

block b3

block b3



block b3

Transformers

block b3

Transformers

Protein Folding-AlphaFold2

# Transformers

**Transformers** → **Natural Languages**

inputs: words in an embedding

inputs are large

inputs are variable in length | no MLPs

parameter sharing

(CNNs RNNs)

connections between words (context)

[self-attention]

**Transformers** → **Biological sequence languages**

# Self-attention

[self-attention] → Transformers

---

**Attention Is All You Need**

2017

---

**Ashish Vaswani\***  
Google Brain  
avaswani@google.com

**Noam Shazeer\***  
Google Brain  
noam@google.com

**Niki Parmar\***  
Google Research  
nikip@google.com

**Jakob Uszkoreit\***  
Google Research  
usz@google.com

**Llion Jones\***  
Google Research  
llion@google.com

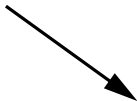
**Aidan N. Gomez\* †**  
University of Toronto  
aidan@cs.toronto.edu

**Lukasz Kaiser\***  
Google Brain  
lukaszkaizer@google.com

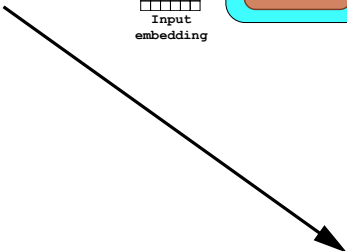
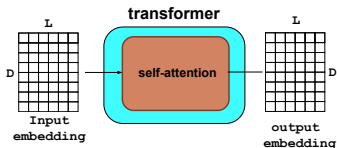
**Illia Polosukhin\* ‡**  
illia.polosukhin@gmail.com

3 transformers for natural language translation

self-attention



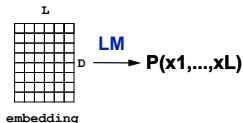
**Transformers**



**Language Models (LLMs)**

## [b4] Language Models (LLMs)

LM outputs are probability distributions over sequences  $P(x_1, \dots, x_L)$



$$P(\text{AUG}, \text{AUG}) = 0.0001$$

$$P(\text{AUG}, \text{stop-codon}) = 0.0000001$$

$$P(\text{AUG}, \text{codon}_1, \dots, \text{codon}_n, \text{stop-codon}) = 0.1$$

LMs are generative models

can sample new examples (synthetic biology)

LMs usually autoregressive ( which is exact, not an approximation)

$$x_1 \sim P(x_1)$$

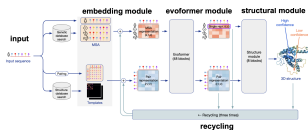
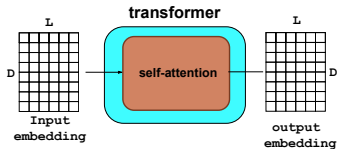
$$x_2 \sim P(x_2 \mid x_1)$$

$$x_3 \sim P(x_3 \mid x_1, x_2)$$

$$x_4 \sim P(x_4 \mid x_1, x_2, x_3)$$

[b3] self-attention

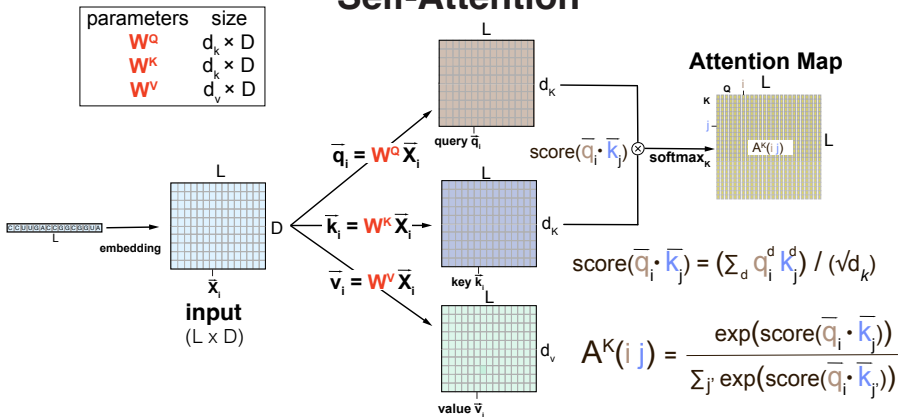
[b3] Transformers



[b3] AlphaFold2

[b4] Language Models (LLMs)

# Self-Attention



# Dot product self-attention maps

## 1. Queries

$$q[L, d_k] = x[L, D] * W^Q[D, d_k] + b^Q[d_K], \quad q_i(d)$$

## 2. Keys

$$k[L, d_k] = x[L, D] * W^K[D, d_k] + b^K[d_k], \quad k_i(d)$$

## 3. Values

$$v[L, d_v] = x[L, D] * W^V[D, d_v] + b^V[d_v], \quad v_i(d)$$

Linear dot-product score

$$\text{score}(q_i \cdot k_j) = \sum_{d=1}^{d_K} q_i(d) k_j(d) / \sqrt{(d_K)}.$$

**Attention** between  $i, j$  is the softmax of the scores respect to the keys

$$A_{i,j} = \frac{e^{\text{score}(q_i \cdot k_j)}}{\sum_{l=1}^L e^{\text{score}(q_i \cdot k_l)}}.$$

## Products between vectors $q[D]$ and $k[D]$

1. **scalar product**  $(q \cdot k)[1]$

$$q \cdot k = \sum_{d=1}^D q_d k_d$$

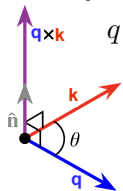
2. **element-wise product**  $(q \odot k)[D]$

$$q \odot k = [q_1 k_1, \dots, q_D k_D]$$

3. **outer product**  $(q \otimes k)[D, D]$

$$q \otimes k = \begin{bmatrix} q_1 k_1 & \dots & q_1 k_D \\ \vdots & \ddots & \vdots \\ q_D k_1 & \dots & q_D k_D \end{bmatrix}$$

4. **cross product** for  $D = 3$ ,  $(q \times k)[3]$



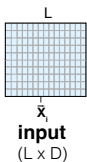
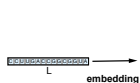
$$q \times k = \|q\| \|k\| \sin(\theta) \hat{n}$$

$$= \begin{vmatrix} \mathbf{i} & \mathbf{j} & \mathbf{k} \\ q_1 & q_2 & q_3 \\ k_1 & k_2 & k_3 \end{vmatrix}$$

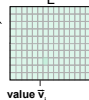
$$= (q_2 k_3 - q_3 k_2, -(q_1 k_3 - q_3 k_1), q_1 k_2 - q_2 k_1)$$

# Self-Attention (1d input)

| parameters | size           |
|------------|----------------|
| $W^Q$      | $d_k \times D$ |
| $W^K$      | $d_k \times D$ |
| $W^V$      | $d_v \times D$ |

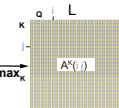


$$\begin{aligned} \bar{q}_i &= W^Q \bar{X}_i \\ \bar{k}_i &= W^K \bar{X}_i \\ \bar{v}_i &= W^V \bar{X}_i \end{aligned}$$



$$\text{score}(\bar{q}_i \cdot \bar{k}_j)$$

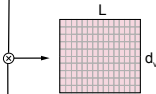
## Attention Map



$$A^K(i, j) = \frac{\exp(\text{score}(\bar{q}_i \cdot \bar{k}_j))}{\sum_j \exp(\text{score}(\bar{q}_i \cdot \bar{k}_j))}$$

$$\text{score}(\bar{q}_i \cdot \bar{k}_j) = (\sum_d q_i^d k_j^d) / (\sqrt{d_k})$$

$\Theta(L^2)$  time  
 $\Theta(L^2)$  memory



$$\bar{u}_i = \sum_j A^K(i, j) \bar{v}_j$$

update  
 $(L \times d_v)$

## Self-Attention (code)

```
# Simple self-attention
class SelfAttention(nn.Module):
    def __init__(self, d_model):
        super().__init__()
        self.d_model = d_model

        # uses all dimensions equal to d_model
        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)

    def forward(self, x):
        # x: (B, alen, d_model)
        B, L, D = x.shape

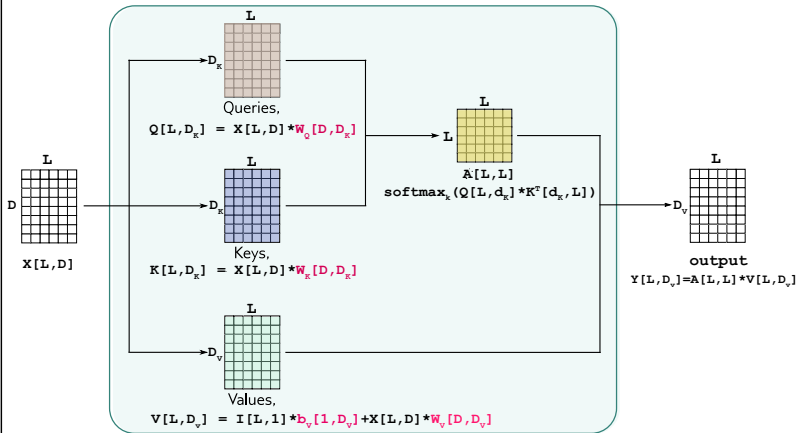
        # q[B, L, d_model]
        # k[B, L, d_model]
        # v[B, L, d_model]
        q = self.q_proj(x)
        k = self.k_proj(x)
        v = self.v_proj(x)

        # q [B, L, d_model]
        # k.transpose(-2,-1)[B, d_model, L]
        # attn[B, L, L]
        # v [B, L, d_model]
        # out[ B, L, d_model]
        # softmax wrt the L dimension coming from the keys
        attn = torch.softmax(q @ k.transpose(-2, -1) / (self.d_model ** 0.5), dim=-1)
        out = attn @ v # (B, L, d_model)

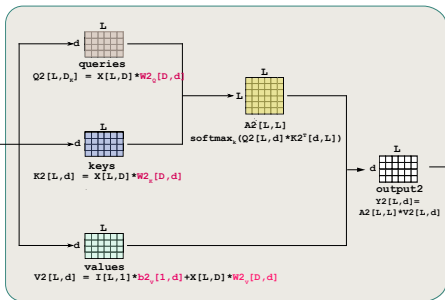
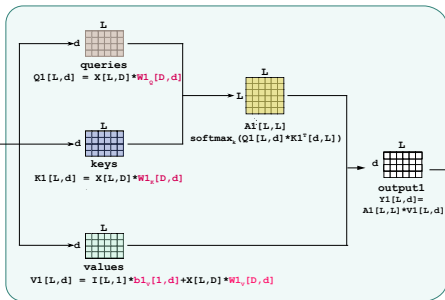
        return out
```

**sequence input = X[B, L, D]**

# Self Attention (matrix view)



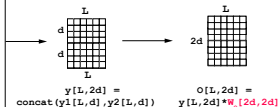
adapted from: UDL, Prince



$D = d\_model$

$H = n\_heads = 2$

$d = d\_head = d\_model / n\_heads$



adapted from: UDL, Prince

```

# multi-head self-attention
class MultiHeadAttention(nn.Module):
    def __init__(self, d_model, num_heads):
        super().__init__()
        assert d_model % num_heads == 0
        self.d_model = d_model
        self.num_heads = num_heads
        self.d_head = self.d_model // self.num_heads # H = D/n_heads

        # uses all dimensions equal to d_model
        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.o_proj = nn.Linear(d_model, d_model)

    def forward(self, x):
        # x: (B, alen, d_model)
        B, L, D = x.shape

        # H x d_head = d_model
        #
        # q[B, H, L, d_head]
        # k[B, H, L, d_head]
        # v[B, H, L, d_head]
        q = self.q_proj(x).reshape(B, L, self.num_heads, self.d_head).transpose(1, 2)
        k = self.k_proj(x).reshape(B, L, self.num_heads, self.d_head).transpose(1, 2)
        v = self.v_proj(x).reshape(B, L, self.num_heads, self.d_head).transpose(1, 2)

        # q[B, H, L, d_head]
        # k.transpose(-2, -1)[B, H, d_head, L]
        # attn[B, H, L, L]
        # v [B, H, L, d_head]
        # out [B, H, L, d_head]
        attn = torch.softmax(q @ k.transpose(-2, -1) / (self.d_head ** 0.5), dim=-1)
        out = attn @ v # (B, heads, L, d_head)

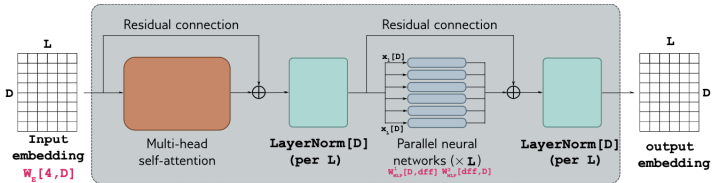
        # concatenate all heads together
        # out.transpose(1, 2)[B, L, H, d_head]
        # out[B, L, H*d_head]
        out = out.transpose(1, 2).reshape(B, L, D)

        # apply one last FC layer
        out = self.o_proj(out)

    return out

```

# Transformer Block



adapted from: UDL, Prince

```
# simple Transformer block = attention -> MLP
class SimpleTransformerBlock(nn.Module):
    def __init__(self, d_model, d_ff):
        super().__init__()

        self.attn = SelfAttention(d_model)

        self.mlp = nn.Sequential(
            nn.Linear(d_model, d_ff),
            nn.ReLU(),
            nn.Linear(d_ff, d_model),
        )

        self.norm_att = nn.LayerNorm(d_model)
        self.norm_ff = nn.LayerNorm(d_model)

    def forward(self, x):
        x = x + self.attn(x)
        x = x + self.mlp(self.norm_att(x))
        x = self.norm_ff(x)
        return x
```

# Layer Normalization

For an input embedding  $x[L, D]$ , the LayerNorm, calculates:

1. the **mean**

$$m_i = \frac{1}{D} \sum_{d=1}^D x_i^d$$

2. the **variance**

$$v_i = \frac{1}{D} \sum_{d=1}^D (x_i^d - m_i)^2$$

3. Then, the **values get updated** as

$$x_i[D] \leftarrow \frac{x_i[D] - m_i}{\sqrt{(v_i + \epsilon)}} \gamma + \delta, \gamma \text{ and } \delta \text{ are trained by the model.}$$

After the normalization layer, each  $x_i[D]$  vector has mean  $\delta$  and standard deviation  $\gamma$ .

# Positional encoding

$$vbins = [-32, -31, \dots, 32]$$

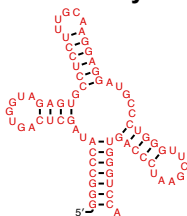
$$d_{ij} = j - i$$

$$p_{ij} = \text{LinearLayer}(\text{onehot}(d_{ij}, vbins))$$

Why Self-attention on  
multiple sequence alignments (MSAs)?

# Why Self-attention on multiple sequence alignments (MSAs)?

## RNA secondary structures



|    | 3' |   |    |   |
|----|----|---|----|---|
|    | A  | C | G  | U |
| 5' | 0  | 0 | 0  | 1 |
| C  | 0  | 0 | 10 | 0 |
| G  | 0  | 5 | 0  | 2 |
| U  | 4  | 0 | 0  | 0 |

maximum likelihood training

structural  
information

X-ray, NMR, CryoEM

## RNA alignments

```

GGCCCCAUGCCUCAG--GG--AGAGUCGCCUCCUUGCAAGGAGAUCC--CCUG--GGUCCGAADCCCAGGGGUCCA
CAUUGAGUAGCCGAA--AG--CAAGUACGGGUCUUAACAACUUUAA--UAGUAAUUAGCACUACUUCUUAAGGA
GGCCCCGUGACCAADUGGADGAGCGUUGACUACGGAAACAAAGGG--UAGG--GGUUCGACUCCUCCGGGGCCG
GGCCCCGUGACCCAGGUGGUGAGGCGACCCUUAUAGCGGAGGG--CGG--GGUUCGAGUCCUCCGGGGCCCA
AGGGCCAUUGUUAA--GG--AGAACAGAGGUCUCCAAACCUCGG--UGUG--GGUUCGAGUCCUCCGGCCCCG
AGGGCCGUGUUAUUUGG--AGAGCACCGGUCUCCAAACCGGGUG--UGGG--AGUUCGAGUCCUCCGGCCCCG
CGUGAUUGCCUCAGUUGG--AGAGCCACCUCUUGGUAAGGGUGAGG--CCCC--AGUUCGACUCCUCCGGUUCAGCA
GGGGCCGUGGCCAAGG--GG--AAGGCAGCGGDUUGGGUCCUUAUCU--CGGA--GGUUCGAADCCUCCGGCCCCG
GCCUUAUACACAG--GG--AGAGUUAAGCCUUGGUAAGCADAAG--GAG--GGUUCGAADCCUUAAGGGCCU
GCCUCCUUGGCUAAG--GGAGAGGACAGAGGUCUUAACAACUUUGG--UADA--GGUUCGAADCCUUAAGGGCCCA
GCADCCAGUGGCUAAG--GGUUAAGCGCCACUUAUAGUUAUUUGGC--GGUUCGAADCCUCCGGGGAGCA
AGUAAGGUCAGCUAA--UU--AAGCUADCGGGCCCAUACCCCGAATA--CGUU--GGUUAUADCCUCCGGCUAUA
AGAUUAUAGUAAAA--CA--AUUACAUAACUUDGUCAAGUUAUU--UADA--GAUCAAUAUCUUAUADCCUUA
GUUUCUGAGUGAA--UU--ACAACGADGUUUDCAUGUCUADUGG--UCGC--AGUUGAAGGCGUGGAGAAUA
GGUUAUAGGGGAA-----UUACAUGCCAAUUGCAAUUCAGAG--UGGU--GAGAAAU--CUUACUAGAGUA
ACUCCUUGAUUA--UU--AAUUAACUGACUCCUUAUUAUGA--UUUU--GAUUAAC--CCAGAAAGAGUA
UGAUUGAAGCCAGUA--AU--AGGGUUAUGGUGUUAACUUAUUUU--CGUA--GGUUAUADCCUCCGAUUCUG
CAUADGAAGCUA-----AGAGCGUUAACCUUUUAAGUUAAGUU--AGAG--ACCUAUA--AUUCCUAGUGUA
GUUAUUGAGCUUAUAAC--AAAGCAAGCACUGAAAUUGCUAGA--UGGA--U--AAUUGU--AUCCUUAUAACA
ACUUAUAGGAAUA--AG--AAUCUCCUGGUGGAAACCAAAA--CCUU--GGUCAAADCCUUAUUAAGUA
GGCCUUAAGCCAGCGGUGAGUGGACCCUUGAUAGGGUGAGG--CACA--AGUUCAGUCCUUAUUAAGGCCCA
    
```

evolutionary  
information

genomes, homology





# Why Self-attention on multiple sequence alignments (MSAs)?

## RNA secondary structures



|    | 3' |   |    |   |
|----|----|---|----|---|
|    | A  | C | G  | U |
| 5' | A  | 0 | 0  | 1 |
| C  | 0  | 0 | 10 | 0 |
| G  | 0  | 5 | 0  | 2 |
| U  | 4  | 0 | 0  | 0 |

maximum likelihood training

structural  
information

X-ray, NMR, CryoEM

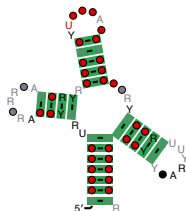
## RNA alignments

```

GGGCCCAUAGCCAGU--GGU--AGAGUCGCCUCCUUUGCAAGGAGGAUCC--CCUG--GGUUCGAUCCCAAGGGGUCCA
GAUAGAUAGACUGAA--AG--CAAGUACGGGUCUCUAAACCAUUAUA--UAGUAAAUAGCACUACUCCUUAUAGA
GGCCCGGAGCCCAUUGGUAAGAGCCGUUGACUACGGACAAAAGGU--UAGG--GGUUCGACCCUCCCGGGCCG
GGCCCGUAGCCUCCAGUUGAGCCAGCCGCGGAAAGCGGGAGGU--CGGU--GGUUCGAGGCCUCCAGGCCCA
AGGGGCAUAGUUUAAC--GGU--AGAACAGGGUCUCCAAACCCUGG--UGUG--GGUUCGAUCCUACUGCCCGC
AGGGGCGUAGUCCAUUGGU--AGAGCCACGGUUCUCCAAAACCGGGUGU--UGGG--AGUUCGAGUCUCCCGCCCGC
GGUGAUUGGCCAGUUGGU--AGAGGCCACCCUUGGUUAGGGGAGGU--CCCC--AGUUCGACUCUGGGUACGCA
GGGGCGGGCCCAAGU--GGU--AAGGAGCGGGUUUGUUGCCUUAUAGU--CGGA--GGUUCGAUCCUCCCGCCCGC
CCUCCUUAAUCAGCGU--GGU--AGAUUAUCGUCAGUAGGACAAAGU--CACU--GGUUCGAUCCUUAAGGGGU
CCUCCUUAUUGCUAAU--GGAUAGGACAGAGGUCUUCUUAACCUUUGG--UUAU--GGUUCGAUCCUUAUGGACGCA
GCUCUCCAUUGGCCUAAU--GGUUAAGGCCCCCAACUCUAAAUUGUUAUUUGUGG--GGUUCGAUCCUCCUGGGAGUCA
AGUAAGGUUAGCUAA--UU--AAGCUAUCGGGCCCAACCCCGAAAA--CGUU--GGUUAUAAUCCUCCCGUACUA
AGAUUAUAGUAAAUAU--CA--AUUACAUAACTUUGCCAAAGUUAUAU--UUAU--GAUCAUAUCCUUAUAUUCUUA
GUUCUUGAGUGUA--UU--ACAACGAUGAUUUUCUUGCCAUUGG--UCGC--AGUUGAAUUGCCUGGUAUUAUA
GGCCUUAAGGUGUA-----UUCAGUUCGAAUUGCAAUUGCAGAG--UGUA--GAUAUAU--CCUUAUAAGUAU
ACCCUCCUAGUUAUA--UU--AAUUAUCUGACUCCUCAAUUUGAGGA--UUUC--GAUUAUAC--CCAGAAGAGUA
AGAUUGAAGCCAGUA--AU--AGGGUUAUUGCCUUAUCAAUUUUUU--CGUA--GGUUAUAAUCCUCCGCAUUCAG
CACUAGAGGCUA-------AGAGCUUAUCCUUUAAGUUAAGUU--AGAG--ACCUUAAA--AUCCUUAUGUA
GUUAAGUGAGCUUAUUAAC--AAAGCAAGCACCGAAAAGUUGUAGA--UGGA--U-AAUGGG--AUCCCAUUAACA
AUUUUAUAGCAUAU--AGU--AAUCCUUGUCUUGGACCCAAUA--CCUU--GGGCAAUCCUUAUUAUAUAUA
GGCCUUAUAGCCAGUCAGUCGGUAGAGUCCACCCUUGAAGGGUGAGGU--CACA--AGUUCGAGUCUUGUUAAGGCCA
    
```

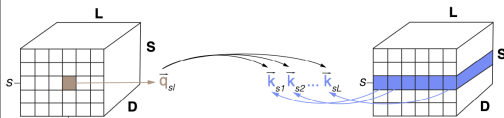
evolutionary  
information

genomes, homology



# MSA attention by row/column

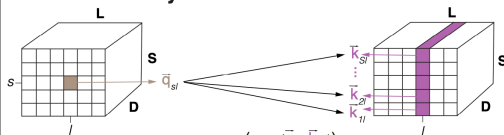
## self-attention by row



$$A(s|s'l) = \frac{\exp(\text{score}(\bar{q}_{sl} \cdot \bar{k}_{s'l}))}{\sum_m \exp(\text{score}(\bar{q}_{sl} \cdot \bar{k}_{sm}))}$$

$$\bar{R}_{sl} = \sum_r A(s|s'l) \bar{v}_{s'l} \quad \Theta(S \times L^2)$$

## self-attention by column



$$A(s|s'l) = \frac{\exp(\text{score}(\bar{q}_{sl} \cdot \bar{k}_{s'l}))}{\sum_r \exp(\text{score}(\bar{q}_{sl} \cdot \bar{k}_{rl}))}$$

$$\bar{C}_{sl} = \sum_r A(s|s'l) \bar{v}_{s'l} \quad \Theta(L \times S^2)$$

parameters size

$W^Q$   $D \times d_q$

$W^K$   $D \times d_k$

$W^V$   $D \times d_v$

$\bar{X}_{sl} = \bar{X}_{sl} W^Q$  **input**, size =  $D$   
 $\bar{q}_{sl} = \bar{X}_{sl} W^Q$  **query**, size =  $d_q$   
 $\bar{k}_{sl} = \bar{X}_{sl} W^K$  **key**, size =  $d_k$   
 $\bar{v}_{sl} = \bar{X}_{sl} W^V$  **value**, size =  $d_v$

```
# Row attention =
# attention across residues for each sequence (MSA row)
# thus: row[S] is fixed, L and D move
class RowAttention(nn.Module):
    def __init__(self, d_model, num_heads):
        super().__init__()
        self.attn = MultiHeadAttention(d_model, num_heads)

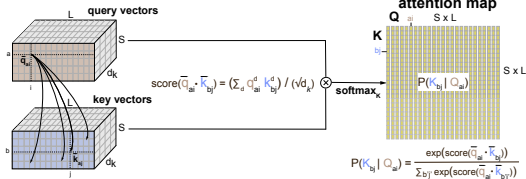
    def forward(self, x):
        # x: (B, S, L, d_model)
        B, S, L, D = x.shape
        x = x.reshape(B * S, L, D)
        x = self.attn(x)
        return x.reshape(B, S, L, D)

# Column attention =
# attention across MSA sequences at each residue (MSA column)
# L is fixed, we move on S and D.
# Thus we need to put S, D as the last two dimension
class ColumnAttention(nn.Module):
    def __init__(self, d_model, num_heads):
        super().__init__()
        self.attn = MultiHeadAttention(d_model, num_heads)

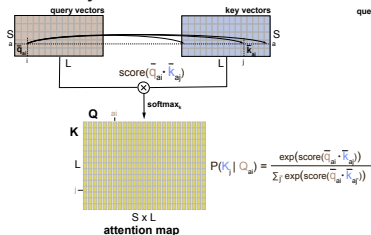
    def forward(self, x):
        # x: (B, S, L, d_model)
        B, S, L, D = x.shape
        x = x.transpose(1, 2).reshape(B * L, S, D)
        x = self.attn(x)
        return x.reshape(B, L, S, D).transpose(1, 2)
```

# Alignment Self-Attention

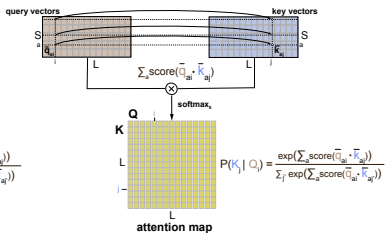
| parameters | size           |
|------------|----------------|
| $W^Q$      | $D \times d_q$ |
| $W^K$      | $D \times d_k$ |
| $W^V$      | $D \times d_v$ |



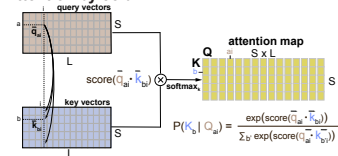
## Attention by row



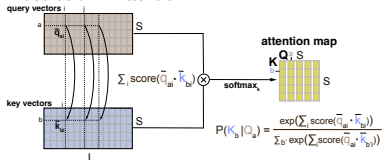
## tied row Attention



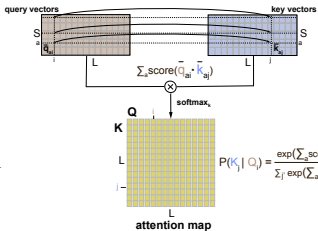
## Attention by column



## tied column Attention



## typed row Attention



## Tied row attention

```

class TiedRowAttention(nn.Module):
    def __init__(self, d_model, num_heads, tie_node="mean"):
        super().__init__()
        self.num_heads = num_heads
        self.head_dim = d_model // num_heads
        self.tie_node = tie_node # "mean" or "sum"

        # shared projections across all rows
        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.o_proj = nn.Linear(d_model, d_model)

    def forward(self, x):
        # x: (B, S, L, D)
        B, S, L, D = x.shape

        # flatten rows so projections are tied
        x_flat = x.reshape([B*S, L, D])

        # queries, keys, values [B*S, L, D] (reshape)
        Q = self.q_proj(x_flat)
        K = self.k_proj(x_flat)
        V = self.v_proj(x_flat)

        # multi-head split
        # queries, keys, values [B*S, L, H, DH] (reshape)
        # queries, keys, values [B*S, H, L, DH] (transpose 1,2)
        Q = Q.reshape([B*S, L, self.num_heads, self.head_dim]).transpose(1, 2)
        K = K.reshape([B*S, L, self.num_heads, self.head_dim]).transpose(1, 2)
        V = V.reshape([B*S, L, self.num_heads, self.head_dim]).transpose(1, 2)

        # per-row attention scores
        # R_scores [B*S, H, L, L] = Q [B*S, H, L, DH] * K^T [B*S, H, DH, L]
        R_scores = torch.matmul(Q, K.transpose(-2, -1)) / (self.head_dim ** 0.5)
        # R_scores [B, S, H, L, L]
        R_scores = R_scores.reshape([B, S, self.num_heads, L, L])

        # tie at the attention-weight level
        # TR_scores [B, 1, H, L, L]
        TR_scores = R_scores.mean(dim=1, keepdim=True) # (B, 1, H, L, L)

        # softmax AFTER tying
        # TR_attn [B, 1, H, L, L]
        TR_attn = torch.softmax(TR_scores, dim=-1) # (B, 1, H, L, L)

        # expand tied attention to all rows
        # TR_attn [B, S, H, L, L]
        TR_attn = TR_attn.repeat(1, S, 1, 1, 1) # (B, S, H, L, L)

        # print TR_attn shape
        print("\ntied-row attention (TR_attn) maps dimensions [B, S, H, L, L]", TR_attn.shape)

        # reshape V back to [B, S, H, L, DH] (from [B*S, H, L, DH])
        V = V.reshape([B, S, self.num_heads, L, self.head_dim])

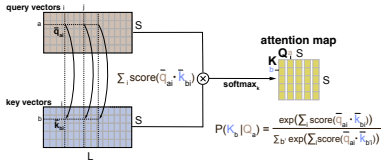
        # out [B, S, H, L, DH] = TR_attn [B, S, H, L, L] * V [B, S, H, L, DH]
        out = torch.matmul(TR_attn, V) # (B, S, H, L, DH)

        # merge heads
        # out [B, S, H, L, DH]
        # out [B, S, 1, H, DH] (transposed 2,3)
        # out [B, S, L, D] (merged all heads)
        out = out.transpose(2, 3).reshape([B, S, L, D])

        # project
        out = self.o_proj(out)

        return out
    
```

## typed column Attention



## Tied column attention

```
class TiedColumnAttention(nn.Module):
    def __init__(self, d_model, num_heads, tie_mode="mean"):
        super().__init__()
        self.num_heads = num_heads
        self.head_dim = d_model // num_heads
        self.tie_mode = tie_mode # "mean" or "sum"

        # shared projections across all rows
        self.q_proj = nn.Linear(d_model, d_model)
        self.k_proj = nn.Linear(d_model, d_model)
        self.v_proj = nn.Linear(d_model, d_model)
        self.o_proj = nn.Linear(d_model, d_model)

    def forward(self, x):
        # x: [B, S, L, D]
        B, S, L, D = x.shape

        # transpose
        # x: [B, S, L, D] -> x: [B, L, S, D]
        x = x.transpose(1,2)

        # flatten rows so projections are tied
        x_flat = x.reshape([B*L, S, D])

        # queries, keys, values [B*L, S, D]
        Q = self.q_proj(x_flat)
        K = self.k_proj(x_flat)
        V = self.v_proj(x_flat)

        # Multi-head split
        # queries, keys, values [B*L, S, H, DH] (reshape)
        # queries, keys, values [B*L, H, S, DH] (transpose 1,2)
        Q = Q.reshape([B*L, S, self.num_heads, self.head_dim].transpose(1, 2))
        K = K.reshape([B*L, S, self.num_heads, self.head_dim].transpose(1, 2))
        V = V.reshape([B*L, S, self.num_heads, self.head_dim].transpose(1, 2))

        # per-column attention scores
        # C_scores = torch.matmul(Q, K.transpose(0, -1)) / (self.head_dim ** 0.5)
        C_scores = torch.matmul(Q, K.transpose(0, -1)) / (self.head_dim ** 0.5)
        # scores [B, L, H, S, S]
        C_scores = C_scores.reshape([B, L, self.num_heads, S, S])

        # tie at the attention-weight level
        # TC_scores [B, 1, H, S, S]
        TC_scores = C_scores.mean(dim=1, keepdim=True) # (B, 1, H, S, S)

        # softmax AFTER tying
        # TC_attn [B, 1, H, S, S]
        TC_attn = torch.softmax(TC_scores, dim=-1) # (B, 1, H, S, S)

        # expand tied attention to all rows
        # TC_attn [B, L, H, S, S]
        TC_attn = TC_attn.repeat([1, L, 1, 1, 1]) # (B, L, H, S, S)

        # print TC_attn shape
        print("\ntied-column attention (TC_attn) maps dimensions [B, L, H, S, S]^T, TC_attn.shape")

        # reshape V to [B, L, H, S, DH] (from [B*L, H, S, DH])
        V = V.reshape([B, L, self.num_heads, S, self.head_dim])

        # out [B, L, H, S, DH] = TC_attn [B, L, H, S, S] * V [B, L, H, S, DH]
        out = torch.matmul(TC_attn, V) # (B, L, H, S, DH)

        # merge heads
        # out [B, L, H, S, DH]
        # out [B, L, S, H, DH] (transposed 2,3)
        # out [B, L, S, D] (merged all heads)
        out = out.transpose(2, 3).reshape([B, L, S, D])

        # out [B, S, L, D] <- out [B, L, S, D]
        out = out.transpose(1, 2)

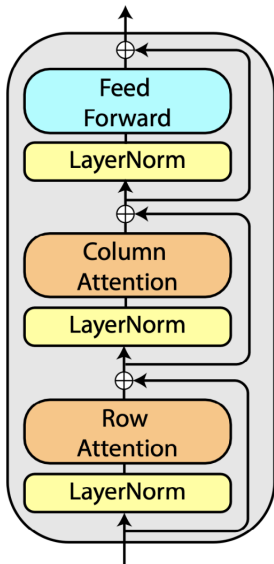
        # project
        out = self.o_proj(out)

        return out
```

# MSA self-attention

MSA Transformer, Rao *et al*, 2021

## MSA transformer



```
# MSA Transformer block = row attention → column attention → MLP  
class MSATransformerBlock(nn.Module):
```

```
    def __init__(self, d_model, num_heads, d_ff):  
        super().__init__()
```

```
        self.row_attn = RowAttention(d_model, num_heads)  
        self.col_attn = ColumnAttention(d_model, num_heads)
```

```
        self.mlp = nn.Sequential(  
            nn.Linear(d_model, d_ff),  
            nn.ReLU(),  
            nn.Linear(d_ff, d_model),  
        )
```

```
        self.norm_row = nn.LayerNorm(d_model)  
        self.norm_col = nn.LayerNorm(d_model)  
        self.norm_ff = nn.LayerNorm(d_model)
```

```
    def forward(self, x):  
        x = x + self.row_attn(self.norm_row(x))  
        x = x + self.col_attn(self.norm_col(x))  
        x = x + self.mlp(self.norm_ff(x))  
        return x
```

# AlphaFold2

## Article

# Highly accurate protein structure prediction with AlphaFold

2024 Nobel in Chemistry  
Jumper, Hassabis, Baker

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021

Accepted: 12 July 2021

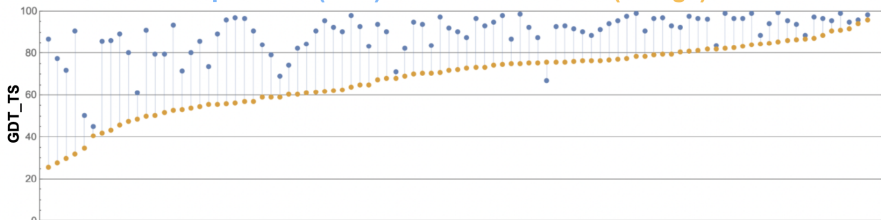
Published online: 15 July 2021

Open access

 Check for updates

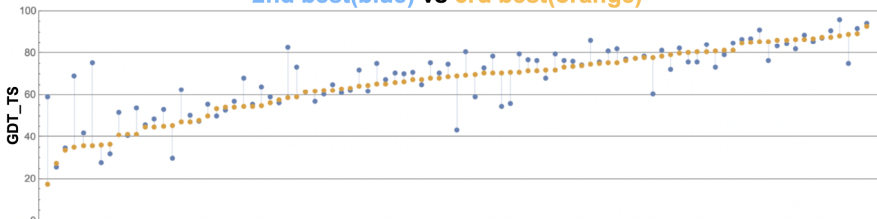
John Jumper<sup>1,4</sup>, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstein<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>

### AlphaFold2(blue) vs 2nd best method(orange)



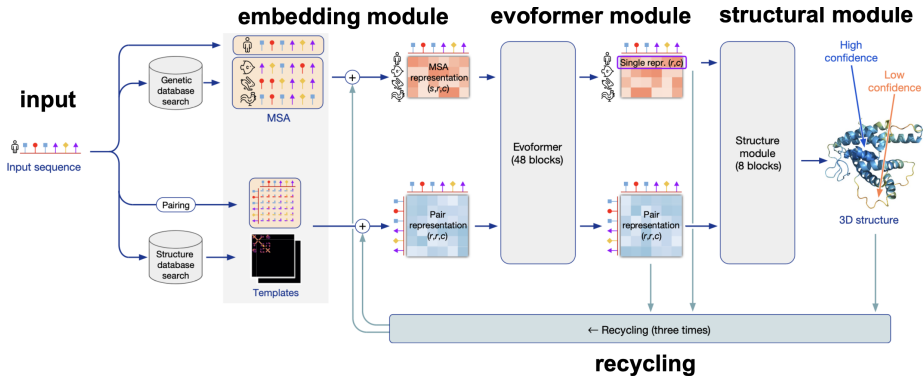
CASP14 targets

### 2nd best(blue) vs 3rd best(orange)



CASP14 targets

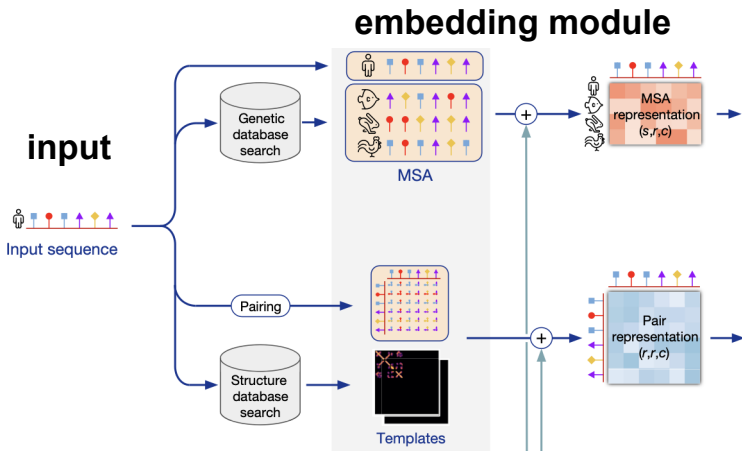
# AlphaFold2 Summary (Figure 1)



# AlphaFold2 (self-reported) Key Innovations

1. **evoformer** Join embedding of alignments and pair representations  
“Direct reasoning between the spatial and evolutionary relationships”
2. **structural module** Residue gas: a new representation of 3D structure as a collection of rotation and translation for each amino acid.
3. **structural module** Equivariant attention
4. **training** End-to-end structure prediction
5. **training** Use of intermediate losses
6. **training** Masked msa loss
7. **training** Distillation: learning from unlabelled sequences

# AlphaFold2 Summary (Figure 1)

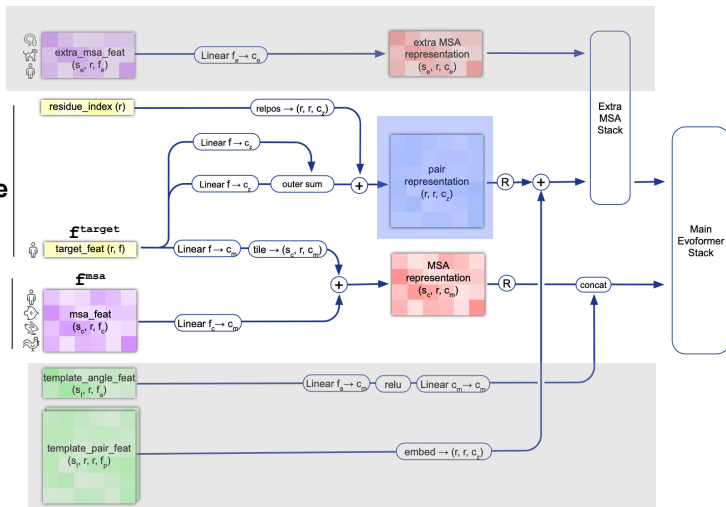


# AF2 Embedding Module

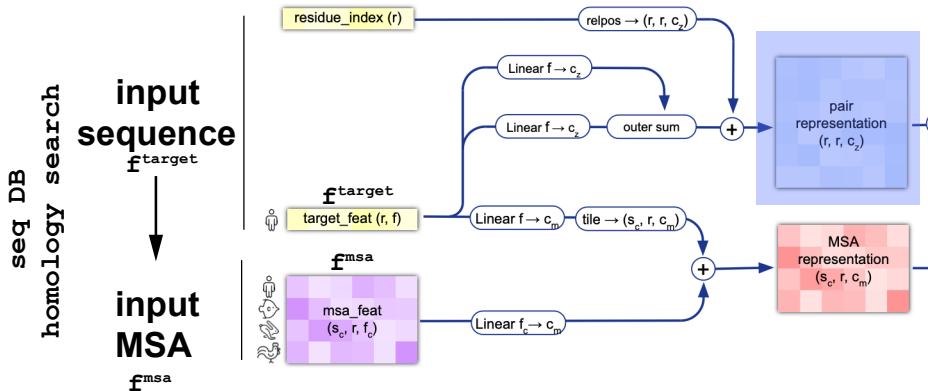
seq DB  
homology search

input  
sequence  
 $f^{\text{target}}$

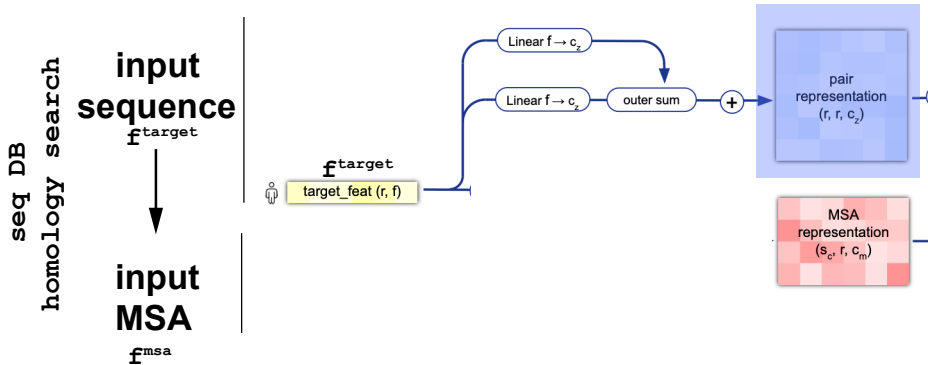
input  
MSA  
 $f^{\text{msa}}$



# AF2 Embedding Module



# AF2 Embedding Module



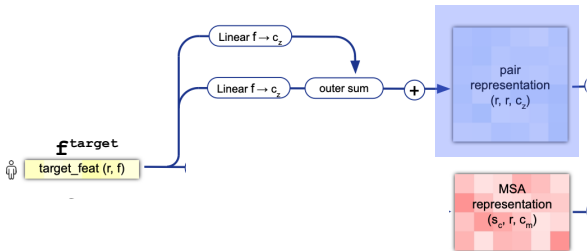
seq DB  
homology search

input  
sequence

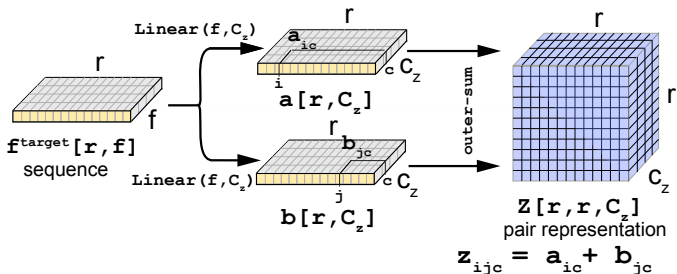
$f^{\text{target}}$

input  
MSA

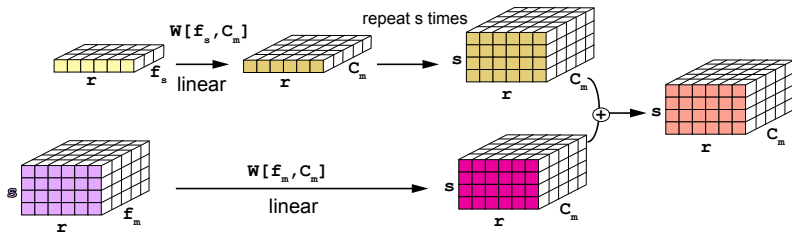
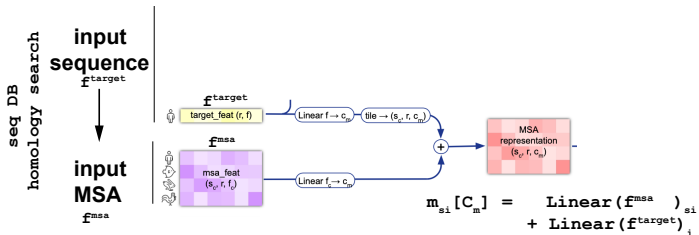
$f^{\text{msa}}$



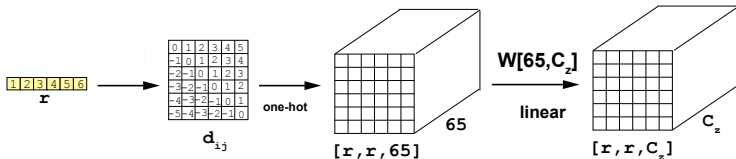
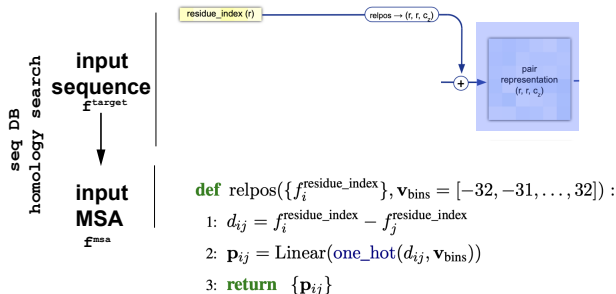
from seq to pair rep by “outer-sum”



# AF2 Embedding Module

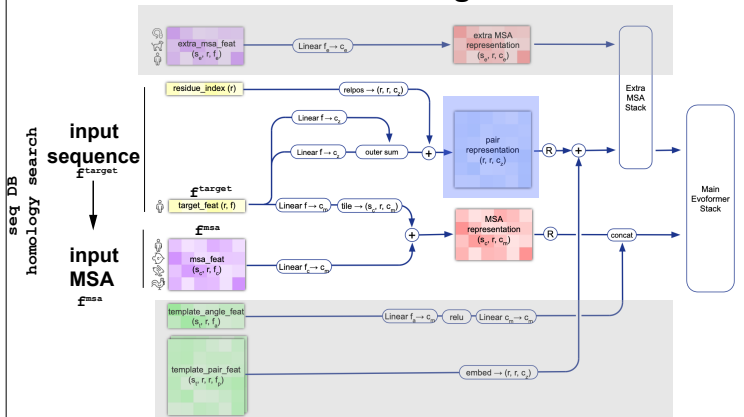


# AF2 Embedding Module



This inductive bias de-emphasizes primary sequence distances. Compared to the more traditional approach of encoding positions in the frequency space, this relative encoding scheme empirically allows the network to be evaluated without quality degradation on much longer sequences than it was trained on.

# AF2 Embedding Module



**InputEmbedding ( $f^{\text{target}}, f^{\text{msa}}, f^{\text{index}}$ ):**

$$a_i = \text{Linear}(f_i^{\text{target}})$$

$$b_i = \text{Linear}(f_i^{\text{target}})$$

$$a_i[C_2], b_i[C_2] \quad C_2 = 128$$

$$z_{ij} = a_i \oplus b_j$$

$$z_{ij} = z_{ij} + \text{relpos}(f^{\text{index}})$$

$$z_{ij}[C_2]$$

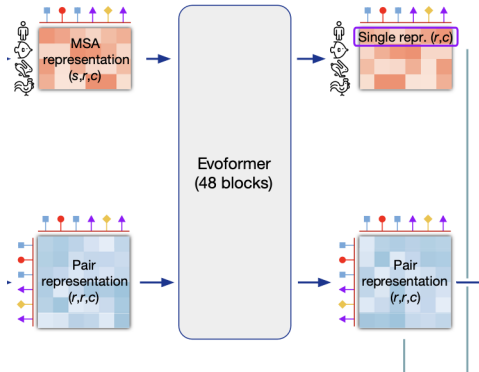
$$m_{si} = \text{Linear}(f_{si}^{\text{msa}}) + \text{Linear}(f_i^{\text{target}})$$

$$m_{si}[C_m] \quad C_m = 256$$

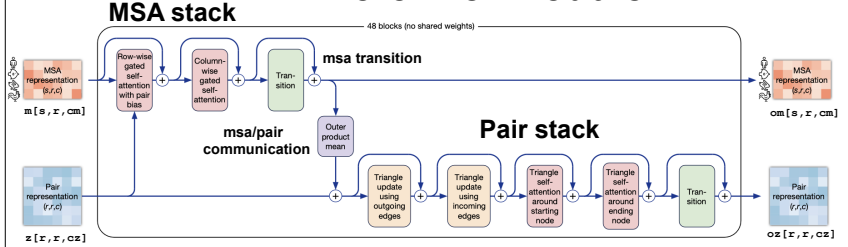
**return  $\{m_{si}\}, \{z_{ij}\}$**

# AlphaFold2 Summary (Figure 1)

## evoformer module



# AF2 Evoformer Module



## MSA stack

RowAttentionBias (m, z) :

$$\begin{aligned}
 m_{sj} &= \text{LayerNorm}(m_{sj}) \\
 q_{si}^h, k_{si}^h, v_{si}^h &= \text{Linear}(m_{si}) \\
 b_{ij}^h &= \text{Linear}(\text{LayerNorm}(z_{ij})) \\
 a_{sij}^h &= \text{softmax}_K(q_{si}^h k_{sj}^h + b_{ij}^h) \\
 g_{si}^h &= \text{sigmoid}(\text{Linear}(m_{si})) \\
 o_{si}^h &= g_{si}^h \sum_j (a_{sij}^h v_{sj}^h) \\
 om_{si} &= \text{Linear}(\text{concat}_h(o_{si}^h))
 \end{aligned}$$

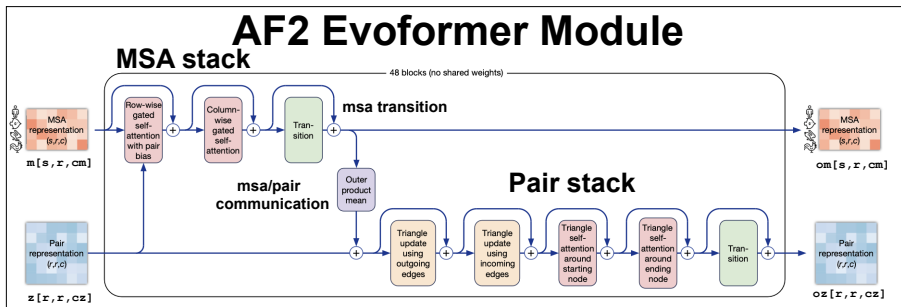
[s, r, r]

ColAttention (m, z) :

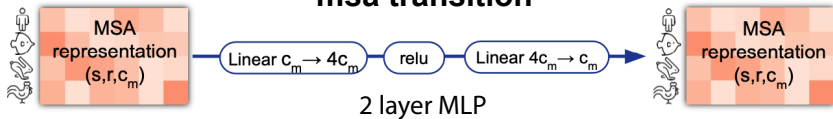
$$\begin{aligned}
 m_{sj} &= \text{LayerNorm}(m_{sj}) \\
 q_{si}^h, k_{si}^h, v_{si}^h &= \text{Linear}(m_{si}) \\
 a_{sti}^h &= \text{softmax}_K(q_{si}^h k_{ti}^h) \\
 g_{si}^h &= \text{sigmoid}(\text{Linear}(m_{si})) \\
 o_{si}^h &= g_{si}^h \sum_t (a_{sti}^h v_{ti}^h) \\
 om_{si} &= \text{Linear}(\text{concat}_h(o_{si}^h))
 \end{aligned}$$

[s, s, r]

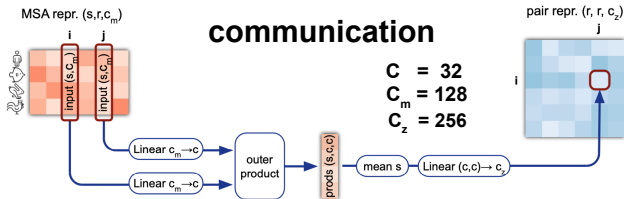
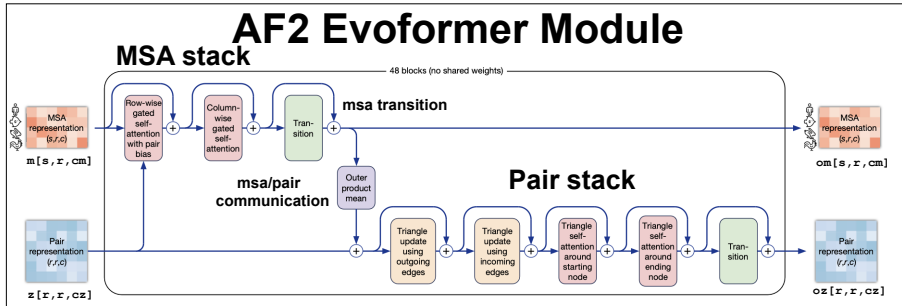
# AF2 Evoformer Module



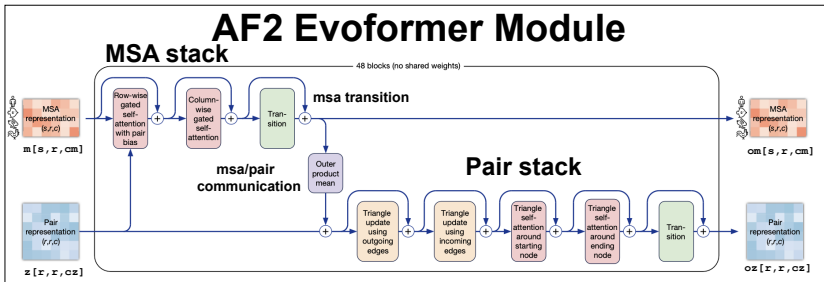
## msa transition



# AF2 Evoformer Module



# AF2 Evoformer Module



## Pair stack

### triangle updates

(parameters not shared)

$$a_{i,j}, b_{i,j} = \text{sig}(\text{Linear}(z_{i,j})) \odot \text{Linear}(z_{i,j})$$

$$g_{i,j} = \text{sigmoid}(\text{Linear}(z_{i,j}))$$

outgoing



$$oz_{i,j} = \text{Linear}(\text{sum}_k (a_{ik} \odot b_{jk}))$$

incoming



$$oz_{i,j} = \text{Linear}(\text{sum}_k (a_{ki} \odot b_{kj}))$$

### triangle self-attentions $a[r, r, r]$

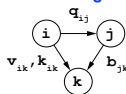
(parameters not shared)

$$q_{i,j}, k_{i,j}, v_{i,j}^h = \text{Linear}(z_{i,j})$$

$$b_{i,j} = \text{Linear}(z_{i,j})$$

$$g_{i,j} = \text{sigmoid}(\text{Linear}(z_{i,j}))$$

starting

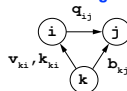


$$a_{i,j,k}^h = \text{softmax}_k (q_{i,j}^h k_{ik}^h + b_{jk}^h)$$

$$o_{i,j}^h = g_{i,j} \odot \text{sum}_k (a_{i,j,k}^h v_{ik}^h)$$

$$oz_{i,j} = \text{Linear}(\text{concat}_h (o_{i,j}^h))$$

ending

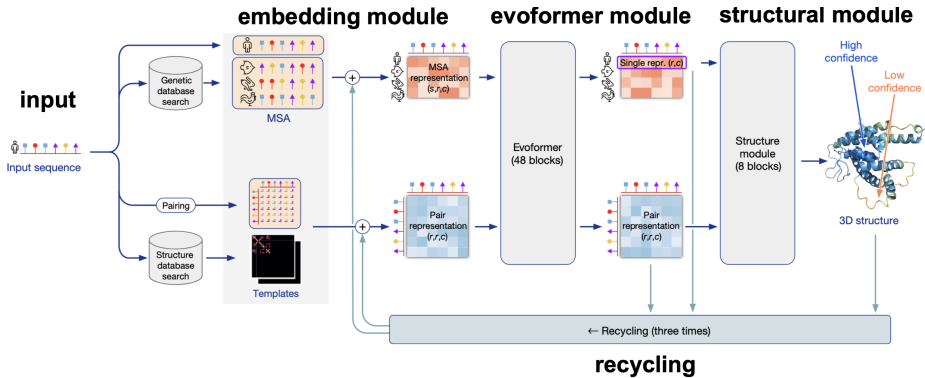


$$a_{i,j,k}^h = \text{softmax}_k (q_{i,j}^h k_{kj}^h + b_{ki}^h)$$

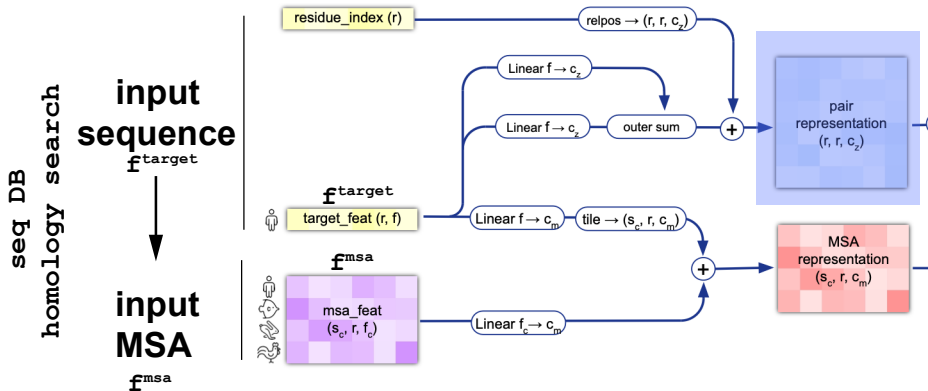
$$o_{i,j}^h = g_{i,j} \odot \text{sum}_k (a_{i,j,k}^h v_{ki}^h)$$

$$oz_{i,j} = \text{Linear}(\text{concat}_h (o_{i,j}^h))$$

# AlphaFold2 Summary



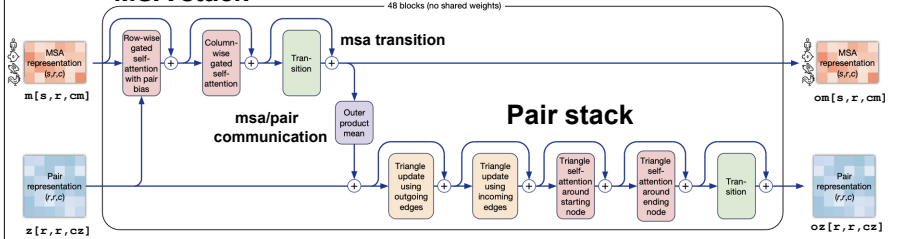
# AF2 Embedding Module



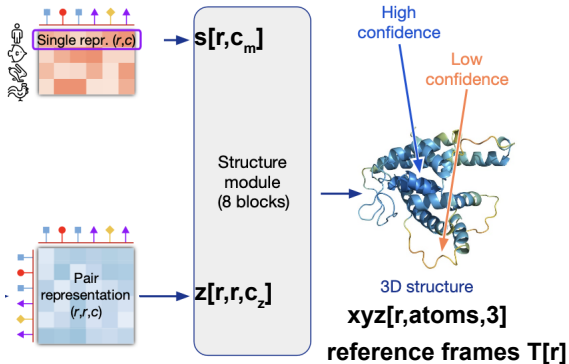
# AF2 Evoformer Module

## MSA stack

48 blocks (no shared weights)

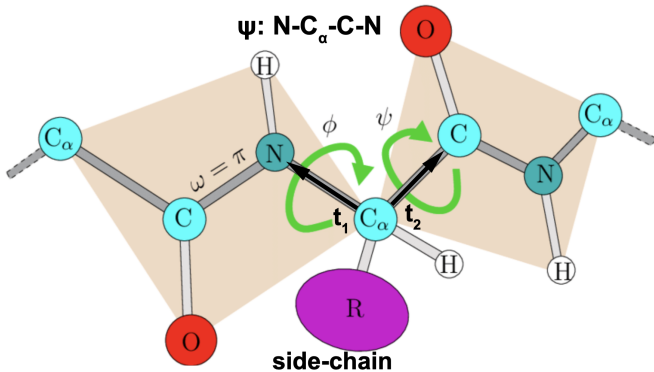


# structural module



$$\text{reference frames } T_i = t_i[3] \text{ translation} \\ + R_i[3,3] \text{ rotation}$$

# amino acid backbone reference frame



$$\vec{t}_1 = \vec{N} - \vec{C}_\alpha$$

$$\vec{t}_2 = \vec{C} - \vec{C}_\alpha$$

$$\vec{e}_1 = \text{norm}(\vec{t}_1)$$

$$\vec{e}_2 = \text{norm}(\vec{t}_2 - (\vec{t}_2 \cdot \vec{e}_1) \vec{e}_1)$$

$$\vec{e}_3 = \vec{e}_1 \times \vec{e}_2$$

$$R = \begin{bmatrix} \vec{e}_1 \\ \vec{e}_2 \\ \vec{e}_3 \end{bmatrix} \quad \vec{t} = \vec{C}_\alpha$$

$$T = (R, \vec{t})$$

1 Reference frame:  $\text{N}-\text{C}_\alpha-\text{C}$

1 backbone torsion angle  $\psi$  to determine O

4 side-chain torsion angles

amino acid chain = sequence of reference frames

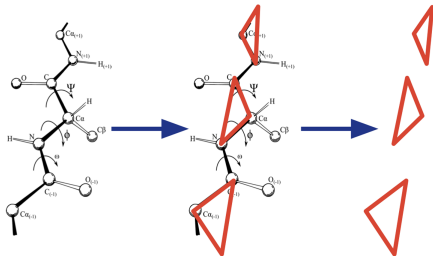
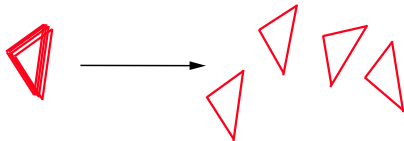


Image: Dcrjsr, vectorised Adam Rędzikowski (CC BY 3.0, Wikipedia)

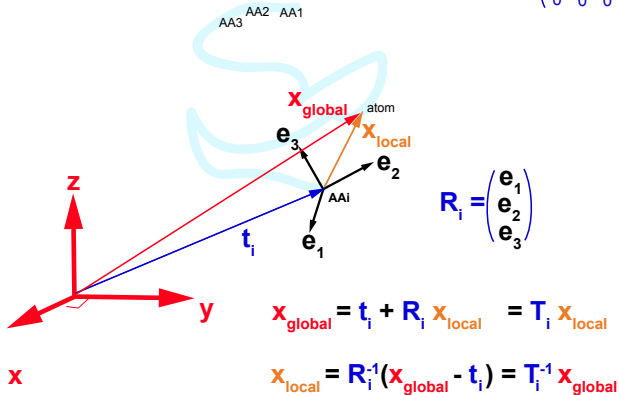


# Special Euclidean Group

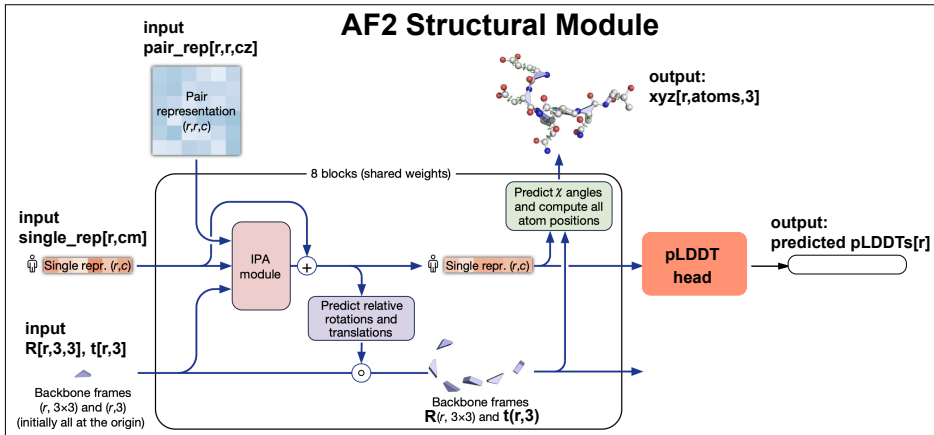
SE(3) = Translation (t) + Rotation (R)

$$\mathbf{t} = (t_1, t_2, t_3) \quad \mathbf{R} = \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix} \text{SO(3) = Special Orthogonal group}$$

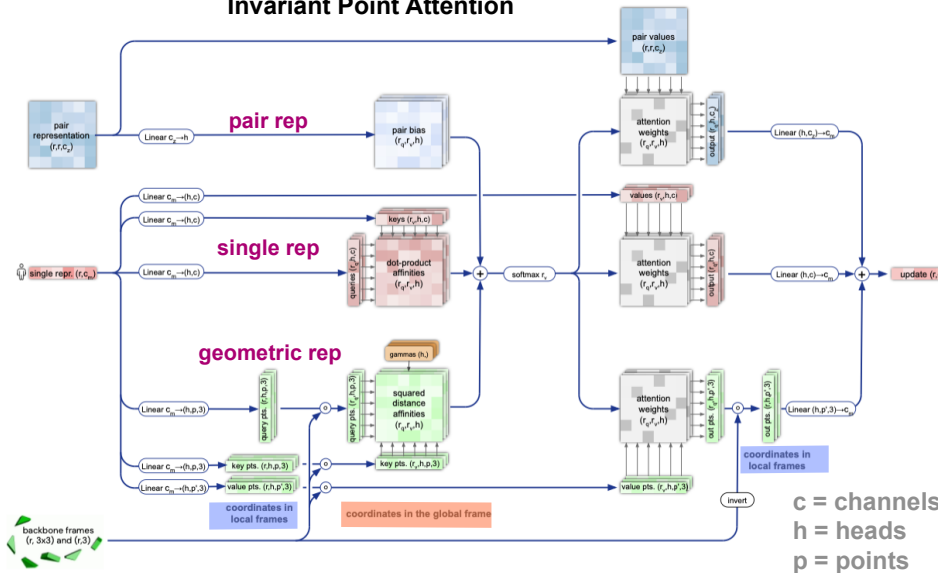
$$\mathbf{T} = (\mathbf{R}, \mathbf{t}) = \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$



# AF2 Structural Module



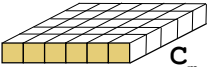
# Invariant Point Attention





$$s[r, C_m] \xrightarrow{\text{IPA}} \text{update}[r, C_m]$$

After IPA.....



$$\text{update}[r, C_m] \xrightarrow[\text{linear}]{W[C_m, 6]} \text{quaternion}[r, 6] = \text{b}[r], \text{c}[r], \text{d}[r], \text{t}[r, 3]$$

$T_i = \text{rotation } R_i[3,3] + \text{translation } \vec{t}_i[3] \text{ per position } 1 \leq i \leq r$

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

$$R_i = \begin{pmatrix} a_i^2 + b_i^2 - c_i^2 - d_i^2 & 2b_i c_i - 2a_i d_i & 2b_i d_i + 2a_i c_i \\ 2b_i c_i + 2a_i d_i & a_i^2 - b_i^2 + c_i^2 - d_i^2 & 2c_i d_i - 2a_i b_i \\ 2b_i d_i - 2a_i c_i & 2c_i d_i + 2a_i b_i & a_i^2 - b_i^2 - c_i^2 + d_i^2 \end{pmatrix}$$

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

# AF2 training

**Labels:** xyz coordinates of all heavy atoms in the 3D structure

## The Protein Data Bank

RCSB PDB Deposit Search Visualize Analyze Download Learn About Careers COVID-19 Help Contact us MyPDB

RCSB PDB PROTEIN DATA BANK 250,441 Structures from the PDB archive 1,068,577 Computed Structure Models (CSM)

Enter search term(s), Ligand ID or sequence  Include CSM

Advanced Search | Chemical Search | Browse Annotations Help

PDB-101 PDB EMDataResource NAKB wwPDB Foundation PDB-IHM

Redesigned PDB Statistics Support Enhanced Functionality Explore Statistics

- Welcome
- Deposit
- Search
- Visualize
- Analyze
- Download
- Learn

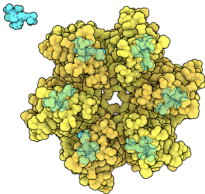
RCSB Protein Data Bank (RCSB PDB) enables breakthroughs in science and education by providing access and tools for exploration, visualization, and analysis of:

- Experimentally-determined 3D structures from the Protein Data Bank (PDB) archive
- Integrative 3D Structures from the PDB Archive
- Computed Structure Models (CSM) from AlphaFold DB and ModelArchive

**NEW** Explore Integrative Structures

PDB-101 Training Resources

March Molecule of the Month



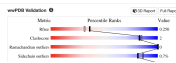
Lenacapavir

# PDB entry: 6V2F

# 6V2F.mmcif

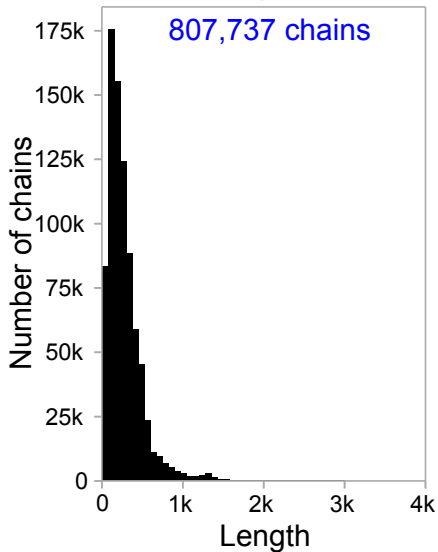
**6V2F** **pdb\_00006v2f**  
 Crystal structure of the HIV capsid hexamer bound to the small molecule long-acting inhibitor, GS-6037  
 Depositor: 2019-11-02 | Released: 2020-07-01  
 Deposition Authority: Gilead, LLC, Lila, J.O., Wei, S.H., Villano, A.G., Sommadossi, J.P., Hu, E.Y., Schneider, S.O., Chiu, T.

| ATOM | 1  | N | N   | . | PRO A | 1 | 2  | ?        | 15.01600 | -6.21200 | 10.40700 | 1.000    | 23.70517 | ? 1   | PRO A | N   | 1 |
|------|----|---|-----|---|-------|---|----|----------|----------|----------|----------|----------|----------|-------|-------|-----|---|
| ATOM | 2  | C | CA  | . | PRO A | 1 | 2  | ?        | 15.72400 | -5.12100 | 11.08300 | 1.000    | 24.88439 | ? 1   | PRO A | CA  | 1 |
| ATOM | 3  | C | C   | . | PRO A | 1 | 2  | ?        | 15.53300 | -5.12000 | 12.59400 | 1.000    | 27.78449 | ? 1   | PRO A | C   | 1 |
| ATOM | 4  | O | O   | . | PRO A | 1 | 2  | ?        | 14.76800 | -5.91700 | 13.12500 | 1.000    | 27.15848 | ? 1   | PRO A | O   | 1 |
| ATOM | 5  | C | CB  | A | PRO A | 1 | 2  | ?        | 15.18100 | -3.86600 | 10.46300 | 1.000    | 28.08434 | ? 1   | PRO A | CB  | 1 |
| ATOM | 6  | C | CG  | A | PRO A | 1 | 2  | ?        | 13.81200 | -4.32400 | 9.84400  | 1.000    | 31.53571 | ? 1   | PRO A | CG  | 1 |
| ATOM | 7  | C | CD  | A | PRO A | 1 | 2  | ?        | 14.10300 | -5.70700 | 9.37000  | 1.000    | 29.62990 | ? 1   | PRO A | CD  | 1 |
| ATOM | 8  | N | N   | . | ILE A | 1 | 3  | ?        | 16.24500 | -4.22700 | 13.27300 | 1.000    | 26.78383 | ? 2   | ILE A | N   | 1 |
| ATOM | 9  | C | CA  | . | ILE A | 1 | 3  | ?        | 16.07600 | -3.90900 | 14.70100 | 1.000    | 31.45083 | ? 2   | ILE A | CA  | 1 |
| ATOM | 10 | C | C   | . | ILE A | 1 | 3  | ?        | 15.42700 | -2.62200 | 14.85900 | 1.000    | 32.50870 | ? 2   | ILE A | C   | 1 |
| ATOM | 11 | O | O   | . | ILE A | 1 | 3  | ?        | 15.98200 | -1.60900 | 14.41600 | 1.000    | 36.22115 | ? 2   | ILE A | O   | 1 |
| ATOM | 12 | C | CB  | . | ILE A | 1 | 3  | ?        | 17.41200 | -4.85800 | 15.45700 | 1.000    | 29.32785 | ? 2   | ILE A | CB  | 1 |
| ATOM | 13 | C | CG1 | . | ILE A | 1 | 3  | ?        | 16.11400 | -5.36100 | 15.13200 | 1.000    | 31.84898 | ? 2   | ILE A | CG1 | 1 |
| ATOM | 14 | C | CG2 | . | ILE A | 1 | 3  | ?        | 17.18500 | -3.92200 | 16.95800 | 1.000    | 36.91482 | ? 2   | ILE A | CG2 | 1 |
| ATOM | 15 | C | CD1 | . | ILE A | 1 | 3  | ?        | 17.55500 | -6.58300 | 15.79800 | 1.000    | 31.74506 | ? 2   | ILE A | CD1 | 1 |
| ATOM | 16 | N | VAL | A | 1     | 4 | ?  | 14.25100 | -2.59200 | 15.47600 | 1.000    | 35.11352 | ? 3      | VAL A | N     | 1   |   |
| ATOM | 17 | C | CA  | . | VAL A | 1 | 4  | ?        | 13.52800 | -1.34300 | 15.66600 | 1.000    | 40.31307 | ? 3   | VAL A | CA  | 1 |
| ATOM | 18 | C | C   | . | VAL A | 1 | 4  | ?        | 13.18100 | -1.15000 | 17.13600 | 1.000    | 39.94655 | ? 3   | VAL A | C   | 1 |
| ATOM | 19 | O | O   | . | VAL A | 1 | 4  | ?        | 13.18200 | -2.18300 | 17.91400 | 1.000    | 41.95933 | ? 3   | VAL A | O   | 1 |
| ATOM | 20 | C | CB  | . | VAL A | 1 | 4  | ?        | 12.25800 | -1.29600 | 14.79800 | 1.000    | 40.74538 | ? 3   | VAL A | CB  | 1 |
| ATOM | 21 | C | CG1 | . | VAL A | 1 | 4  | ?        | 12.62300 | -1.30500 | 13.31900 | 1.000    | 40.23604 | ? 3   | VAL A | CG1 | 1 |
| ATOM | 22 | C | CG2 | . | VAL A | 1 | 4  | ?        | 11.34200 | -2.46300 | 15.13200 | 1.000    | 46.80582 | ? 3   | VAL A | CG2 | 1 |
| ATOM | 23 | N | N   | . | MET A | 1 | 11 | ?        | 14.53600 | -2.31000 | 21.90800 | 1.000    | 44.72441 | ? 10  | MET A | N   | 1 |
| ATOM | 24 | C | CA  | . | MET A | 1 | 11 | ?        | 15.02700 | -2.68200 | 20.56800 | 1.000    | 42.41279 | ? 10  | MET A | CA  | 1 |
| ATOM | 25 | C | C   | . | MET A | 1 | 11 | ?        | 14.76900 | -4.16300 | 20.31000 | 1.000    | 41.99084 | ? 10  | MET A | C   | 1 |
| ATOM | 26 | O | O   | . | MET A | 1 | 11 | ?        | 15.33700 | -5.82500 | 20.98100 | 1.000    | 40.16552 | ? 10  | MET A | O   | 1 |
| ATOM | 27 | C | CB  | . | MET A | 1 | 11 | ?        | 16.51000 | -2.36600 | 20.38300 | 1.000    | 42.16612 | ? 10  | MET A | CB  | 1 |
| ATOM | 28 | C | CG  | . | MET A | 1 | 11 | ?        | 16.96900 | -8.90000 | 20.65200 | 1.000    | 45.27050 | ? 10  | MET A | CG  | 1 |
| ATOM | 29 | S | SD  | . | MET A | 1 | 11 | ?        | 15.98700 | 0.24500  | 19.59100 | 1.000    | 54.71942 | ? 10  | MET A | SD  | 1 |
| ATOM | 30 | C | CE  | . | MET A | 1 | 11 | ?        | 16.68100 | -0.17700 | 17.99500 | 1.000    | 46.37049 | ? 10  | MET A | CE  | 1 |



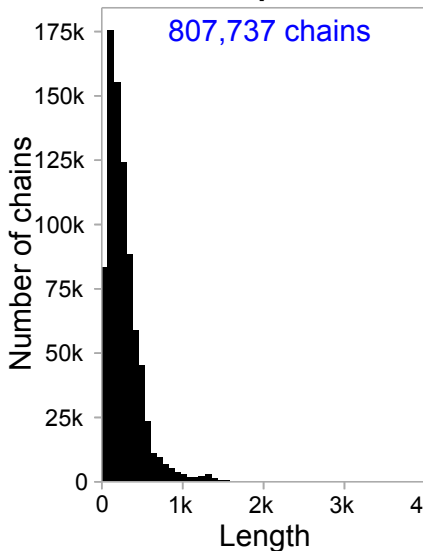
# AF2 Training data (2024)

PDB proteins

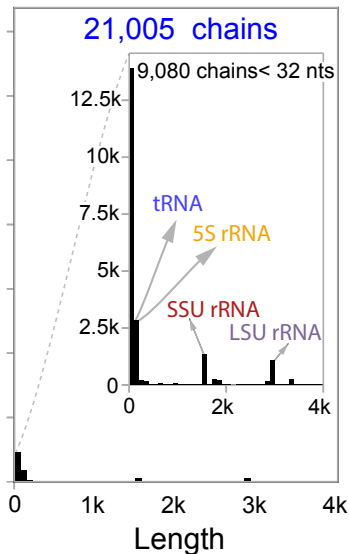


# AF2 Training data

PDB proteins



PDB RNAs



# AF2 losses

## 1. **evoformer losses**

1.1 Distograms

1.2 MSA loss

## 2. **structural losses**

2.1 FAPE = Frame aligned position error

2.2 Torsion angles loss

2.3 Model confidence prediction (pLDDT) local distance difference test

# Distogram loss

To associate the pair representation for the structural model

## supervised training

Pair representation  $z_{ij} + z_{ji}$

one-hot 64 distance bins covering 2-22 Amstrongs

$$p_{ij} = \text{softmax}(\text{one-hot}(z_{ij} + z_{ji}, 64\text{bins}))$$

$$y_{ij} = \text{softmax}(\text{one-hot}(\text{distance}_{ji}, 64\text{bins}))$$

$$Loss_{dist} = -\frac{1}{r^2} \sum_{b=1}^{64} y_{ij}^b \log p_{ij}^b$$

# MSA loss

To encourage inter-sequence and phylogenetic relationships

**self-supervised training**

$$m_{si}[C_m] \rightarrow p_{si}[23]$$

where  $23 = 20aa + unknown + gap + mask$

$y_{si}[23]$  are the one-hot ground-truth values in the msa.

$$Loss_{msa} = -\frac{1}{N_{masked}} \sum_{s,i \in mark} \sum_{c=1}^{23} y_{si}^c \log p_{si}^c$$

# FAPE loss (Frame aligned point error)

## Structural loss

predicted coordinates of a atom  $\bar{x}_j$  under a predicted local frame  $T_i$

$$T_i^{-1}\bar{x}_j = R_i^{-1}(\bar{x}_j - \bar{t}_i)$$

compared to the true ones

$$(T^{true})_i^{-1}\bar{x}_j^{true}$$

$$d_{ij} = \|T_i^{-1}\bar{x}_j - (T^{true})_i^{-1}\bar{x}_j^{true}\|^2$$

### AF2 Losses

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

# Training tricks

Recycling (3 times)

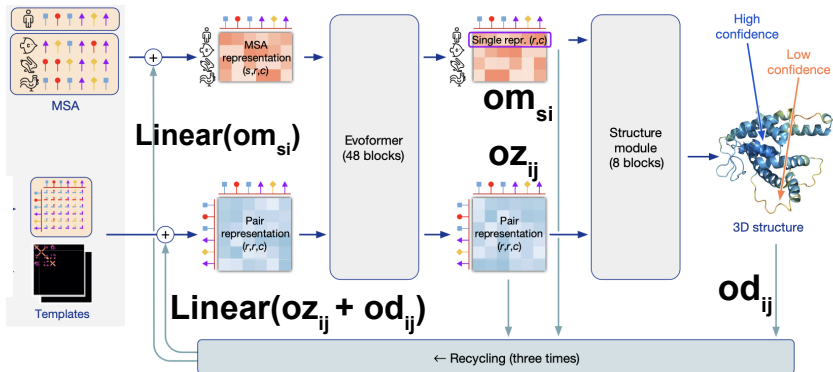
Iterative refinement using the whole network

Self-distillation

# Training tricks

Iterative refinement using the whole network

**Recycling (3 times)**  
(no losses applied during recycling)



# Training tricks

self-distillation

training with labeled and unlabeled data

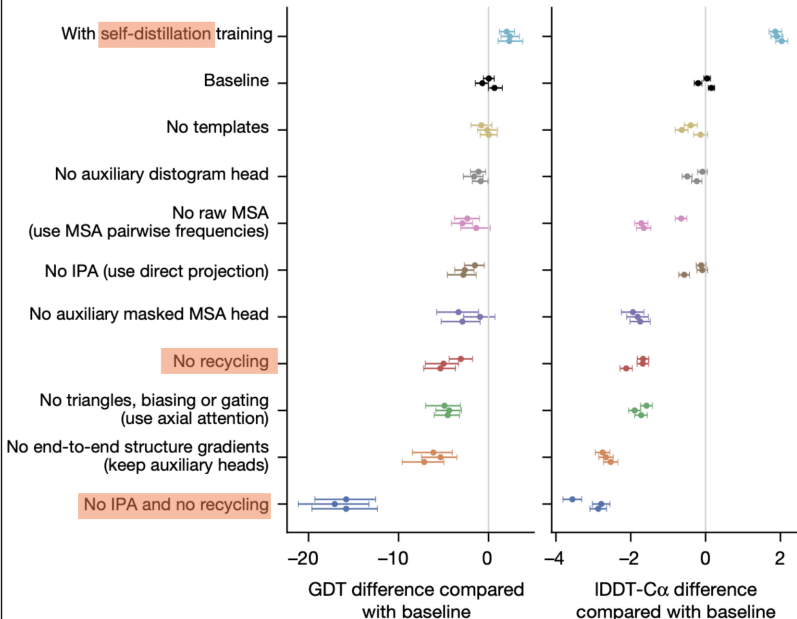
- ▶ Train AF2 supervised on proteins with PDB structures (labeled data)
- ▶ Predict structures of proteins without known structures (unlabeled data)
- ▶ Retrain model from scratch on  
(proteins with known PDB structures) +  
(proteins with high-confidence predictions)

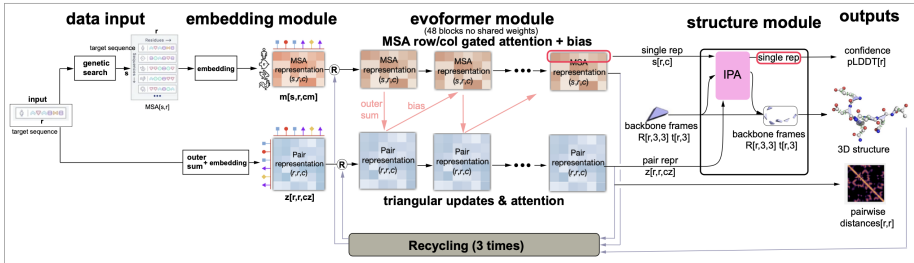
effective use of the unlabelled sequence data  
considerably improves the accuracy of the resulting network

# AlphaFold2 ablation results

Test set of CASP14 domains

Test set of PDB chains





Some important characteristics of AlphaFold2 are

- ▶ **Directly predicts the 3D coordinates of all heavy atoms** from the amino acid sequence using aligned sequences of homologs as inputs.
- ▶ Uses **alignments  $m[r,s]$  and pair representations  $z[r,r]$**  with different attention mechanisms for each of them, but treated jointly such that they exchange information with each other.
- ▶ **trained end-to-end**: all parameters are updated in the one single gradient descent optimization routine
- ▶ Uses many **different loss functions** some used at intermediate stages, and a final global loss function used to optimize all parameters of the model at once.
- ▶ Uses **self-distillation** to learn from unlabeled proteins. Self-distillation uses predictions of the model for further training.
- ▶ It is able to **estimate accuracy**
- ▶ Introduces a **new formalism for protein 3D structure prediction**: a rotation and translation to identify the location of each amino acid in a protein molecule. All residues are trivially initialized at the same location which evolves with training into an accurate description of all atoms of the amino acid chain.