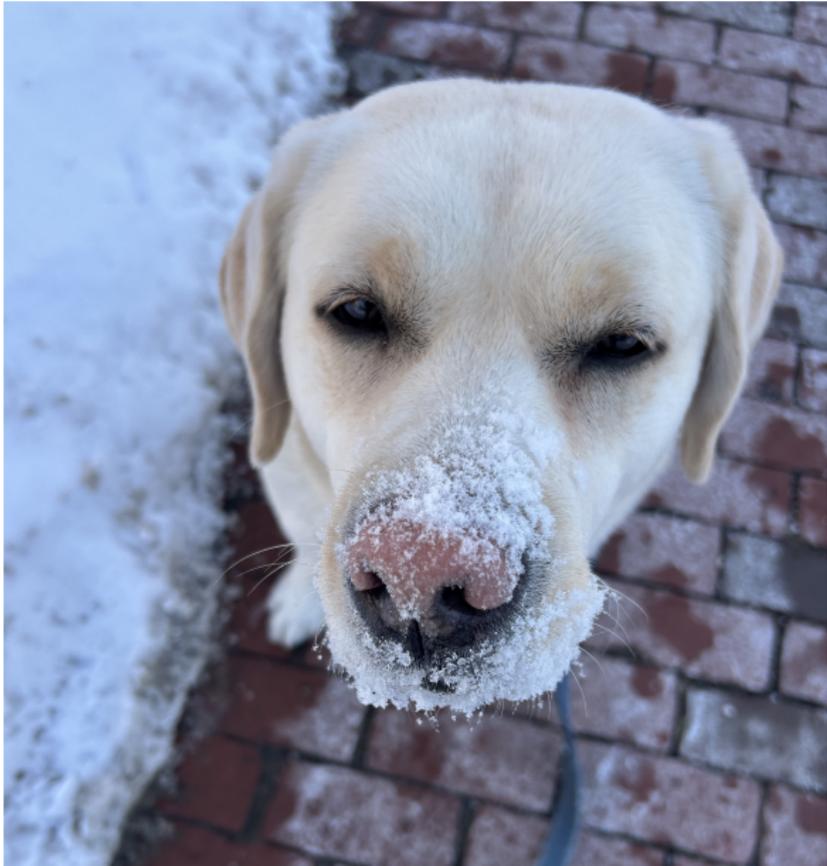


block b3

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

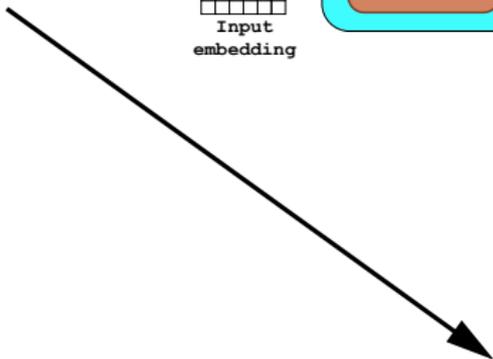
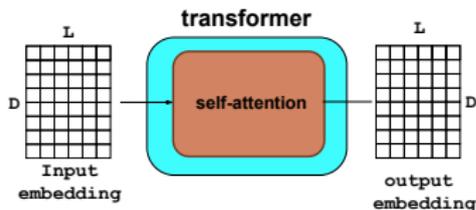
Illia Polosukhin* ‡
illia.polosukhin@gmail.com

3 transformers for natural language translation

self-attention



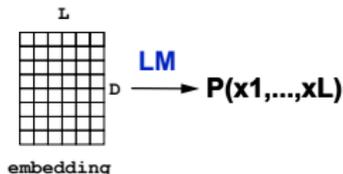
Transformers



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)

$$x_1 \sim P(x_1)$$

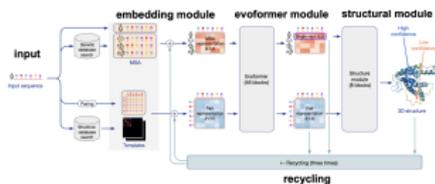
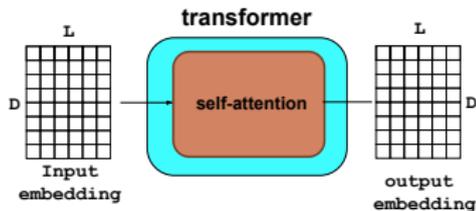
$$x_2 \sim P(x_2 \mid x_1)$$

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

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

[b3] self-attention

[b3] Transformers

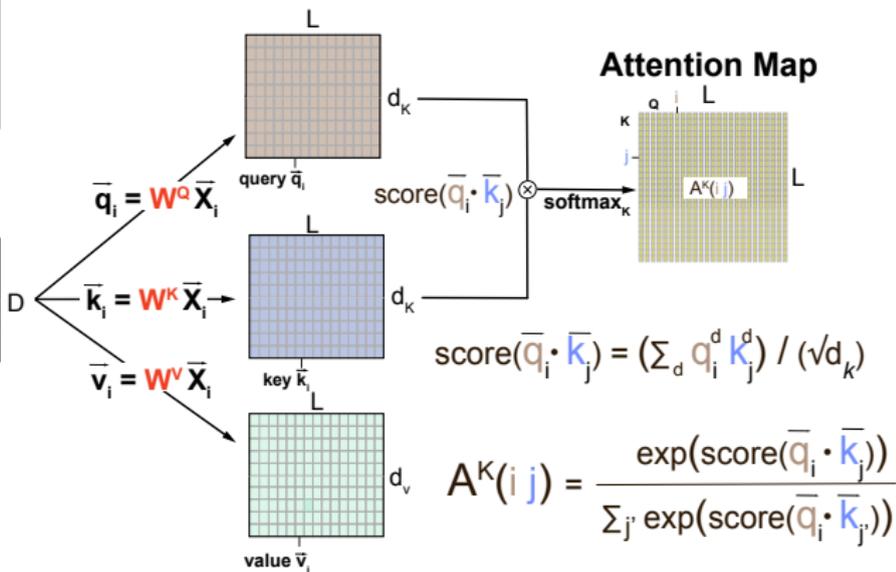
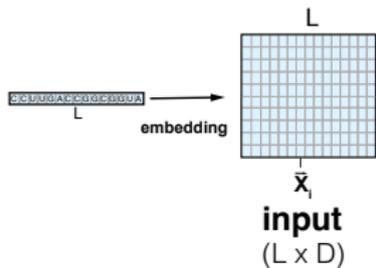


[b3] AlphaFold2

[b4] Language Models (LLMs)

Self-Attention

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



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** $(q \cdot k)[1]$

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

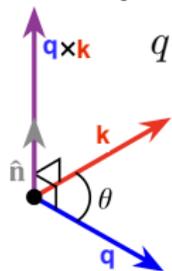
2. **element-wise product** $(q \odot k)[D]$

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

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

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

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



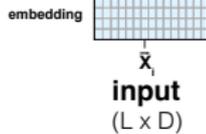
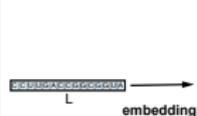
$$q \times k = \|q\| \|k\| \sin(\theta) \hat{n}$$

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

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

Self-Attention (1d input)

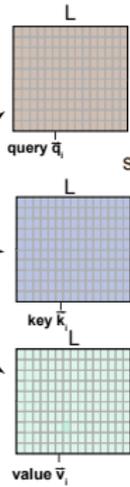
parameters	size
W^Q	$d_k \times D$
W^K	$d_k \times D$
W^V	$d_v \times D$



$$\bar{q}_i = W^Q \bar{X}_i$$

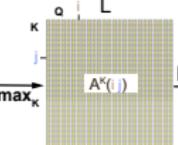
$$\bar{k}_i = W^K \bar{X}_i$$

$$\bar{v}_i = W^V \bar{X}_i$$



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

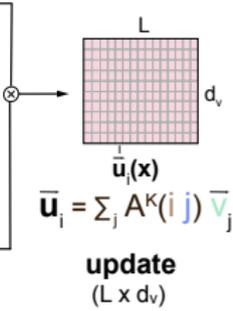
Attention Map



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

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

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



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

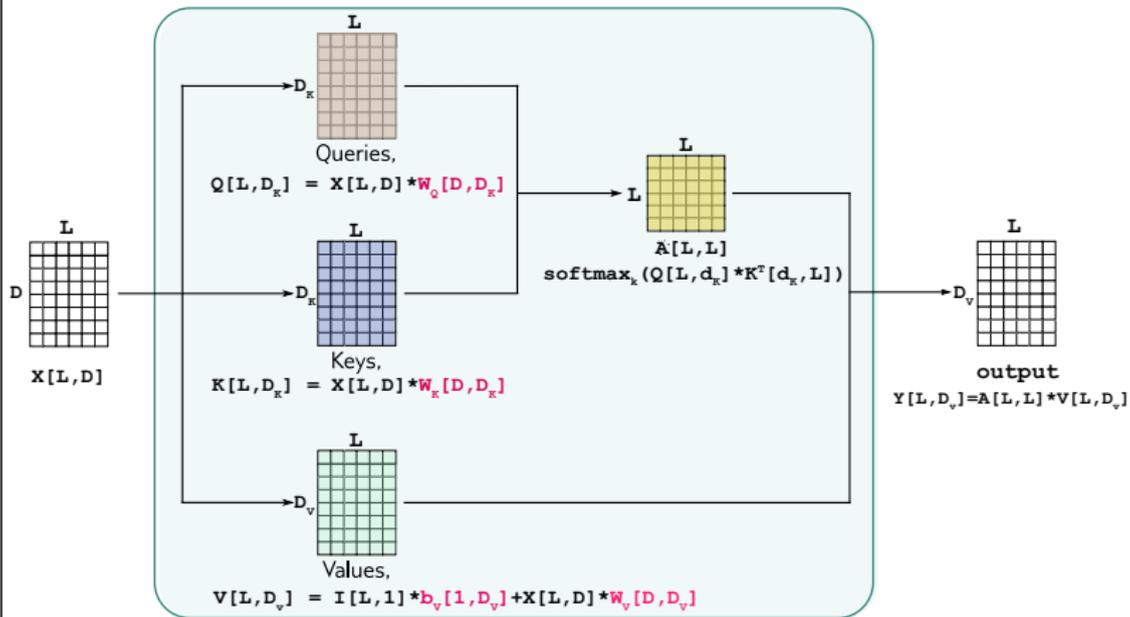
        # q[B, L, d_model]
        # k[B, L, d_model]
        # v[B, L, d_model]
        q = self.q_proj(x)
        k = self.k_proj(x)
        v = self.v_proj(x)

        # 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)

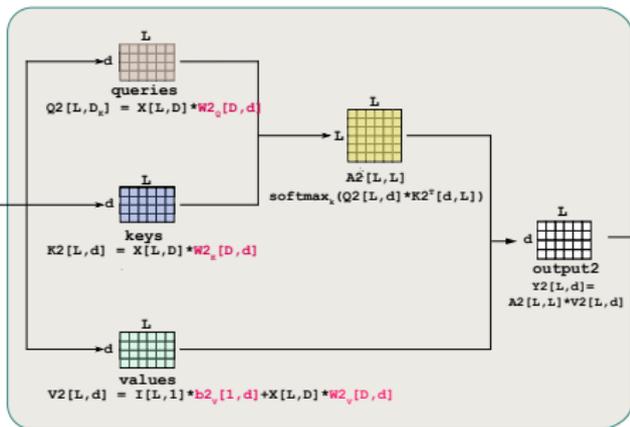
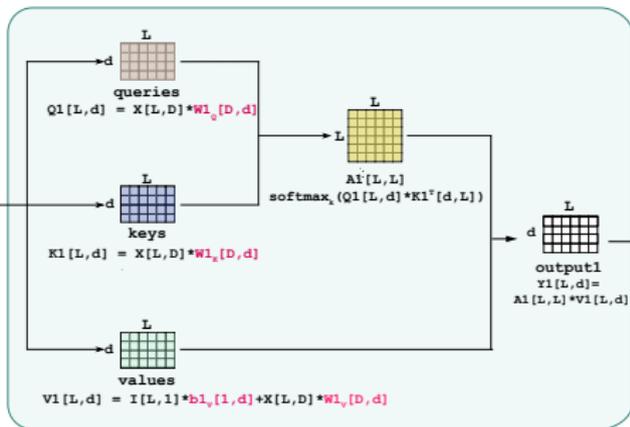
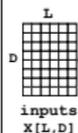
        return out
```

sequence input = X[B, L, D]

Self Attention (matrix view)



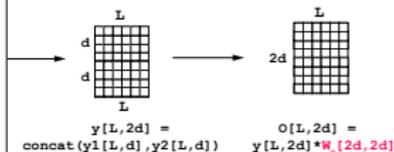
adapted from: UDL, Prince



$D = d_model$

$H = n_heads = 2$

$d = d_head = d_model / n_heads$



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)

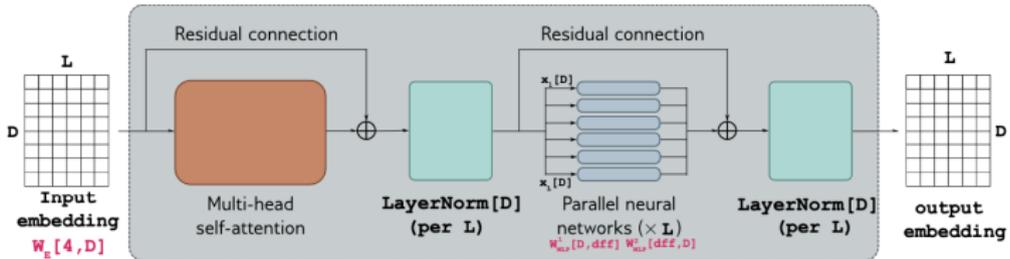
        # concatenate all heads together
        # out.transpose(1, 2)[B, L, H, d_head]
        # out[B, L, H*d_head]
        out = out.transpose(1, 2).reshape(B, L, D)

        # 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 δ and standard deviation γ .

Positional encoding

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

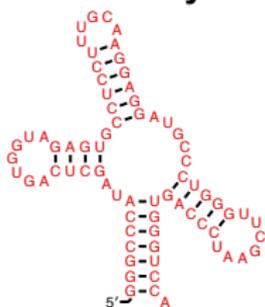
$$dij = j - i$$

$$pij = \text{LinearLayer}(\text{onehot}(dij, vbins))$$

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' A	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

```

GGGCCCAUAGCCUCAG--GG--AGAGUGCCUCCUUGCAAGGAGAUCC--CCUG--GGUCCGAADCCCAGGGGUCCA
CAUUGAUGACCGAA--AG--CAAGUACGGGUCUUAACAACUUUAA--UAGUAAUUAGCACUUAUCUUAAAG
GGCCCCGUGACCAADUGGADGAGCGUUGACUACGGAAACAAAGGG--UAGG--GGUUCGACUCCUCCGGGGCCG
GGCCUUGAGCCAGCGUGGUGAGGCGACCCUUAUAGCGGAGGG--CGG--GGUUCGAGUCUCCUCCGGCCCA
AGGGCCAUUGUUAA--GG--AGAACAGAGGUCUCCAAACCUCGG--UGUG--GGUUCGAGUCUCCUCCGGCC
AGGGCCGUGUUAUUUGGU--AGAGCACCGGUCUCCAAACCGGGUG--UGGG--AGUUCGAGUCUCCUCCGGCC
CGUGAUUGCCUCAGUUGGU--AGAGCCACCCUUGGUAAGGGUGAGGU--CCCC--AGUUCGACUCCUCCGGUUCAGCA
GGGGCCGUGGCCAAGG--GG--AAGGCAGCGGDUUGGUCUCCUUAUCU--CGGA--GGUUCGAAUCCUCCUCCGGCC
GCCUUUAACACAGG--GG--AGAGUUAUGCCUUGGUAAGCADAAG--GAG--GGUUCGAAUCCUCCUCCGGCC
GCCUCCUUGGCUAAG--GGAGAGGACAGAGGUCUUAACAACUUUGG--UADA--GGUUCGAAUCCUCCUCCGGCC
GCACCCAGUGGCUAAG--GGUUAAGCGCCACUUAUUGGUAUUUGGC--GGUUCGAAUCCUCCUCCGGAGCA
AGUAAGGUCAGCUAA--UU--AAGCUAUCGGGCCCAUACCCCGAATA--CGUU--GGUUAUAAUCCUCCUCCGACUA
AGAUUAUAGUAAAAD--CA--AUUACAUAACUUDGUCAAAGUUAUU--UADA--GAUCAAUAUCUUAUUAUCUUA
GUUUCUGAGUGAA--UU--ACAACGADGUUUDCAUGUCUADUGG--UCGC--AGUUGAAGGCGUGGUAAGAAUA
GGUUAUAGGGGAA-----UUACAUGCCAAUUGCAAUUCAGAG--UGGU--GAGAAAU--CUUACUAGAGAA
ACUCCUUUGUAUA--UU--AAUAUAACUGACUCCAAUUAUAGUA--UUUU--GAUUAAC--CCAGAAAGAGUA
UGAUUGAAGCCAGUA--AU--AGGGUUAUGCGUGUUAACUAAUUUU--CGUA--GGUUAUAAUCCUCCUCCUCC
CACUADGAAGCUA-----AGAGCGUUAACCUUUUAAGUUAAGUU--AGAG--ACCUAUA--AUUCCUUAAGUA
GUUAUUGAGCUUAUAAC--AAAGCAAGCACUGAAAUUGCUUAGA--UGGA--U--AAUUGU--AUCCUUAUAACA
ACUUAUAGGAAUAU--AG--AAUUCUCCUUGGUAAGCAAA--CCUU--GGUCCAAUCCUUAUUAAGUA
GGCCUUAAGCCAGCGGUGAGUGGACCCUCCUUAAGGGUGAGG--CACA--AGUUCGAAUCCUCCUCCUCCUCC
    
```

evolutionary
information

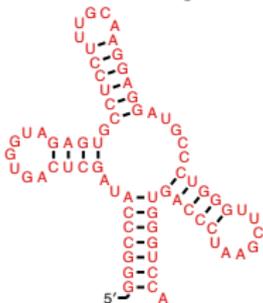
genomes, homology

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

RNA secondary structures

structural
information

X-ray, NMR, CryoEM



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

maximum likelihood training

RNA alignments

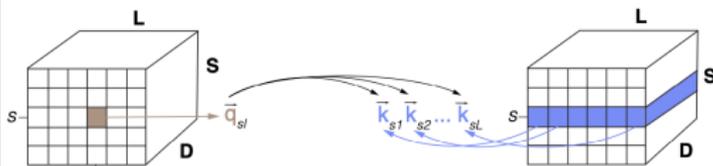
evolutionary
information

genomes, homology

G..CA
C..GA
G..CG
G..CA
A..UG
A..UG
G..CA
U..AG
G..CU
G..CA
G..CA
A..UA
A..UA
G..UA
G..UU
A..UA
U..AG

MSA attention by row/column

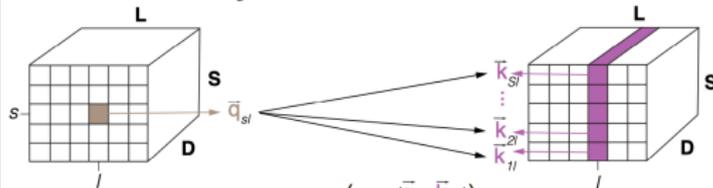
self-attention by row



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

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

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)

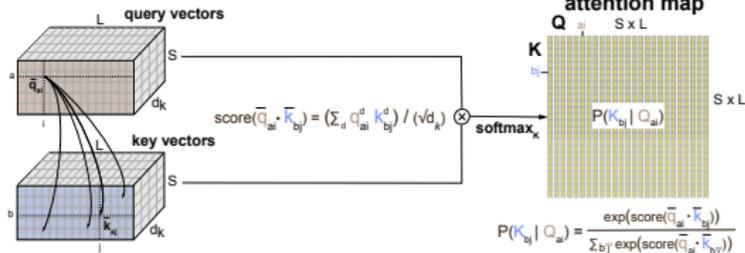
    def forward(self, x):
        # x: (B, S, L, d_model)
        B, S, L, D = x.shape
        x = x.reshape(B * S, L, D)
        x = self.attn(x)
        return x.reshape(B, S, L, D)

# 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)

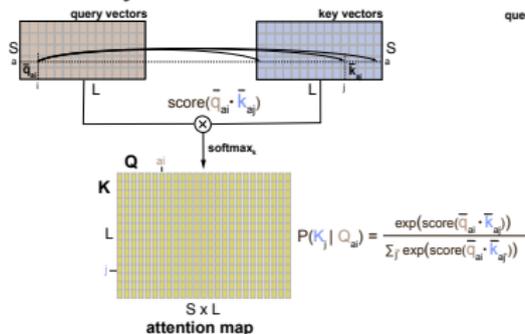
    def forward(self, x):
        # x: (B, S, L, d_model)
        B, S, L, D = x.shape
        x = x.transpose(1, 2).reshape(B * L, S, D)
        x = self.attn(x)
        return x.reshape(B, L, S, D).transpose(1, 2)
```

Alignment Self-Attention

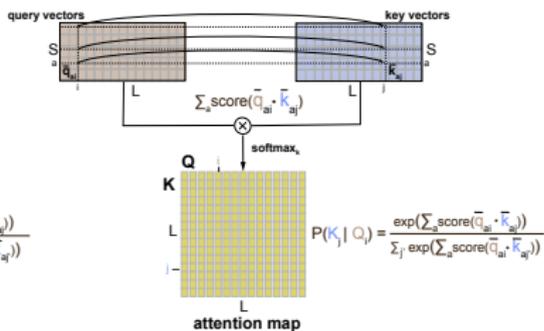
parameters	size
W^Q	$D \times d_q$
W^K	$D \times d_k$
W^V	$D \times d_v$



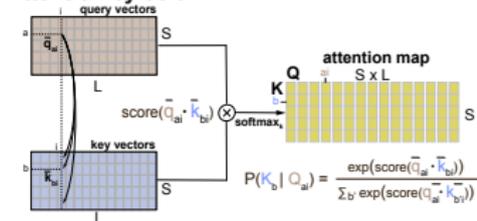
Attention by row



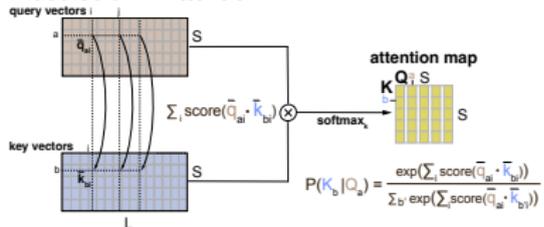
tied row Attention



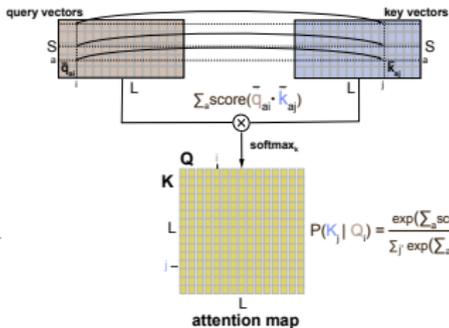
Attention by column



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)

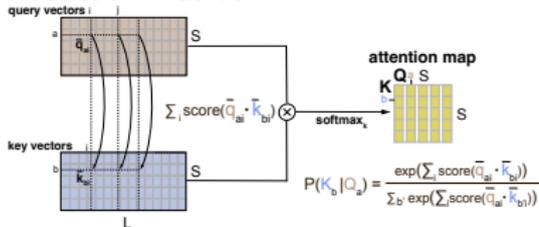
        # out [B, S, H, L, DH] = TR_attn [B, S, H, L, L] * V [B, S, H, L, DH]
        out = torch.matmul(TR_attn, V) # (B, S, H, L, DH)

        # 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("typed-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)

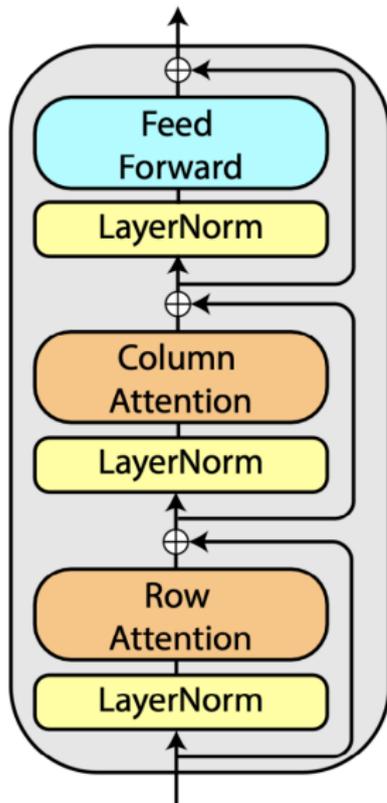
        # project
        out = self.o_proj(out)

        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

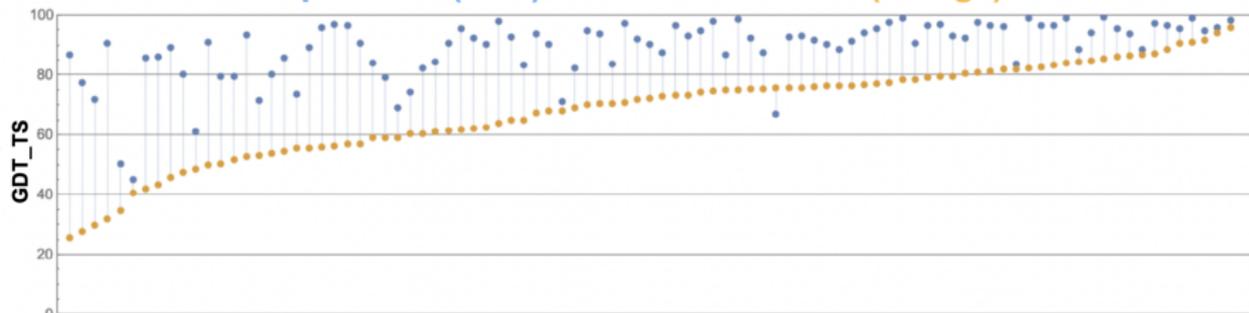
Published online: 15 July 2021

Open access

 Check for updates

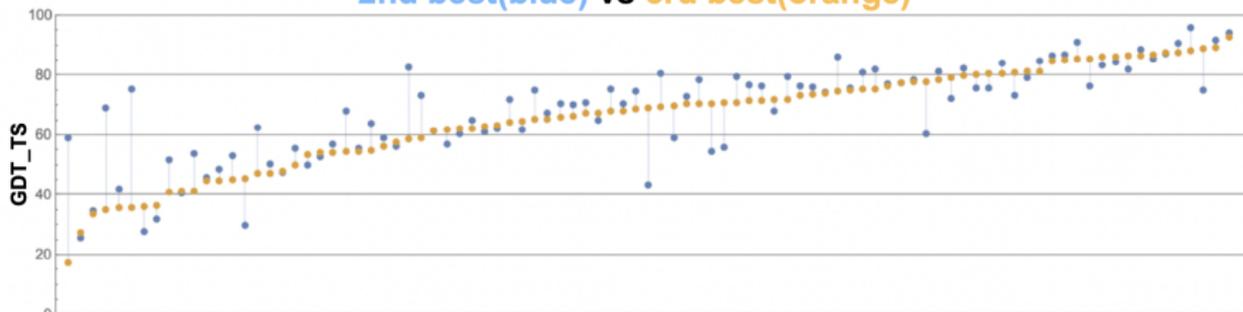
John Jumper^{1,4}, Richard Evans^{1,4}, Alexander Pritzel^{1,4}, Tim Green^{1,4}, Michael Figurnov^{1,4}, Olaf Ronneberger^{1,4}, Kathryn Tunyasuvunakool^{1,4}, Russ Bates^{1,4}, Augustin Židek^{1,4}, Anna Potapenko^{1,4}, Alex Bridgland^{1,4}, Clemens Meyer^{1,4}, Simon A. Kohl^{1,4}, Andrew J. Ballard^{1,4}, Andrew Cowie^{1,4}, Bernardino Romera-Paredes^{1,4}, Stanislav Nikolov^{1,4}, Rishub Jain^{1,4}, Jonas Adler¹, Trevor Back¹, Stig Petersen¹, David Reiman¹, Ellen Clancy¹, Michal Zielinski¹, Martin Steinegger^{2,3}, Michalina Pacholska¹, Tamas Berghammer¹, Sebastian Bodenstein¹, David Silver¹, Oriol Vinyals¹, Andrew W. Senior¹, Koray Kavukcuoglu¹, Pushmeet Kohli¹ & Demis Hassabis^{1,4}

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



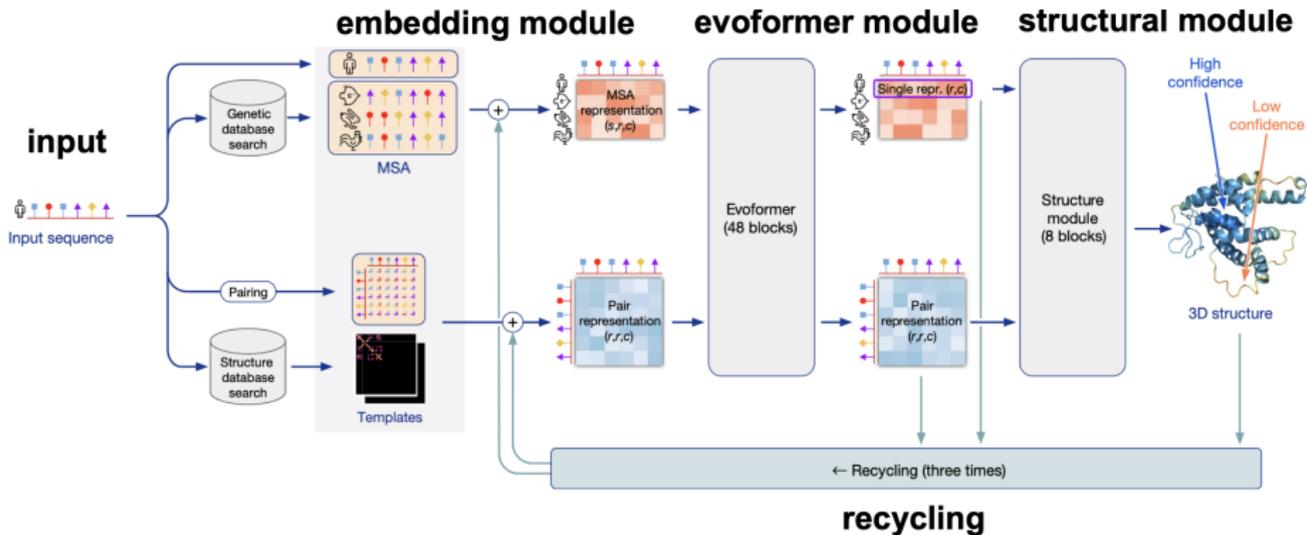
CASP14 targets

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

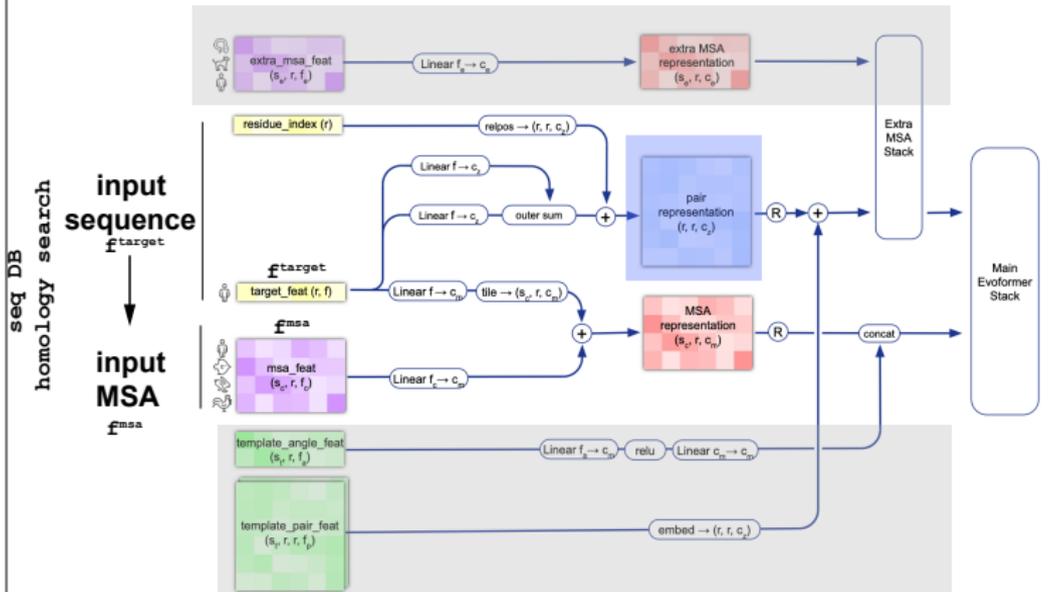


CASP14 targets

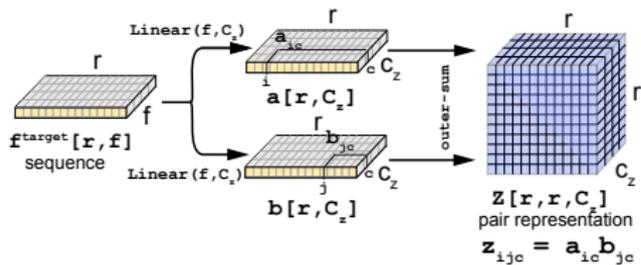
AlphaFold2 Summary (Figure 1)



AF2 Embedding Module



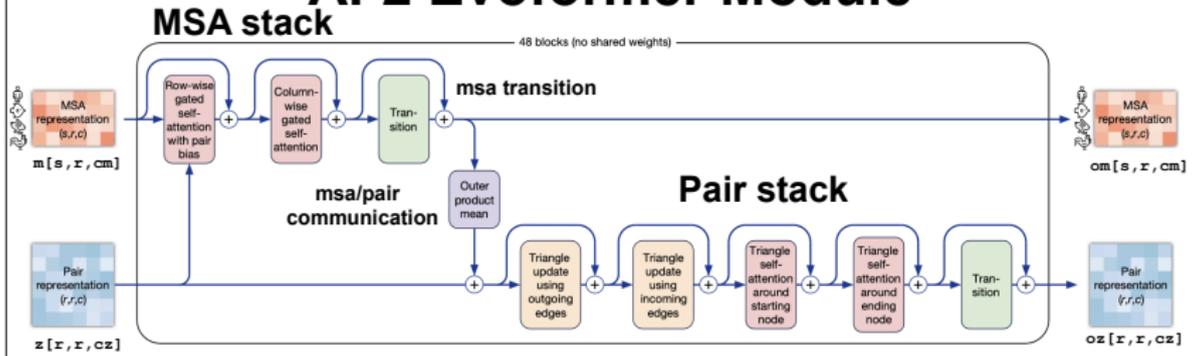
from seq to pair rep by "outer-sum"



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

$$\begin{aligned}
 a_i &= \text{Linear}(f_i^{\text{target}}) & a_i[C_2], b_i[C_2] \quad C_2 = 128 \\
 b_i &= \text{Linear}(f_i^{\text{target}}) \\
 z_{ij} &= a_i + b_j \quad (\text{outer sum}) & z_{ij}[C_2] \\
 m_{s_i} &= \text{Linear}(f_{s_i}^{\text{msa}}) + \text{Linear}(f_i^{\text{target}}) & m_{s_i}[C_m] \quad C_m = 256 \\
 \text{return } & \{m_{s_i}\}, \{z_{ij}\}
 \end{aligned}$$

AF2 Evoformer Module



MSA stack

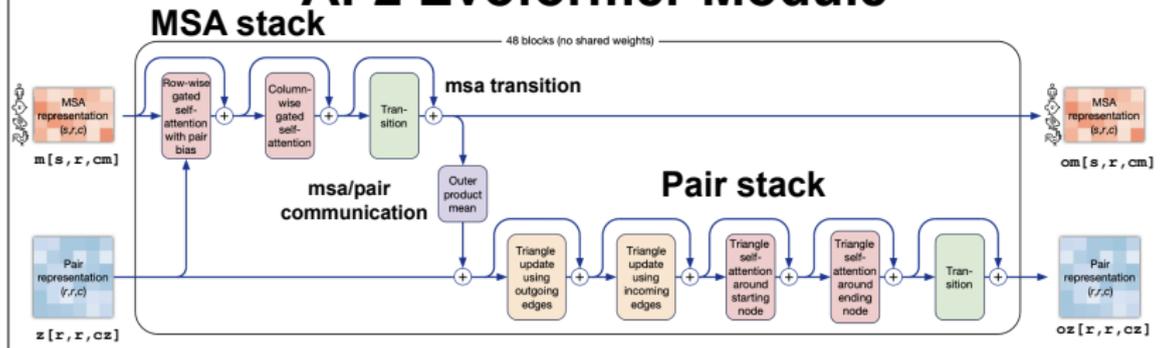
RowAttentionBias (m, z) :

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

ColAttention (m, z) :

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

AF2 Evoformer Module



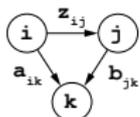
Pair stack

triangle updates

(parameters not shared)

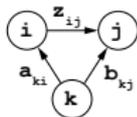
$$a_{i,j}, b_{i,j} = \text{sig}(\text{Linear}(z_{i,j})) \text{Linear}(z_{i,j})$$

outgoing



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

ingoing



$$o_{i,j} = \text{Linear}(\text{sum}_k(a_{k,i} b_{k,j}))$$

triangle attentions

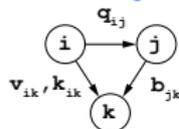
$a[r, r, r]$

(parameters not shared)

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

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

starting

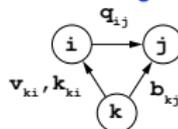


$$a_{i,j}^h = \text{softmax}_k(q_{i,j}^h k_{i,k}^h + b_{j,k})$$

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

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

ending

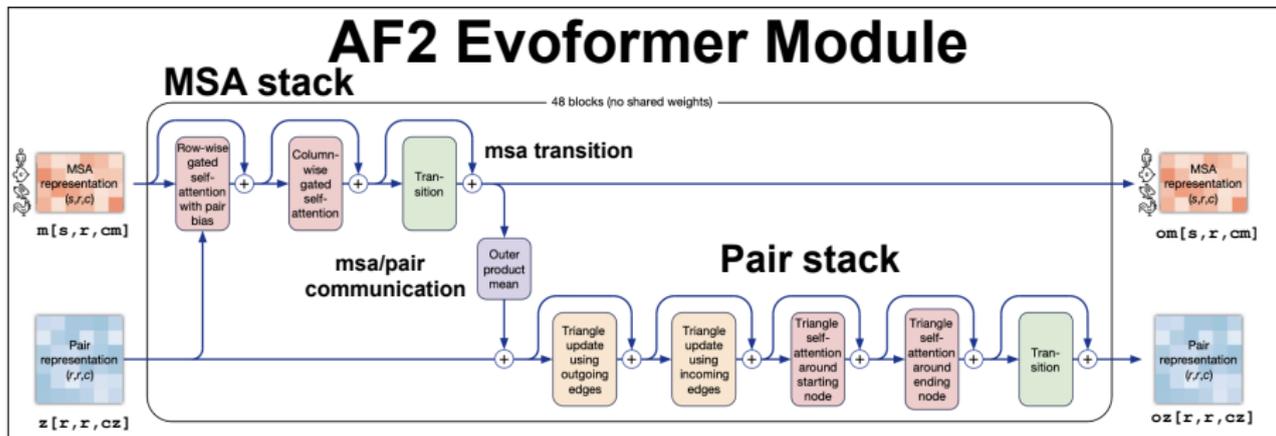


$$a_{i,j}^h = \text{softmax}_k(q_{i,j}^h k_{i,k}^h + b_{k,j})$$

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

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

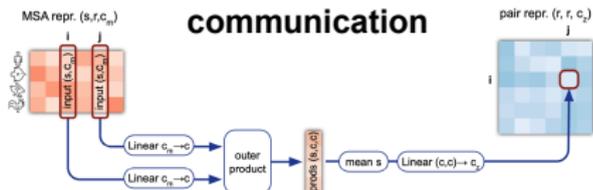
AF2 Evoformer Module



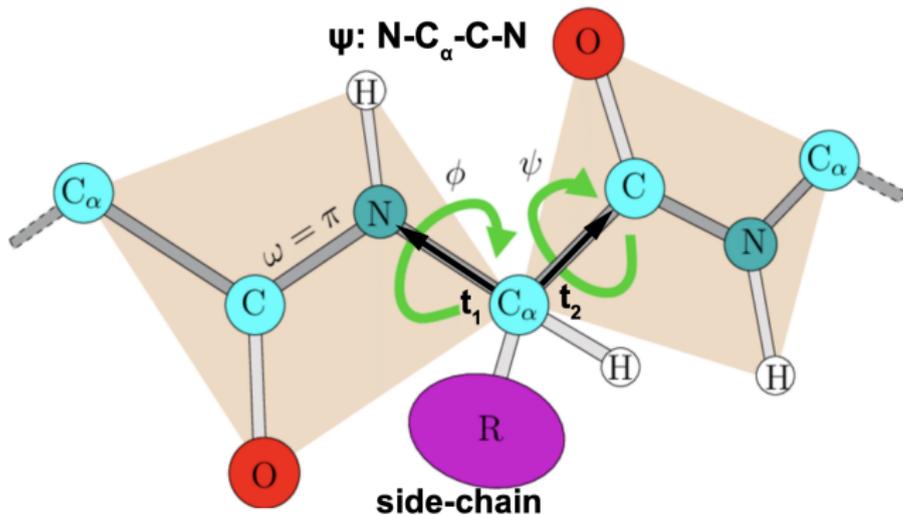
msa transition



communication



amino acid backbone reference frame



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

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

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

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

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

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

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

1 Reference frame: N-C_α-C

1 backbone torsion angle ψ to determine O

4 side-chain torsion angles

AF2 Losses

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

AlphaFold2 ablation results

Test set of CASP14 domains

Test set of PDB chains

