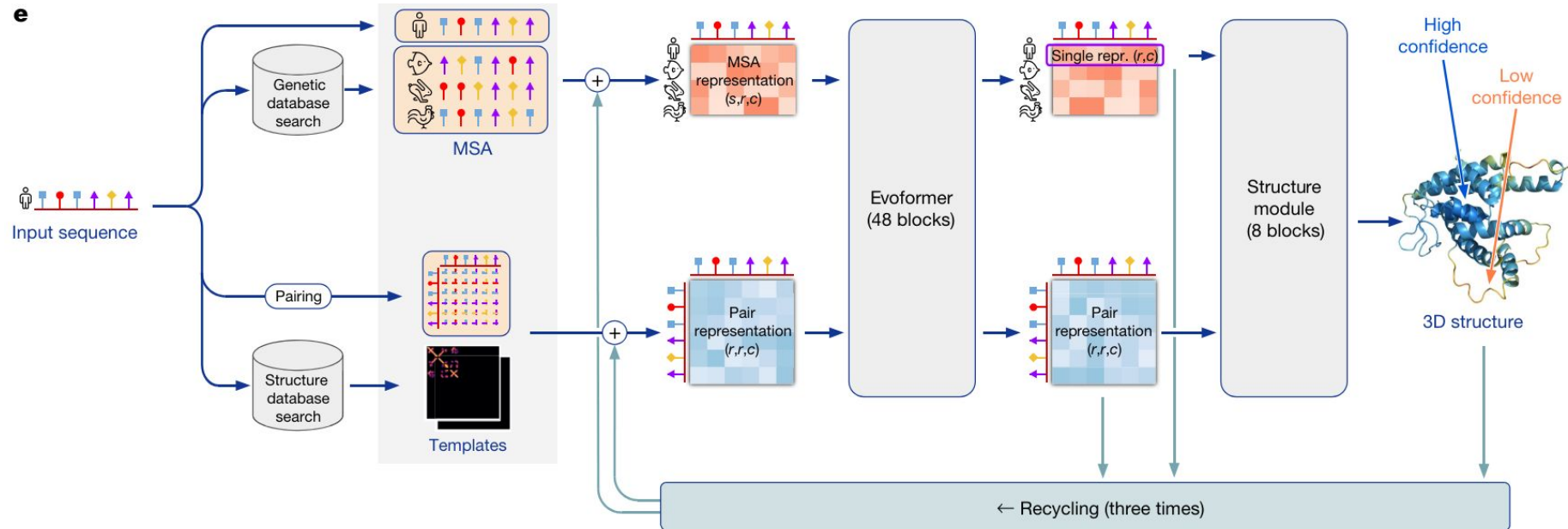
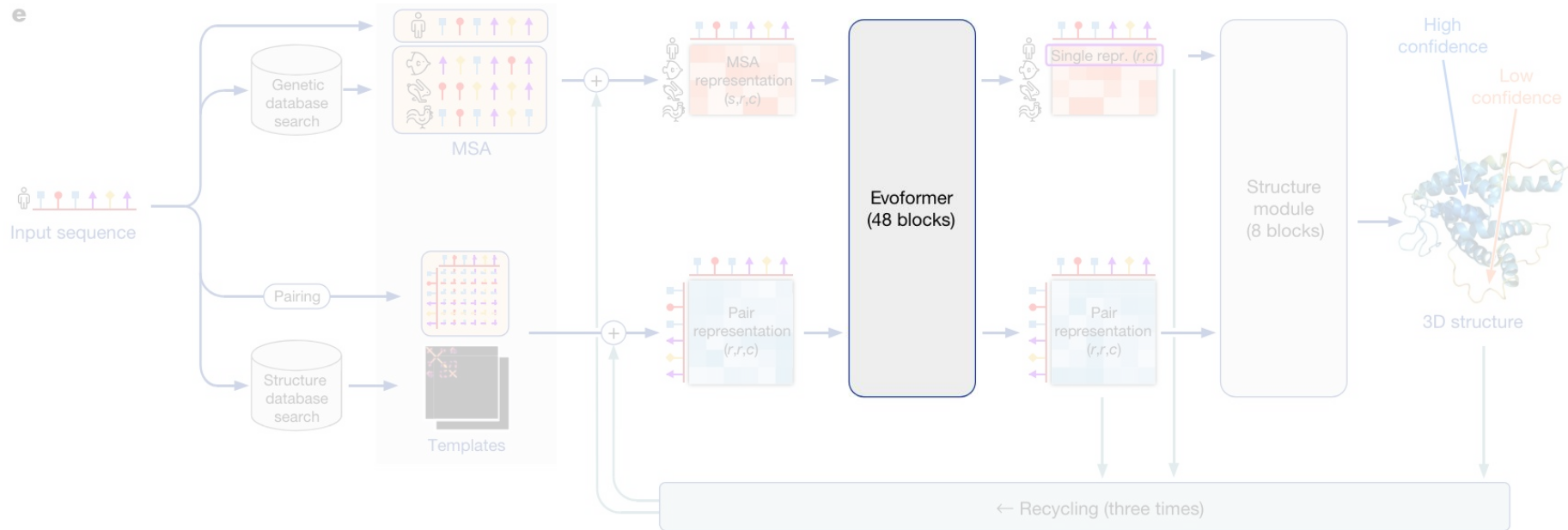


# EvoFormer

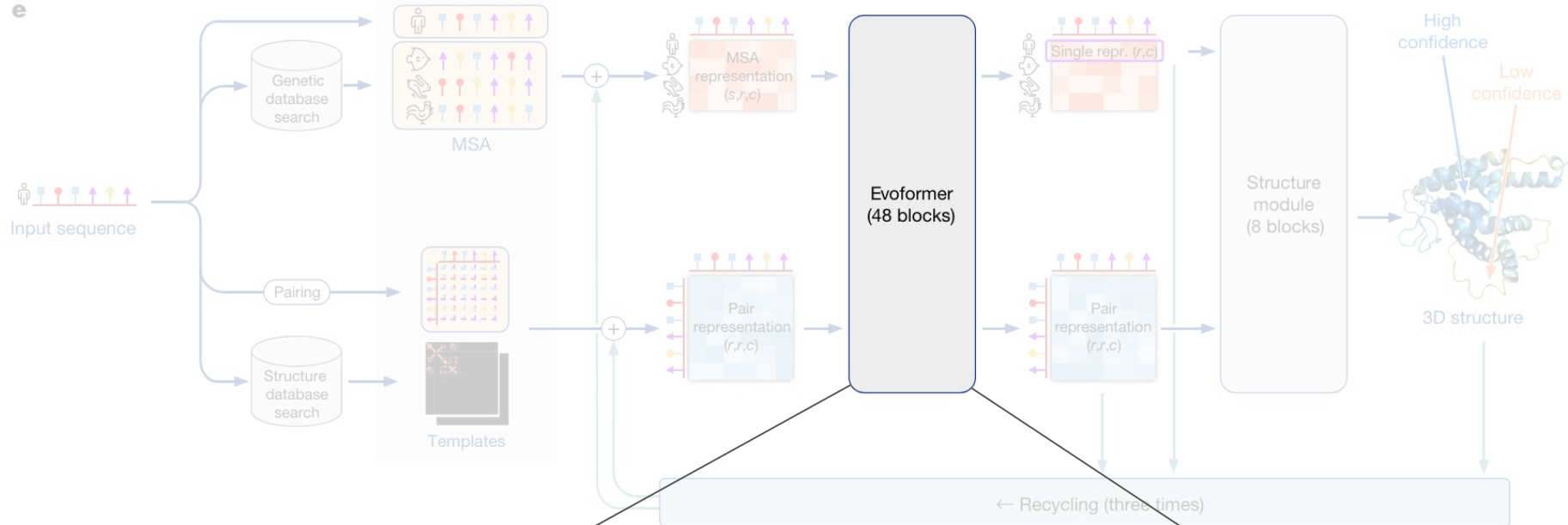
Alpha Fold 2 in details

**e**

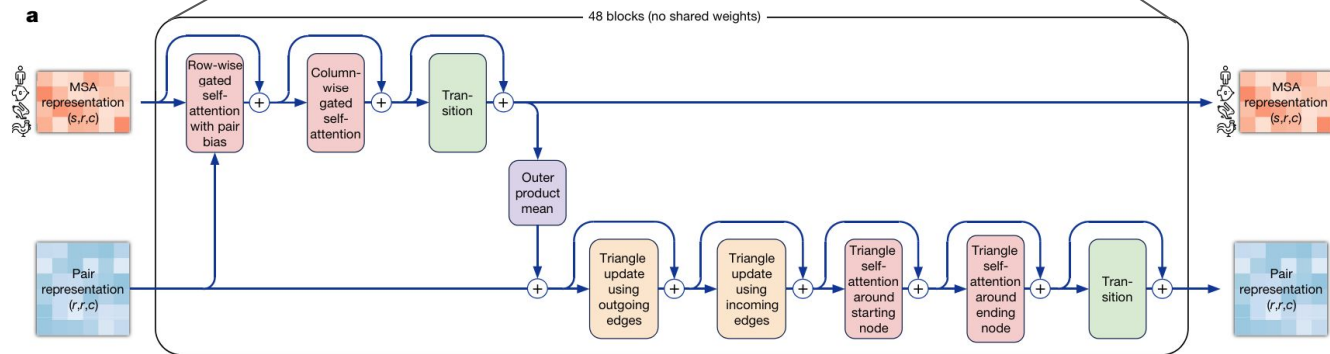
e



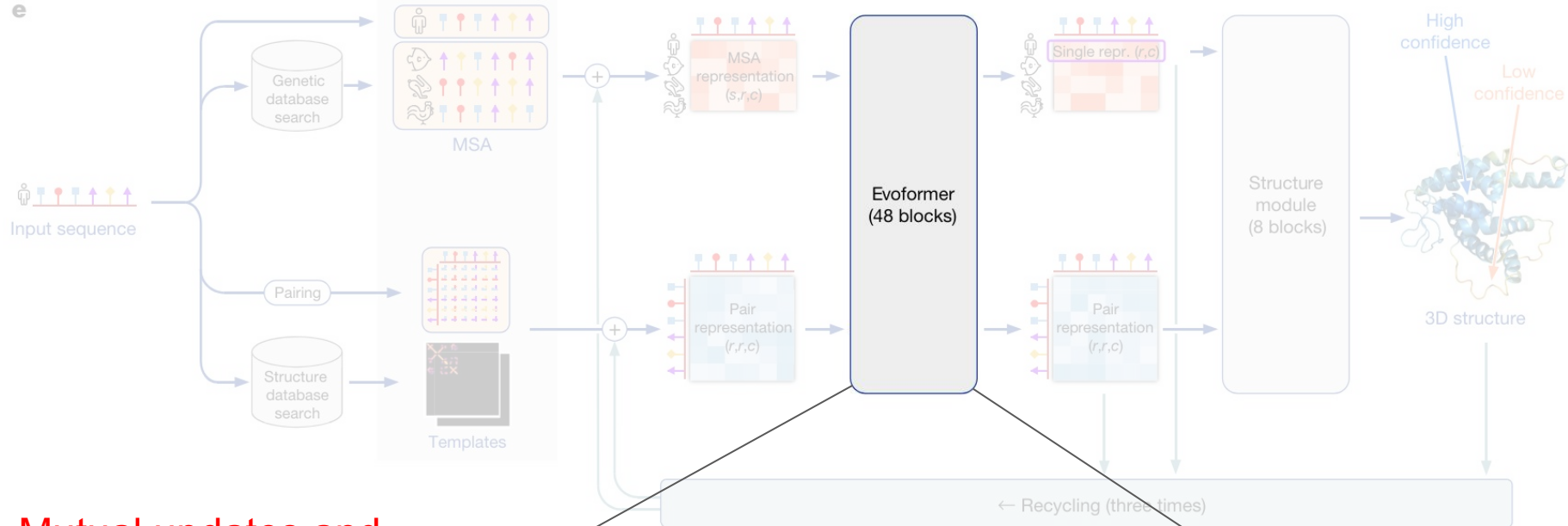
e



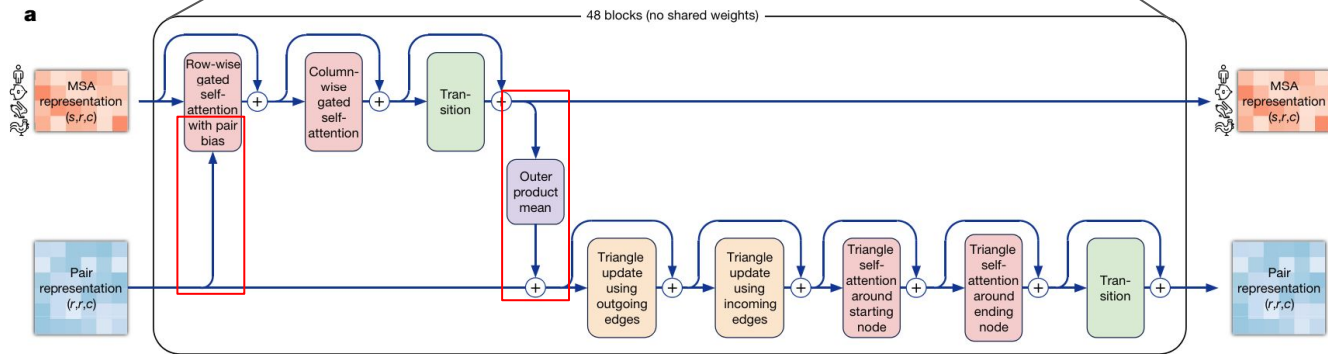
a



e



## Mutual updates and communication

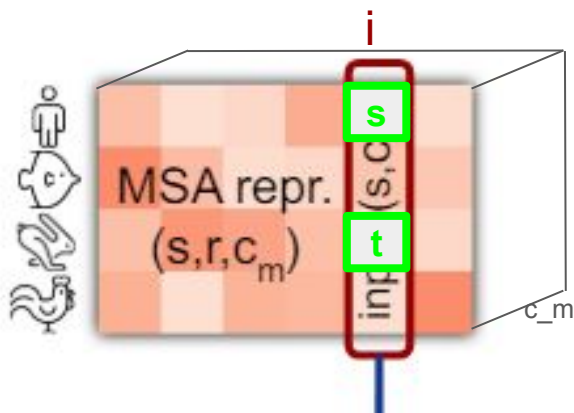


# Notation

We use  $\odot$  for the element-wise multiplication,  $\otimes$  for the outer product,  $\oplus$  for the outer sum, and  $\mathbf{a}^\top \mathbf{b}$  for the dot product of two vectors. Indices  $i, j, k$  always operate on the **residue dimension**, indices  $s, t$  on the **sequence dimension**, and index  $h$  on the attention heads dimension. The channel dimension is implicit and we type the **channel-wise vectors in bold**, e.g.  $\mathbf{z}_{ij}$ . Algorithms operate on sets of such vectors, e.g. we use  $\{\mathbf{z}_{ij}\}$  to denote all pair representations.

## Note: Gated Self-Attention

A **gating vector  $\mathbf{g}$**  is used to further modulate attention scores, effectively controlling how much influence each element should have

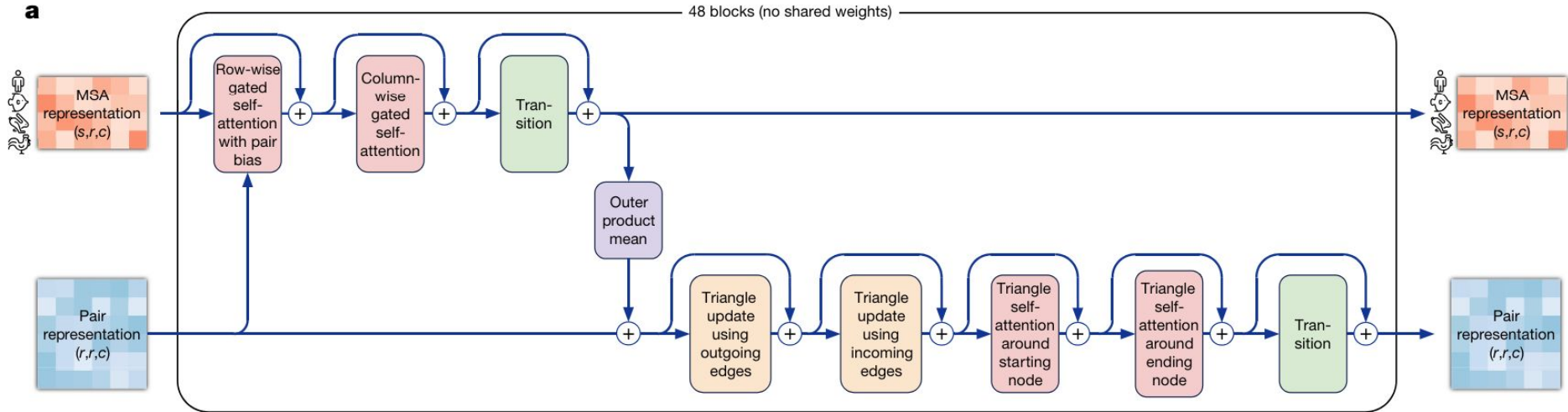


$$a_{sti}^h = \text{softmax}_t \left( \frac{1}{\sqrt{c}} \mathbf{q}_{si}^h{}^\top \mathbf{k}_{ti}^h \right)$$

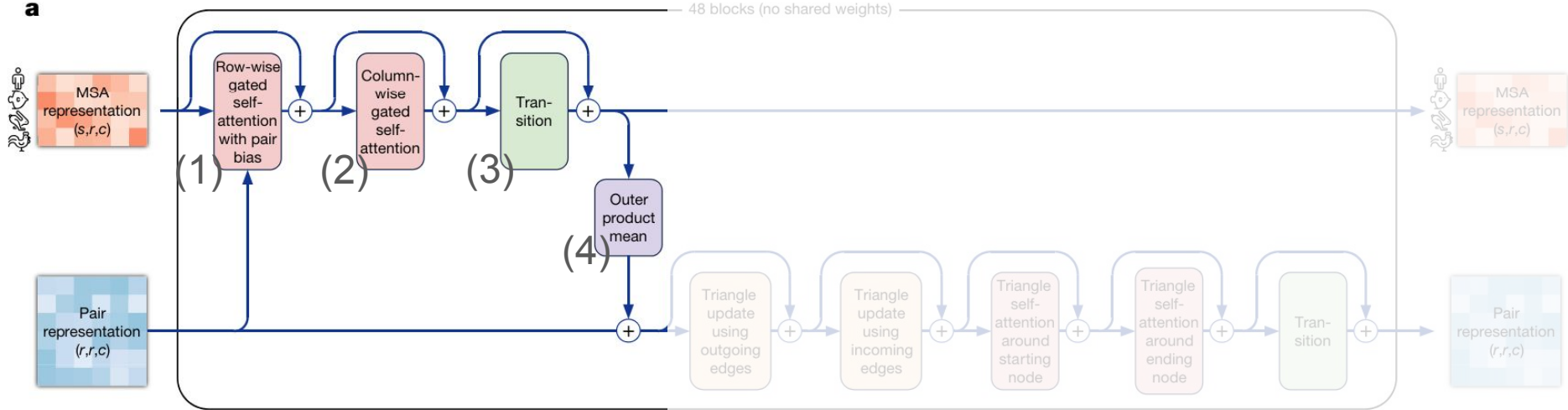
$$\mathbf{o}_{si}^h = \mathbf{g}_{si}^h \odot \sum_t a_{sti}^h \mathbf{v}_{st}^h$$

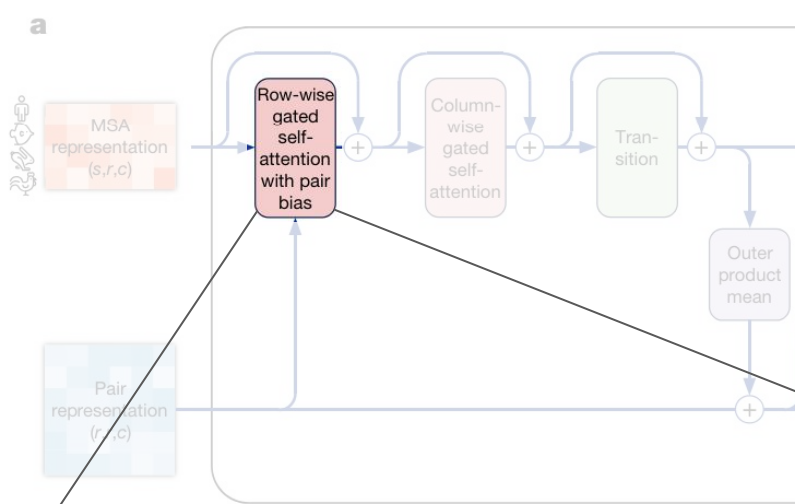
element-wise product

vanilla scaled dot-product attn

**a**



**a**



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

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

*# Input projections*

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

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

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

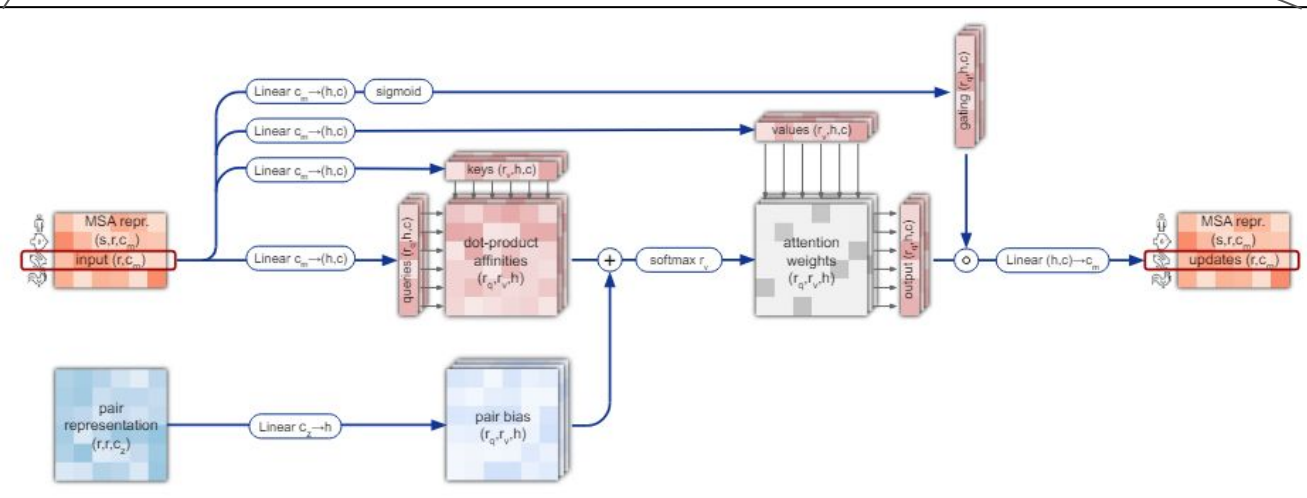
*# Attention*

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

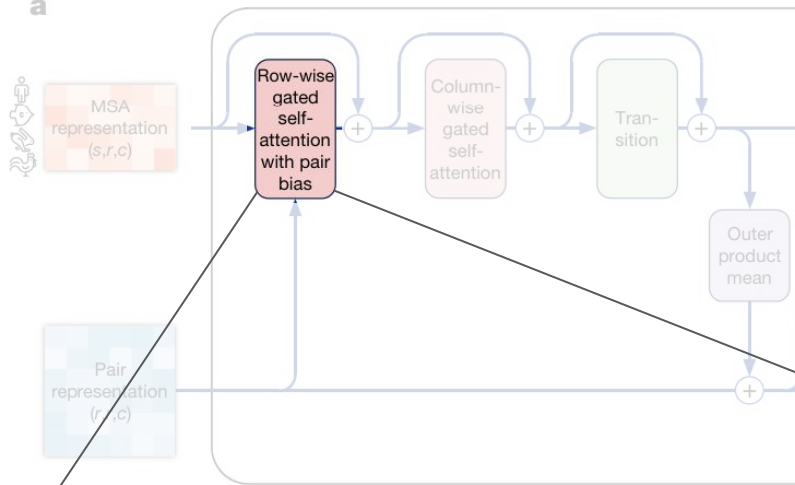
*# Output projection*

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

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



a



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

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

*# Input projections*

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

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

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

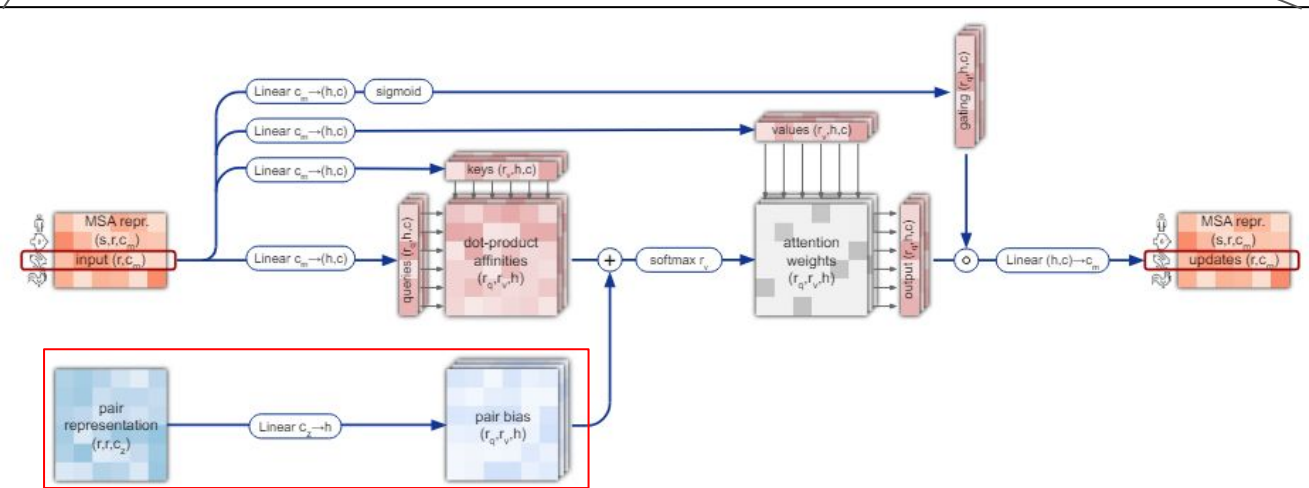
*# Attention*

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

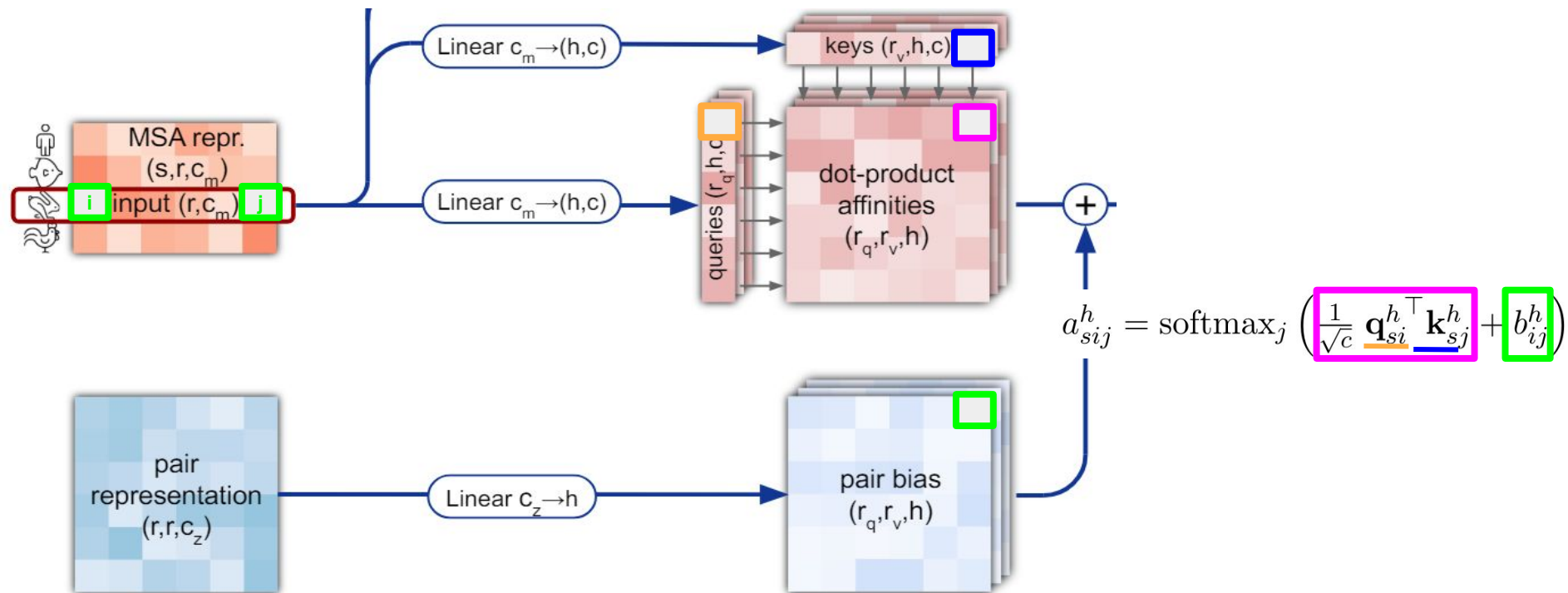
*# Output projection*

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

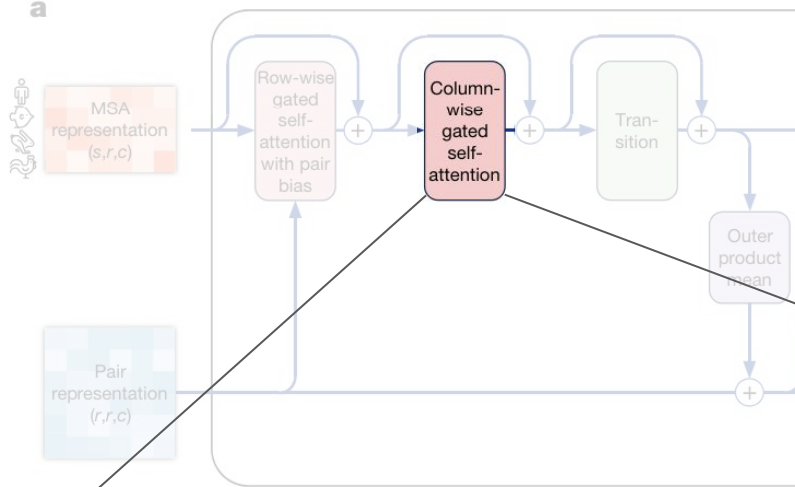
$$\tilde{\mathbf{m}}_{si} \in \mathbb{R}^{c_m}$$



Pair Representation influences MSA via a bias term in the row-wise Gated self-attn



a



### Algorithm 8 MSA column-wise gated self-attention

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

*# Input projections*

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

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

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

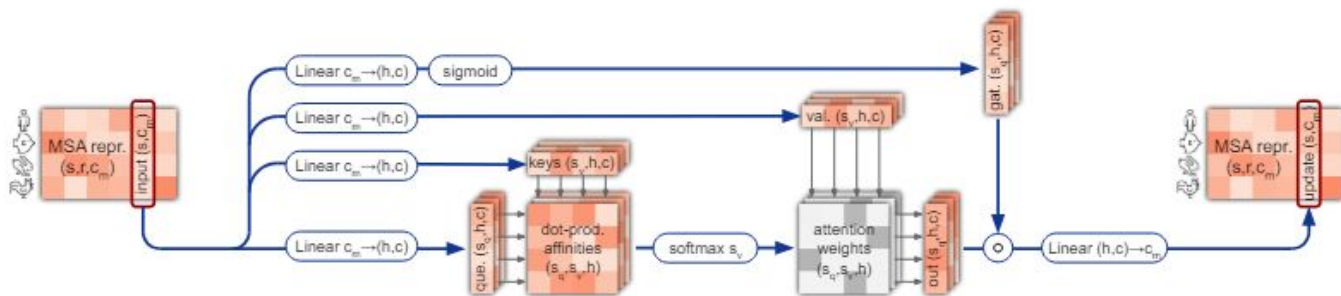
*# Attention*

- 4:  $a_{sti}^h = \text{softmax}_t \left( \frac{1}{\sqrt{c}} \mathbf{q}_{si}^h \mathbf{k}_{ti}^h \right)$
- 5:  $\mathbf{o}_{si}^h = \mathbf{g}_{si}^h \odot \sum_t a_{sti}^h \mathbf{v}_{st}^h$

*# Output projection*

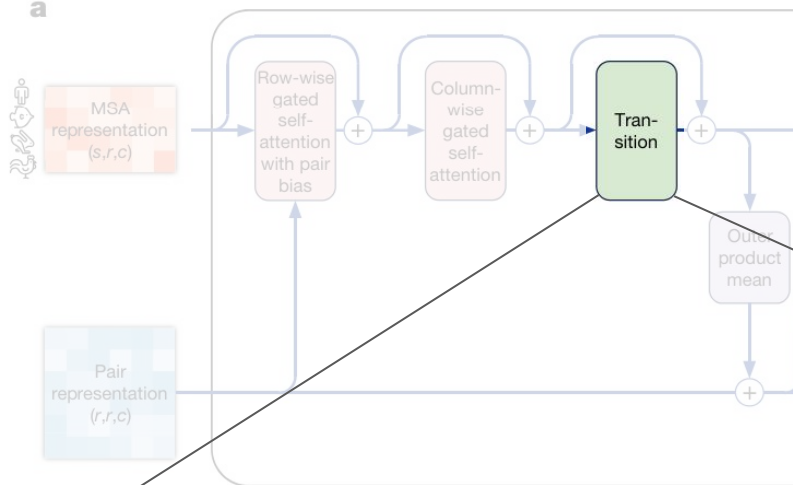
- 6:  $\tilde{\mathbf{m}}_{si} = \text{Linear}(\text{concat}_h(\mathbf{o}_{si}^h))$
- 7: **return**  $\{\tilde{\mathbf{m}}_{si}\}$

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



Vanilla  
Column-wise  
Gated self-attn

**Supplementary Figure 3** | MSA column-wise gated self-attention. Dimensions: s: sequences, r: residues, c: channels, h: heads.

**a**

### Algorithm 9 Transition layer in the MSA stack

**def** MSATransition( $\{\mathbf{m}_{si}\}, n = 4$ ) :

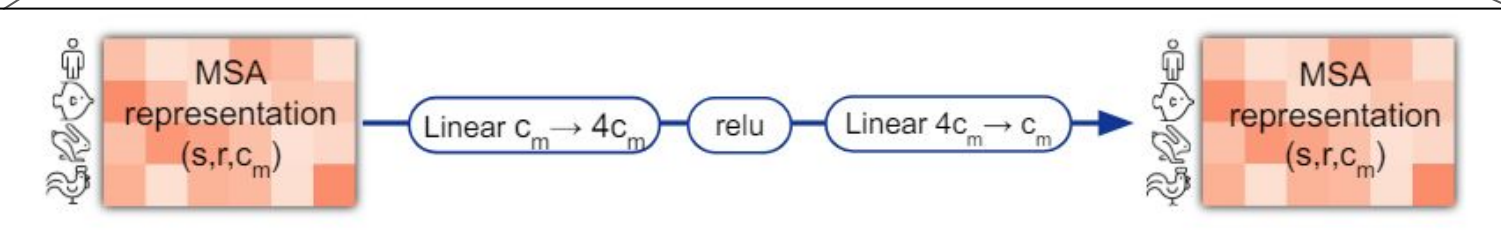
1:  $\mathbf{m}_{si} \leftarrow \text{LayerNorm}(\mathbf{m}_{si})$

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

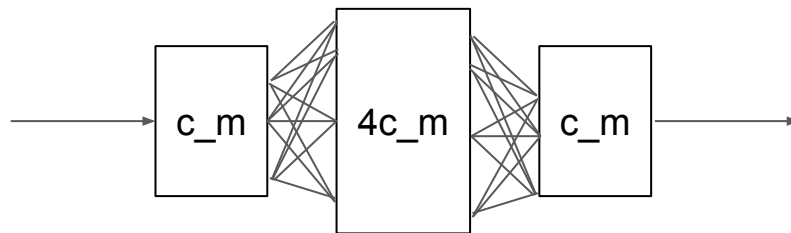
3:  $\mathbf{m}_{si} \leftarrow \text{Linear}(\text{relu}(\mathbf{a}_{si}))$

4: **return**  $\{\mathbf{m}_{si}\}$

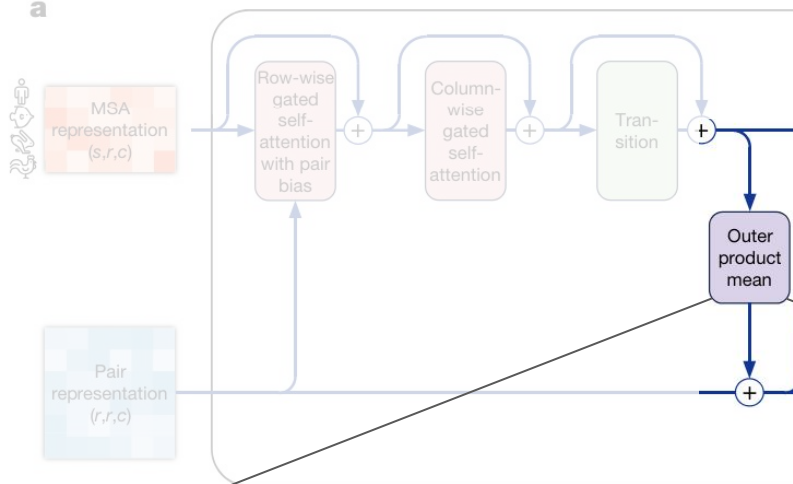
$$\mathbf{a}_{si} \in \mathbb{R}^{n \cdot c_m}$$



1-layer MLP  
with relu  
activation



a



### Algorithm 10 Outer product mean

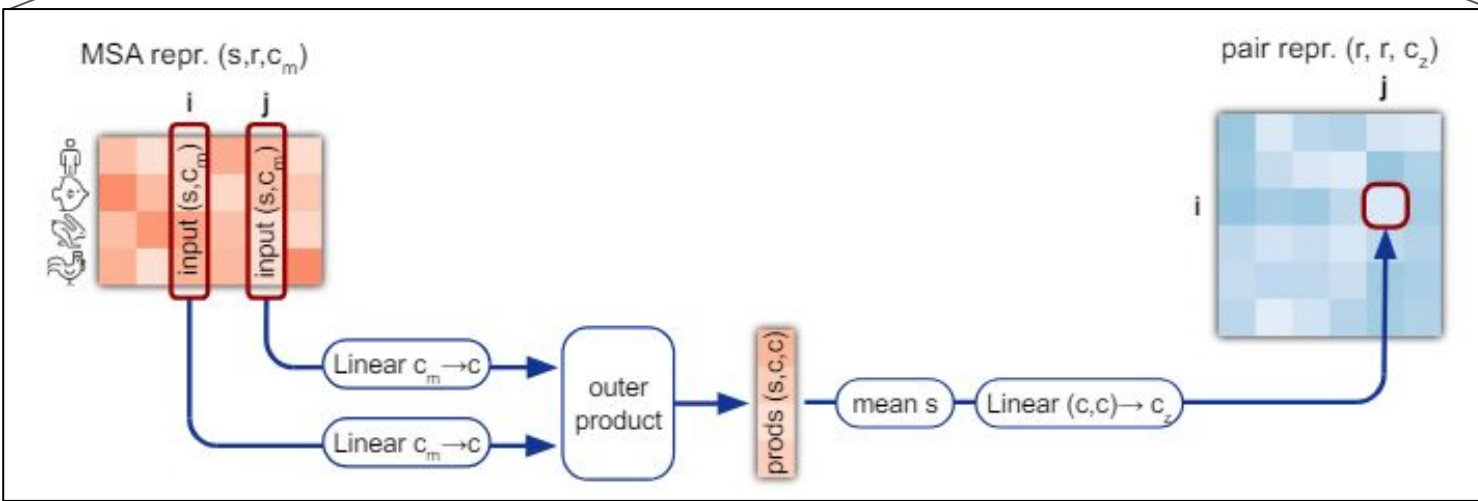
**def** OuterProductMean( $\{\mathbf{m}_{si}\}, c = 32$ ) :

- 1:  $\mathbf{m}_{si} \leftarrow \text{LayerNorm}(\mathbf{m}_{si})$
- 2:  $\mathbf{a}_{si}, \mathbf{b}_{si} = \text{Linear}(\mathbf{m}_{si})$
- 3:  $\mathbf{o}_{ij} = \text{flatten}(\text{mean}_s(\mathbf{a}_{si} \otimes \mathbf{b}_{sj}))$
- 4:  $\mathbf{z}_{ij} = \text{Linear}(\mathbf{o}_{ij})$
- 5: **return**  $\{\mathbf{z}_{ij}\}$

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

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

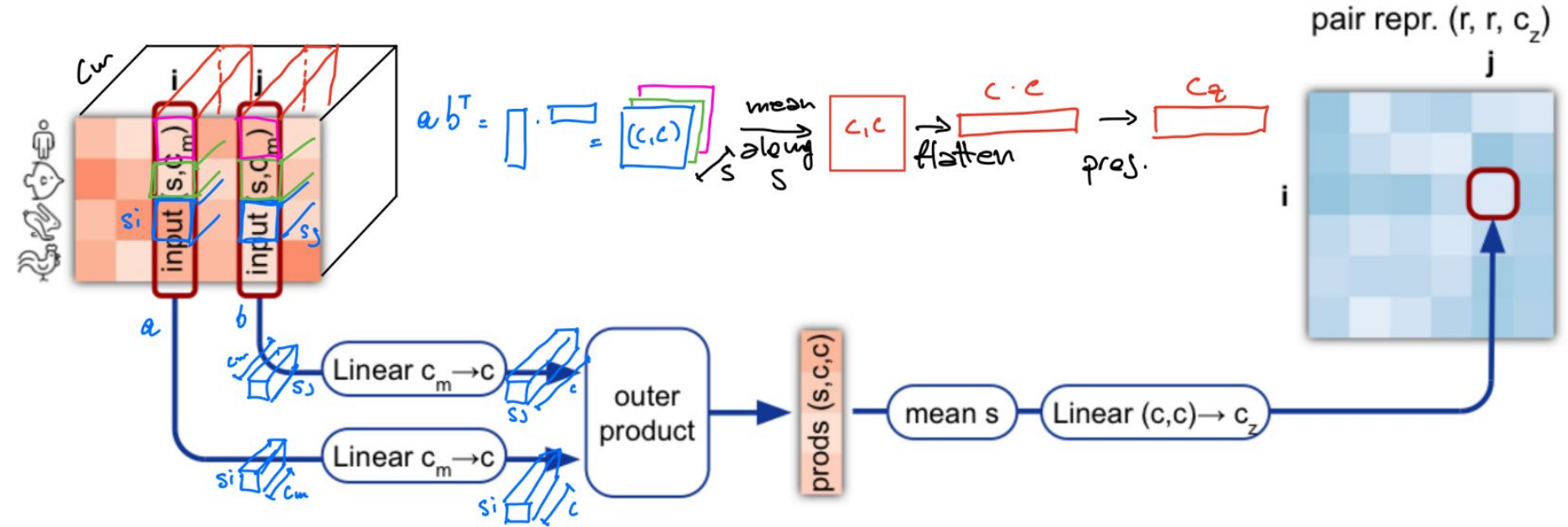
$$\mathbf{z}_{ij} \in \mathbb{R}^{c_z}$$



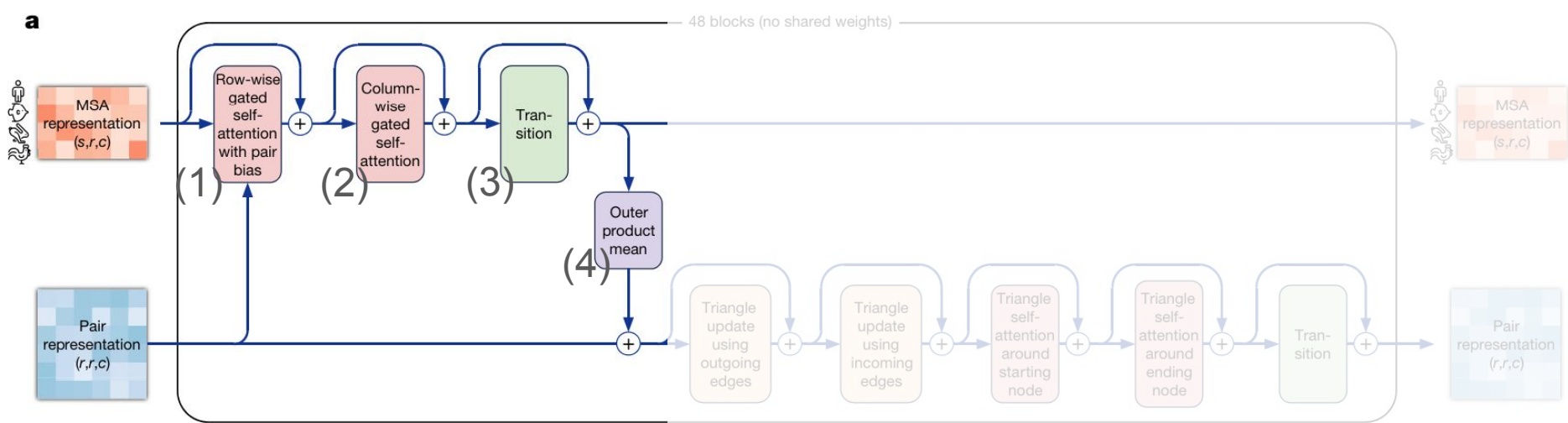
MSA influences  
Pair  
Representation  
via  
outer-product  
update of each  
entry



# Note: Outer Product Mean



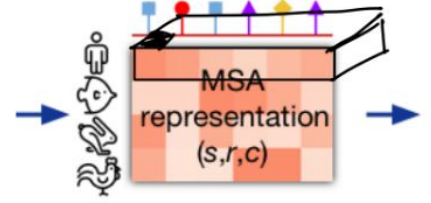
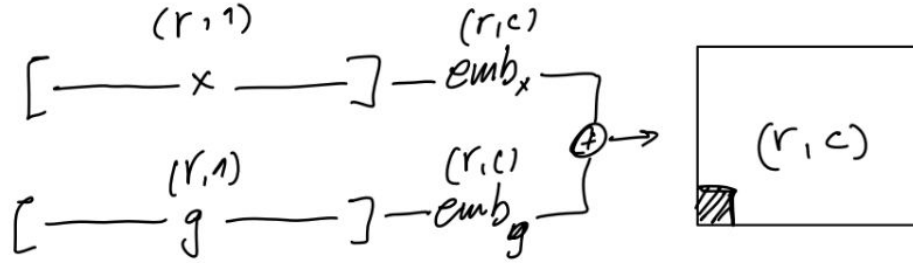




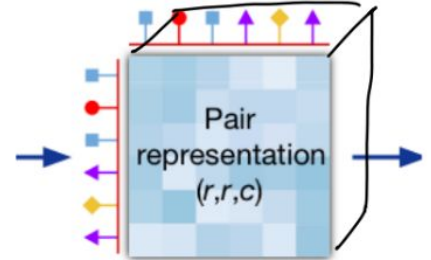
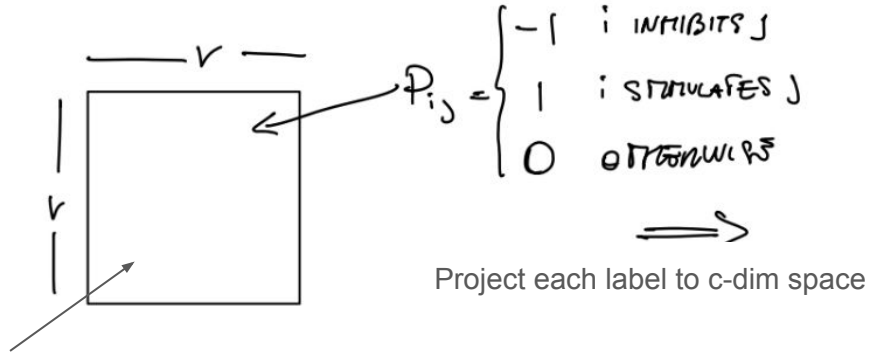
idea

# Idea

scRNA-seq:



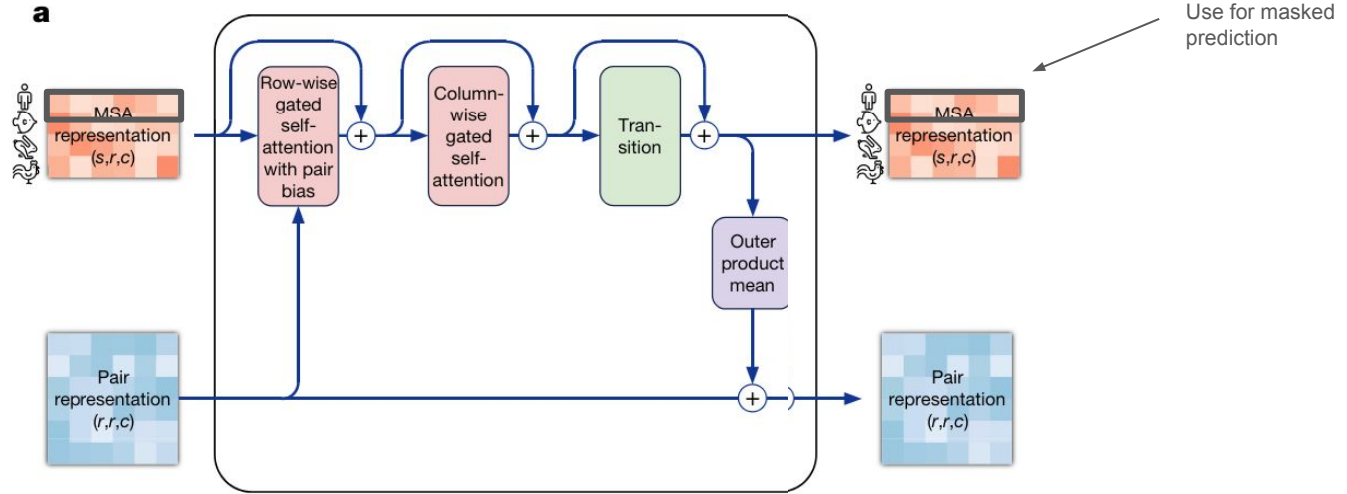
TFs Interactions:



Could add multiple layers (e.g. Tommaso's enrichment score), then project and sum as above

# Idea

scRNA-seq:



TFs Interactions:

Thank You : )