits around a manually curated set of 81 de-novo proteins
- **Amount and diversity:**
  - sampled a minibatch of UniRef50 sequences for each training update
  - replaced each sequence with a sequence sampled uniformly from the corresponding UniRef90 cluster

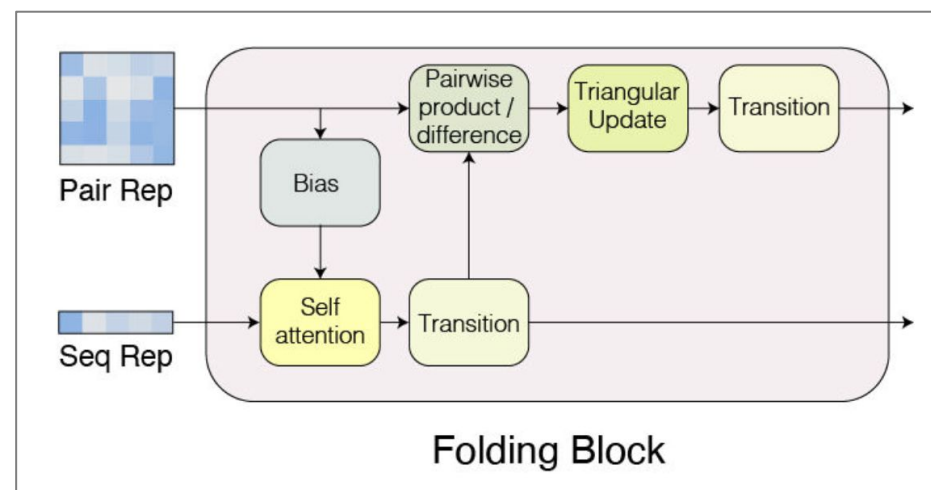


# Folding Block vs Evoformer Block

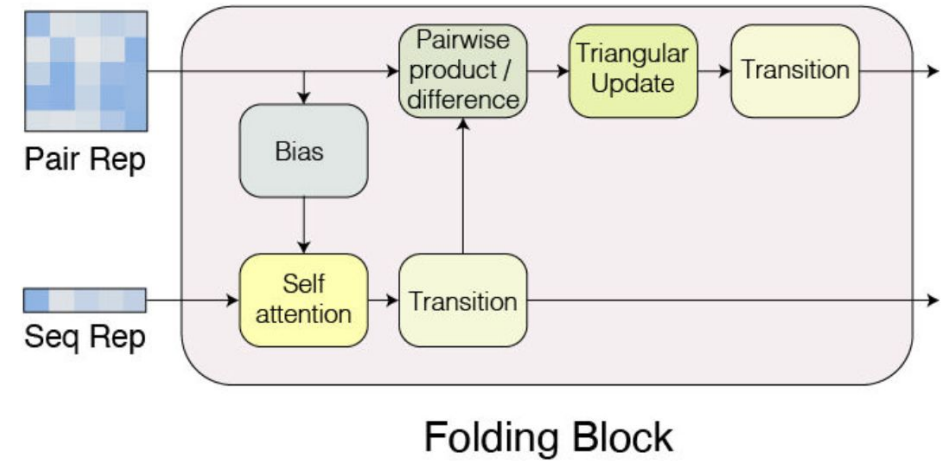


Evoformer  
Block

Folding  
Block

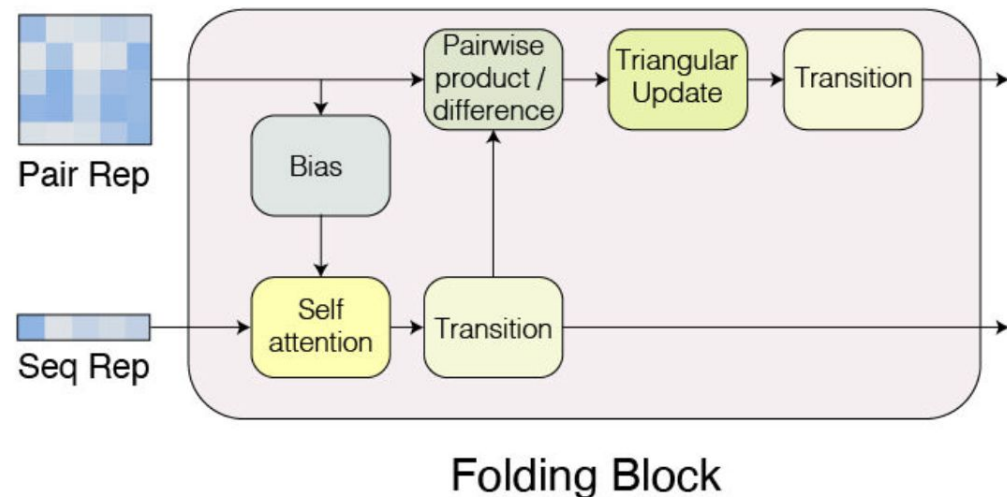


# Two main changes



- use **standard attention** over this feature space, as the language model features are one dimensional
  - *Evoformer block*: employs axial attention over the columns and rows of the MSA, as MSAs are two-dimensional.
- input the **attention maps** from the language model for structure information
  - *Evoformer block*: pass template information to the model as pairwise distances, input to the residue-pairwise embedding

# Folding Block Algorithm



Algorithm 1:

```
FoldingBlock(s, z)
```

```
b = Linear(z)
```

```
s = s + MultiHeadSelfAttention(s, bias=b)
```

```
s = s + MLP(s)
```

```
z = z + Linear(Concat([OuterProduct(s), OuterDifference(s)]))
```

```
z = z + TriangularMultiplicativeUpdateOutgoing(z)
```

```
z = z + TriangularMultiplicativeUpdateIncoming(z)
```

```
z = z + TriangularSelfAttentionOutgoing(z)
```

```
z = z + TriangularSelfAttentionIncoming(z)
```

```
z = z = MLP(z)
```

```
return s, z
```

# ESMFold Algorithm

Algorithm 2:

```
esm_c_s: number of channels in ESM hidden representation
c_s = 1024
c_z = 128
ESMFold(sequence)
s = ESM_hiddens(sequence) # num_layers x Length x esm_c_s
s = (softmax(layer_weights) * s).sum(0)
s = MLP(s)
z = PairwiseRelativePositionalEncoding(Length)
for b in folding_blocks:
    s, z = b(s, z)
return StructureModule(s, z)
```

# ESMFold output

- The IDDT head is output from the hidden representation of the StructureModule.
- The TM head uses the pairwise representation  $z$ .
- The distogram is predicted from the pairwise representation  $z$ .

# ESMFold training loss

- AlphaFold2:

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

- ESMFold:

$$\mathcal{L} = \mathcal{L}_{\text{FAPE}} + \mathcal{L}_{\text{dist}}$$

# Training structure data

- **Real Structure:**

- all PDB chains until 2020-05-01 with resolution greater than or equal to 9Å and length greater than 20
- cluster resulting in sequences at 40% sequence identity

- **Sampling:**

- sampling cluster evenly
- Rejection sampling to train longer proteins more frequently

- **Predicted Structure:**

- 13,477,259 structures predicted using AlphaFold2 on MSAs (predicted IDDT greater than 70)
- 75% predicted structures and 25% real structures during training

# Validation & Test structure data

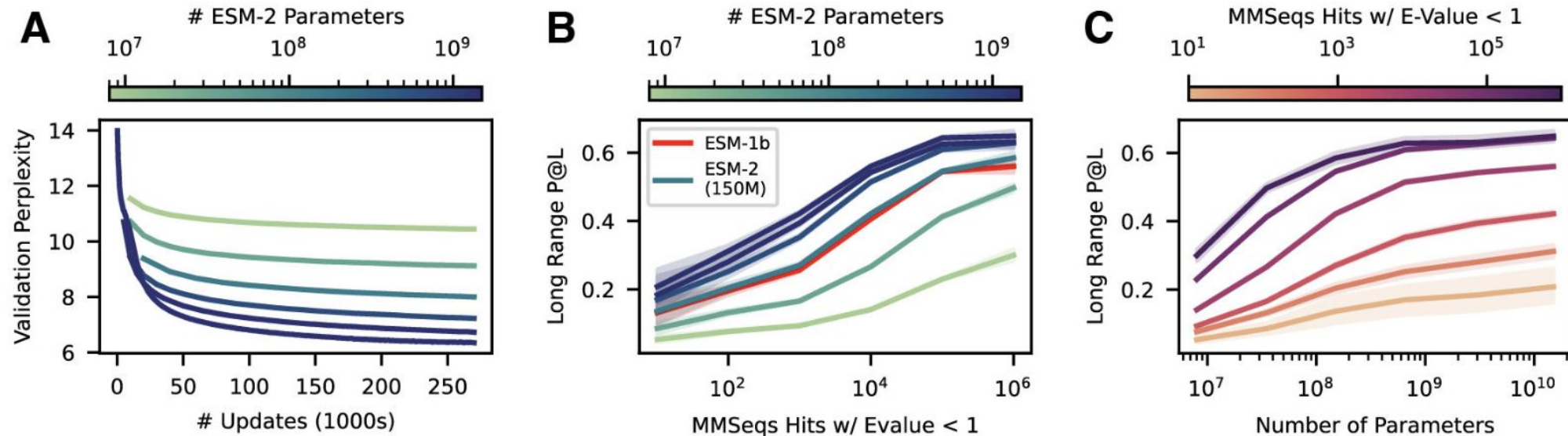
- Validation
  - Continuous Automated Model EvaluatiOn (CAMEO) (August 2021 to January 2022)
- Test:
  - CAMEO (194 test proteins from April 01, 2022 through June 25, 2022)
  - CASP14 competition (51 targets)
  - No filtering is performed on these test sets, even included length-2166 target T1044.



# Metrics

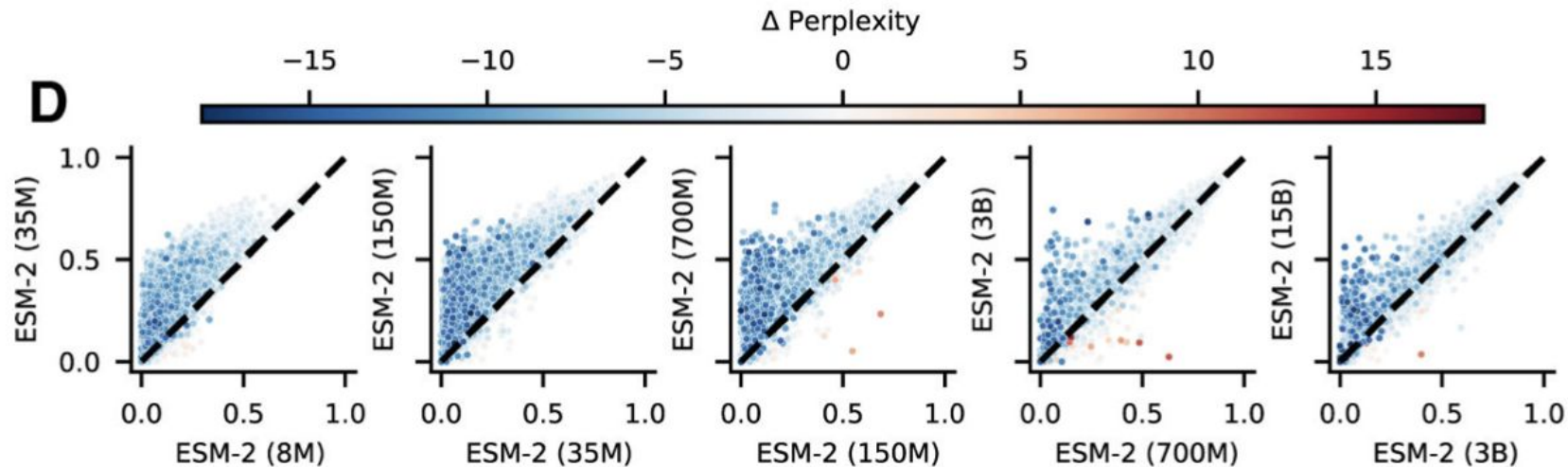
- **Validation Perplexity**: exponential of the negative log-likelihood over the validation set or a single sequence (lower is better)
- **P@L**: long-range precision @ L for unsupervised contact prediction performance (higher is better)
- **RMSD**: Root Mean Square Deviation (smaller is better)
- **TM-score**: Template Modeling score (higher is better)
- **pLDDT**: Model confidence prediction (higher is better)

# Scaling up to 15B parameters



- Larger models perform better at all levels
- 150M parameter ESM-2 model performs comparably with the 650M parameter ESM-1b model.
- The largest improvement is seen for sequences with  $O(10^4)$  MMseqs hits

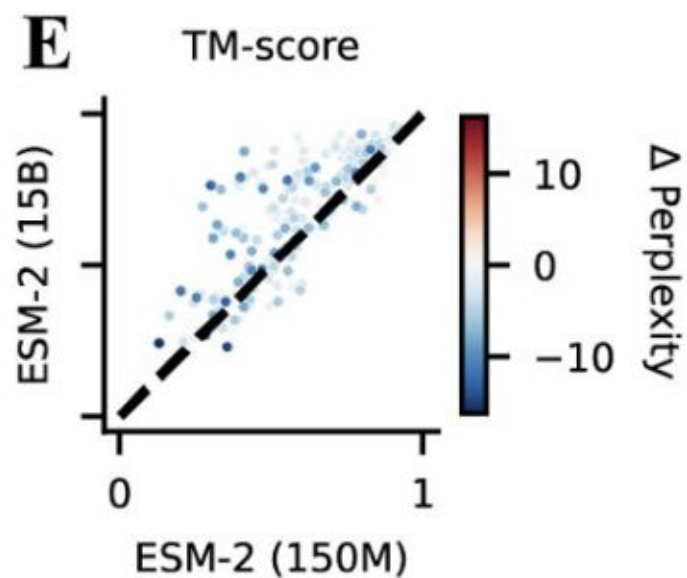
# Scaling up to 15B parameters (cont.)



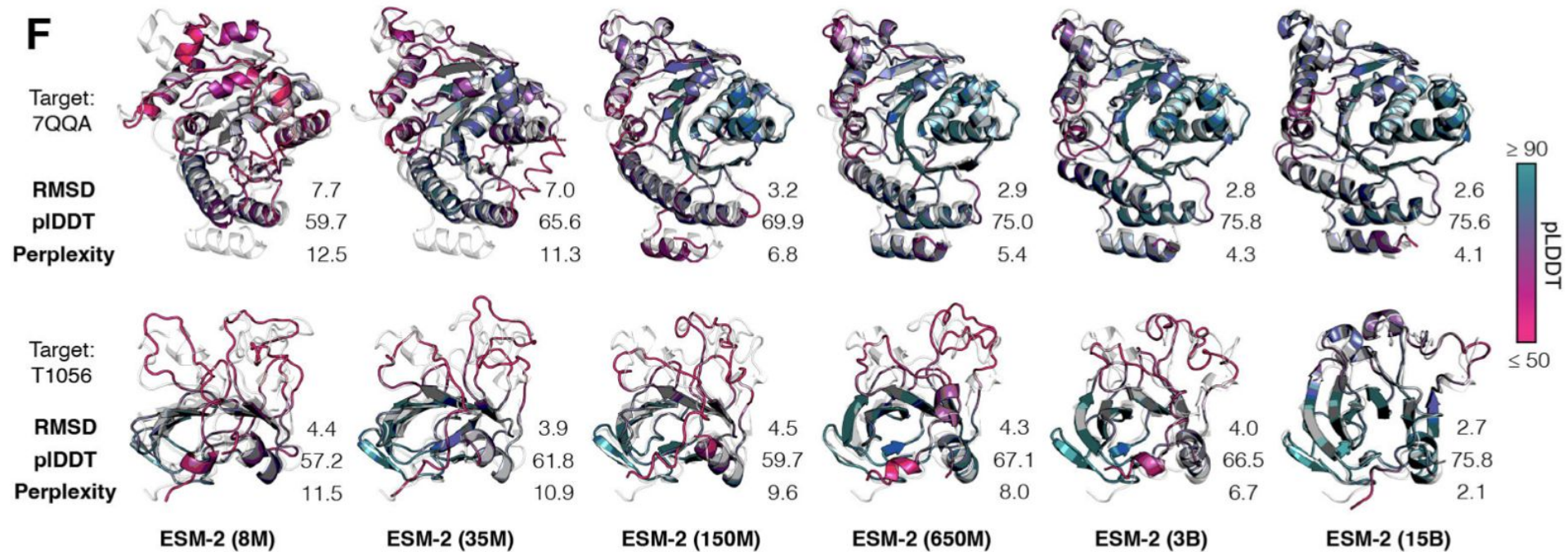
**D:** Left-to-right shows models from 8M to 15B parameters, consecutively comparing the smaller model (x-axis) against the next larger model (y-axis) in terms of unsupervised contact precision.

- Sequences with large changes in contact prediction performance exhibit large changes in language model understanding measured by pseudo-perplexity.

# TM-score on combined CASP14 and CAMEO test sets



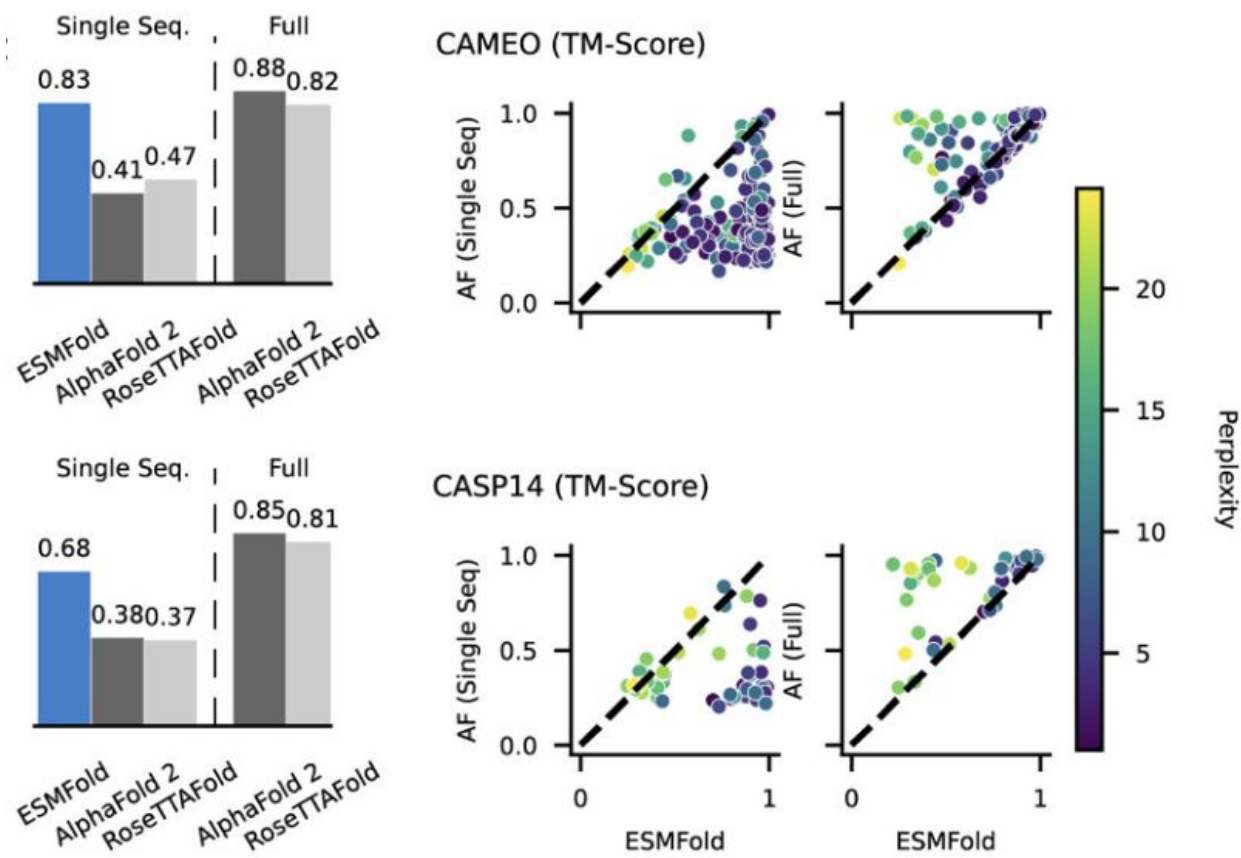
# structure predictions on CAMEO structure 7QQA and CASP target 1056



# More comparisons

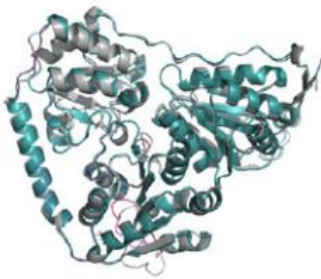
Model	# Params	Validation Perplexity	LR P@L	CASP14	CAMEO
ESM-2	8M	10.33	0.17	0.37	0.48
	35M	8.95	0.30	0.41	0.56
	150M	7.75	0.44	0.49	0.65
	650M	6.95	0.52	0.51	0.70
	3B	6.49	<b>0.54</b>	0.52	<b>0.72</b>
	15B	<b>6.37</b>	<b>0.54</b>	<b>0.55</b>	<b>0.72</b>
ESM-1b <sup>1</sup>	650M	—	0.41	0.42	0.64
Prot-T5-XL-UR50 (19)	3B	—	0.48	0.50	0.69
Prot-T5-XL-BFD (19)	3B	—	0.36	0.46	0.63
CARP (44)	640M	—	—	0.42	0.59

# Comparison with AlphaFold2

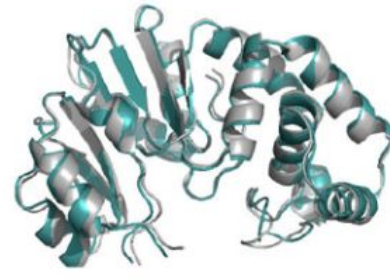




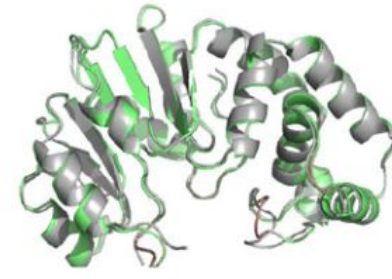
# Structure prediction comparison



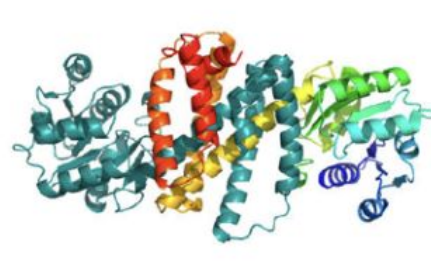
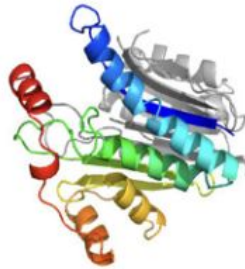
CASP14 T1076 (6XN8)  
TM-score ESMFold: 0.98  
TM-score AlphaFold: 0.99



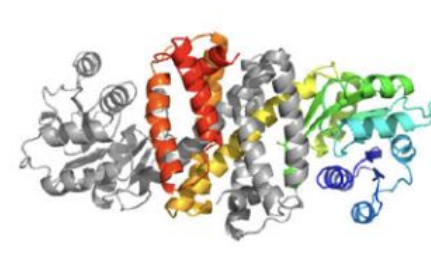
CASP14 T1057 (7M6B)  
TM-score ESMFold: 0.98  
TM-score AlphaFold: 0.97



Imine Reductase (7A3W)  
TM-score ESMFold: 0.956



L-asparaginase (6QQ8)  
TM-score ESMFold: 0.985

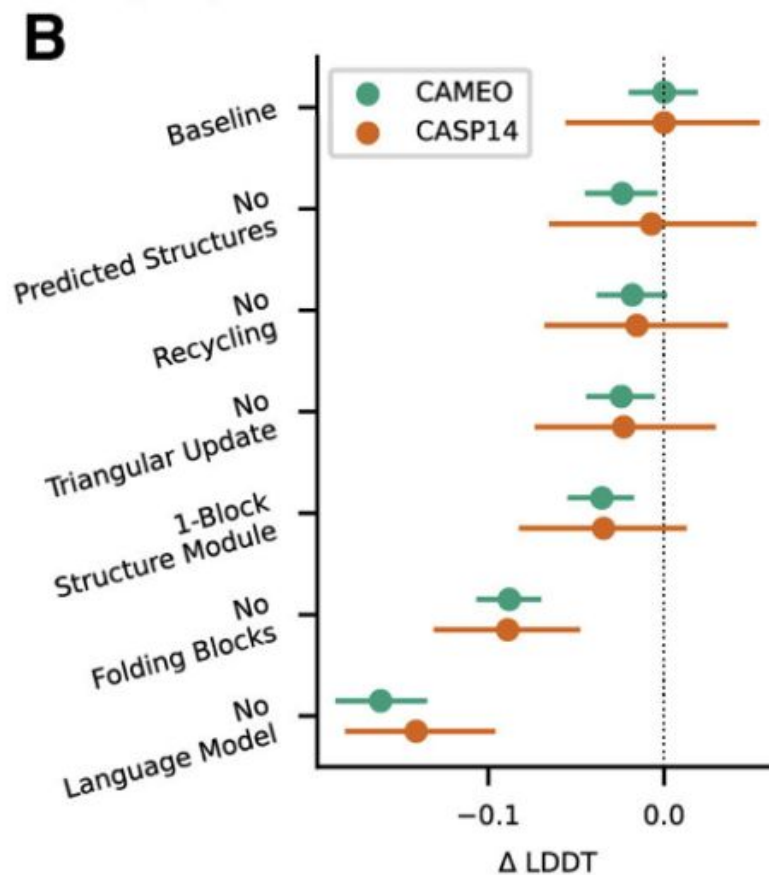




# Ablation Study on ESM-2

	<b>LR P@L</b>	<b>LR P@L/5</b>	<b>Validation Perplexity</b>
Baseline	0.381	0.626	8.42
No RoPE	0.365	0.599	8.62
Older UniRef Data	0.368	0.599	7.98
No UR90 Sampling	0.387	0.631	8.40

# Ablation Studies on ESMFold



# Inference Time

