# Highly accurate protein structure prediction with AlphaFold
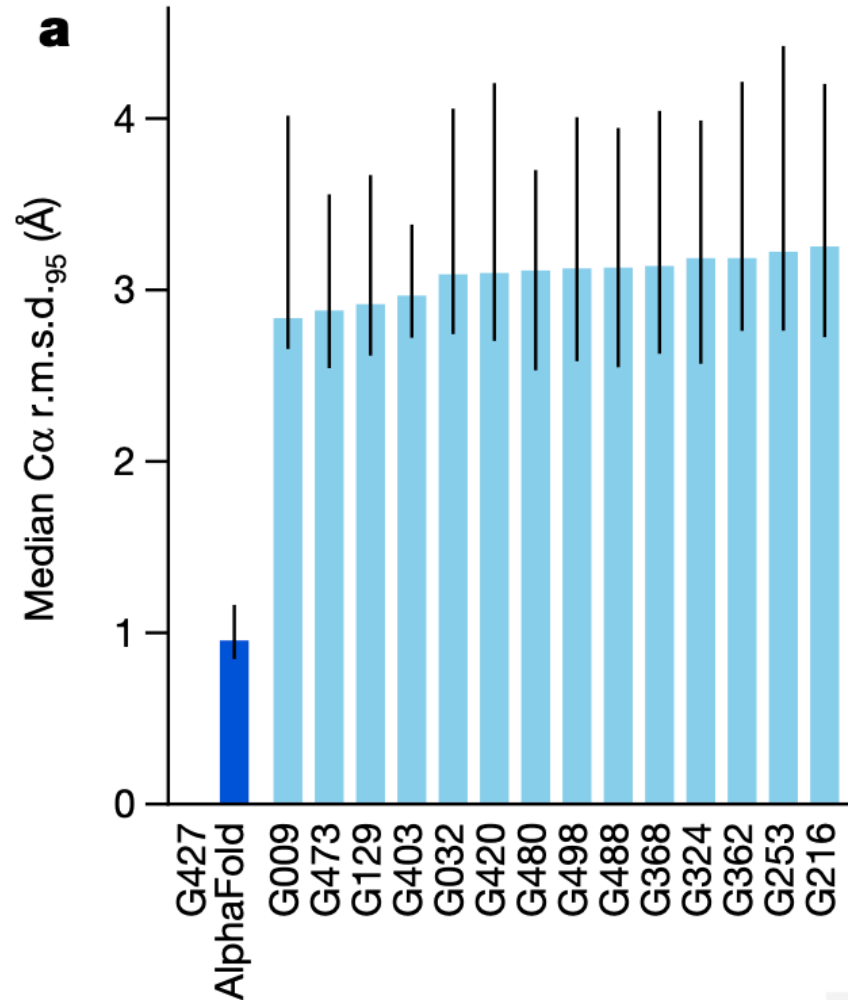
Ning Sun

# High-level impact

Timeline

- Dec 2018: Alphafold 1 wins CASP
    - CASP: Critical Assessment of protein Structure Prediction
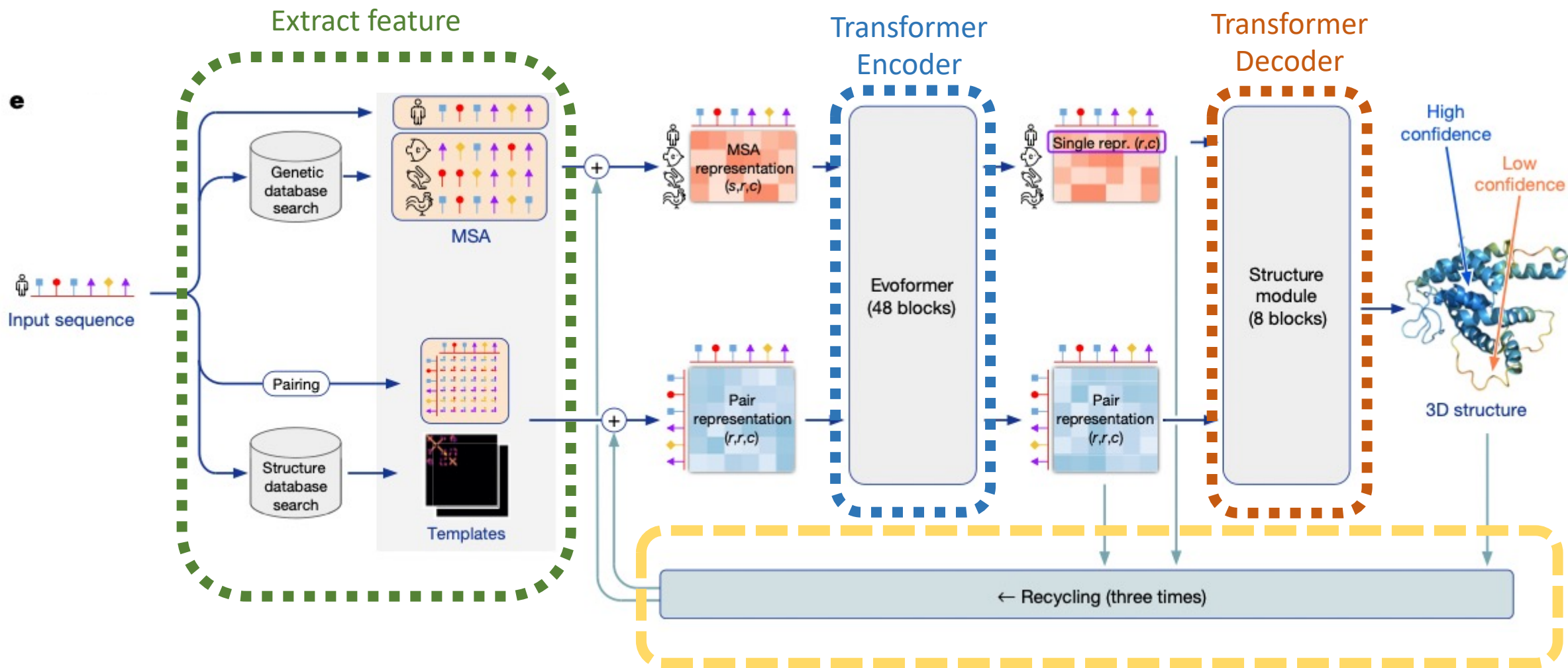- Nov 2020: Alphafold 2 solves CASP


Impact on drug discovery

- Universal end-to-end molecular drug discovery pipeline now available

# Performance at CASP



First computational method that predict protein structures with atomic accuracy.

# Alphafold2

# Data Preprocessing

Dataset: 25% PDB dataset + 75% self-distillation Uniclust30 dataset(unlabeled)

MSA: multiple sequence alignments

Template: 3D atom coordinates of homologous structures

MSA Preprocess:

- Filter

- MSA Block Deletion

- MSA Clustering (cluster_profile, extra_msa, mask_msa )

- Residue Cropping(training)

# Training & Inference details

**MSA resample and ensemble**

- MSA block deletion and clustering is stochastic method

- Resampling at training and inference time

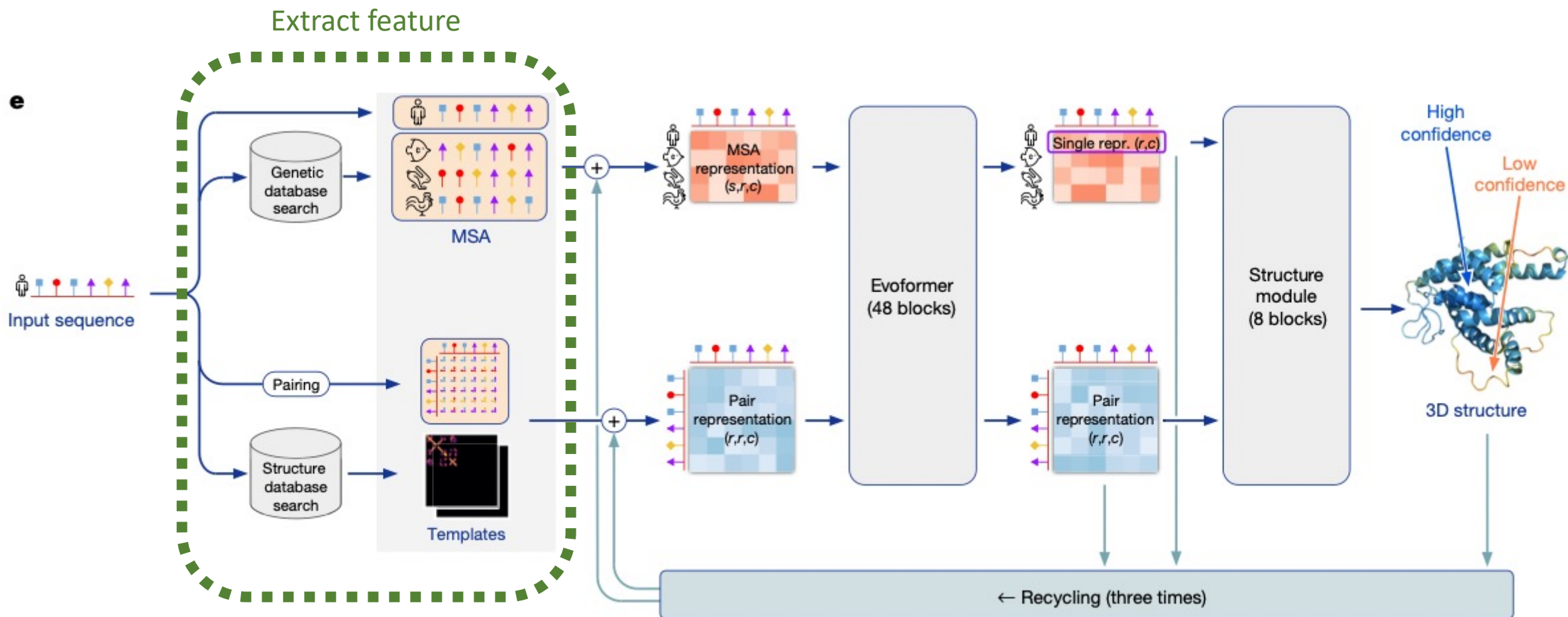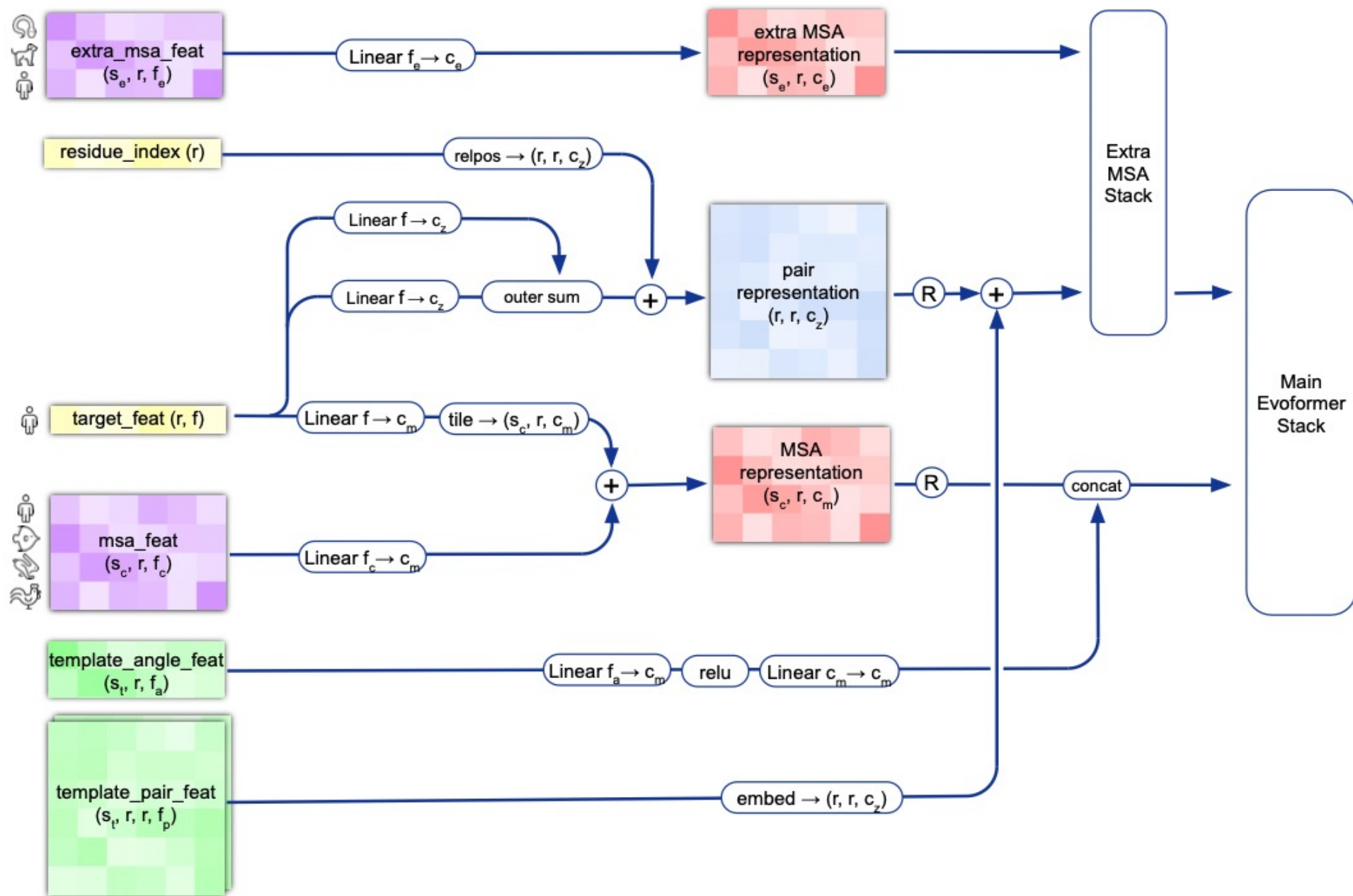- Ensemble at inference time

**Recycling**

- Random $N \in [0, N_{cycle}]$ at training time

- Fixed $N_{cycle}$ at inference time

**Reducing the memory consumption**

- gradient checkpoint at training time

- 'chunk' at inference time

# Alphafold2

**Algorithm 4** Relative position encoding

---

**def** $\text{relpos}(\{f_i^{\text{residue\_index}}\}, \mathbf{v}_{\text{bins}} = [-32, -31, \dots, 32])$ :
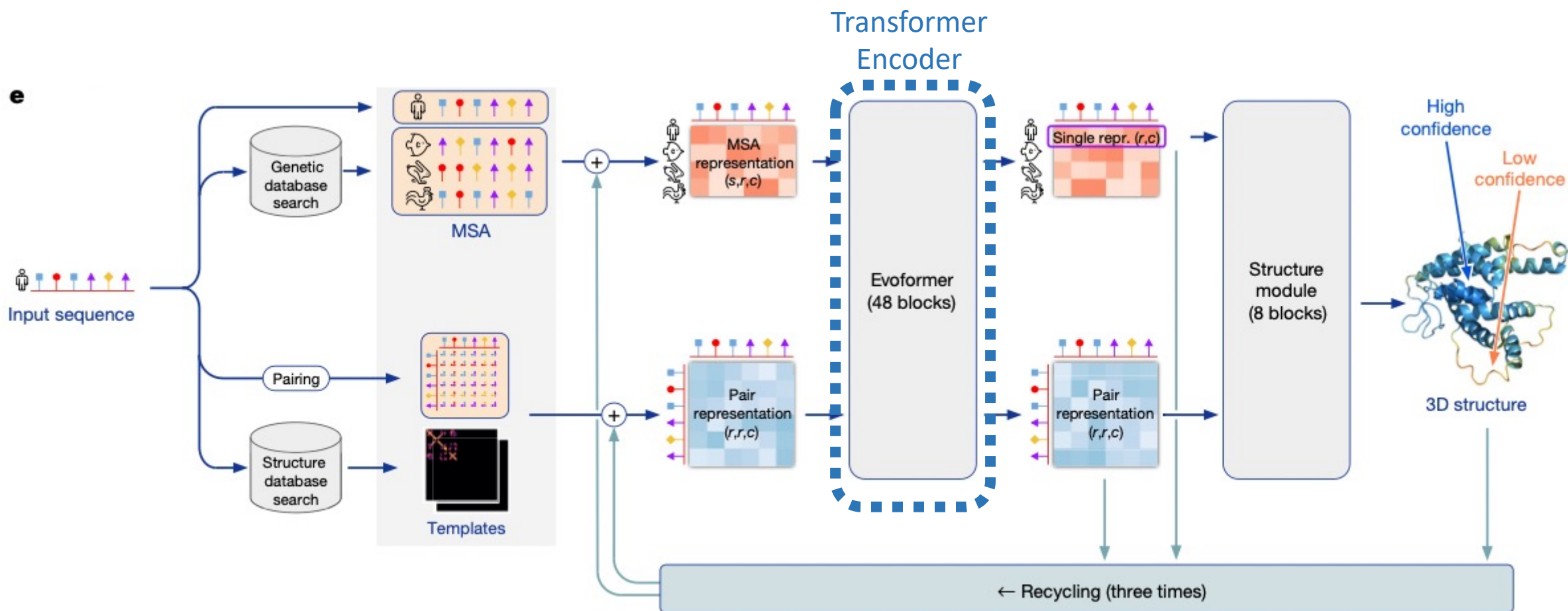
1: $d_{ij} = f_i^{\text{residue\_index}} - f_j^{\text{residue\_index}}$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad d_{ij} \in \mathbb{Z}$

2: $\mathbf{p}_{ij} = \text{Linear}(\text{one\_hot}(d_{ij}, \mathbf{v}_{\text{bins}}))$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \mathbf{p}_{ij} \in \mathbb{R}^{c_z}$

3: **return** $\{\mathbf{p}_{ij}\}$

---

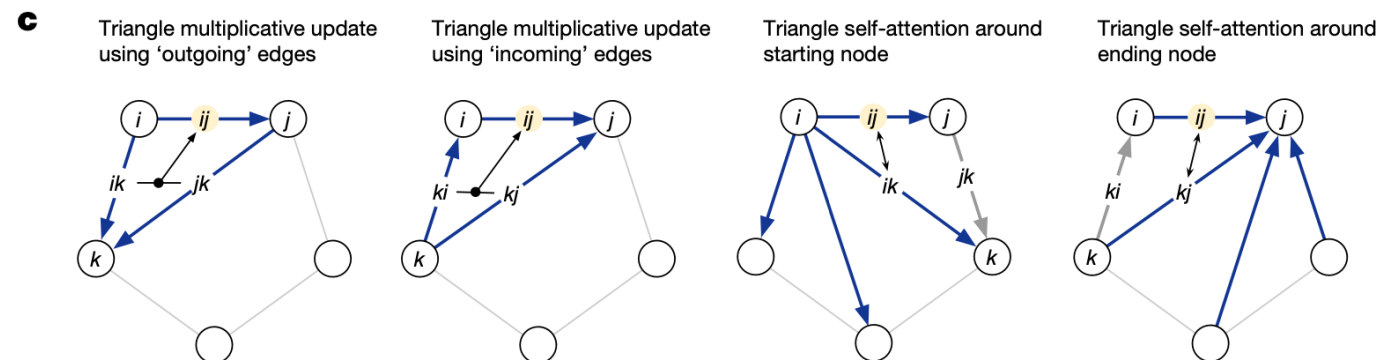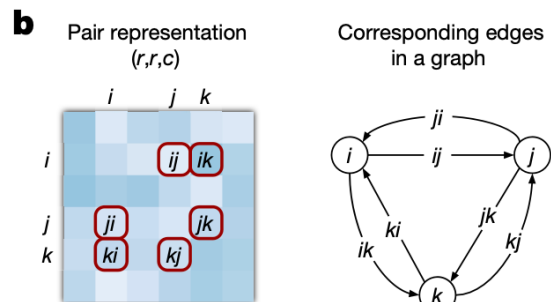**Algorithm 17** Template pointwise attention

---

**def** $\text{TemplatePointwiseAttention}(\{\mathbf{t}_{s_t ij}\}, \{\mathbf{z}_{ij}\}, c = 64, N_{\text{head}} = 4)$ :

1: $\mathbf{q}_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$ $\qquad\qquad\qquad\qquad\qquad\qquad \mathbf{q}_{ij}^h \in \mathbb{R}^c, \ h \in \{1, \dots, N_{\text{head}}\}$

2: $\mathbf{k}_{s_t ij}^h, \mathbf{v}_{s_t ij}^h = \text{LinearNoBias}(\mathbf{t}_{s_t ij})$ $\qquad\qquad\qquad\qquad\qquad\qquad \mathbf{k}_{s_t ij}^h, \mathbf{v}_{s_t ij}^h \in \mathbb{R}^c$

3: $a_{s_t ij}^h = \text{softmax}_{s_t} \left( \frac{1}{\sqrt{c}} \, \mathbf{q}_{ij}^{h\top} \mathbf{k}_{s_t ij}^h \right)$

4: $\mathbf{o}_{ij}^h = \sum_{s_t} a_{s_t ij}^h \mathbf{v}_{s_t ij}^h$

5: $\{\tilde{\mathbf{z}}_{ij}\} = \text{Linear}\left( \text{concat}_h(\{\mathbf{o}_{ij}^h\}) \right)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad \tilde{\mathbf{z}}_{ij} \in \mathbb{R}^{c_z}$

6: **return** $\{\tilde{\mathbf{z}}_{ij}\}$

---

# Alphafold2

# Evoformer

# Row-wise gated self-attention



**Supplementary Figure 2** | MSA row-wise gated self-attention with pair bias. Dimensions: s: sequences, r: residues, c: channels, h: heads.

---

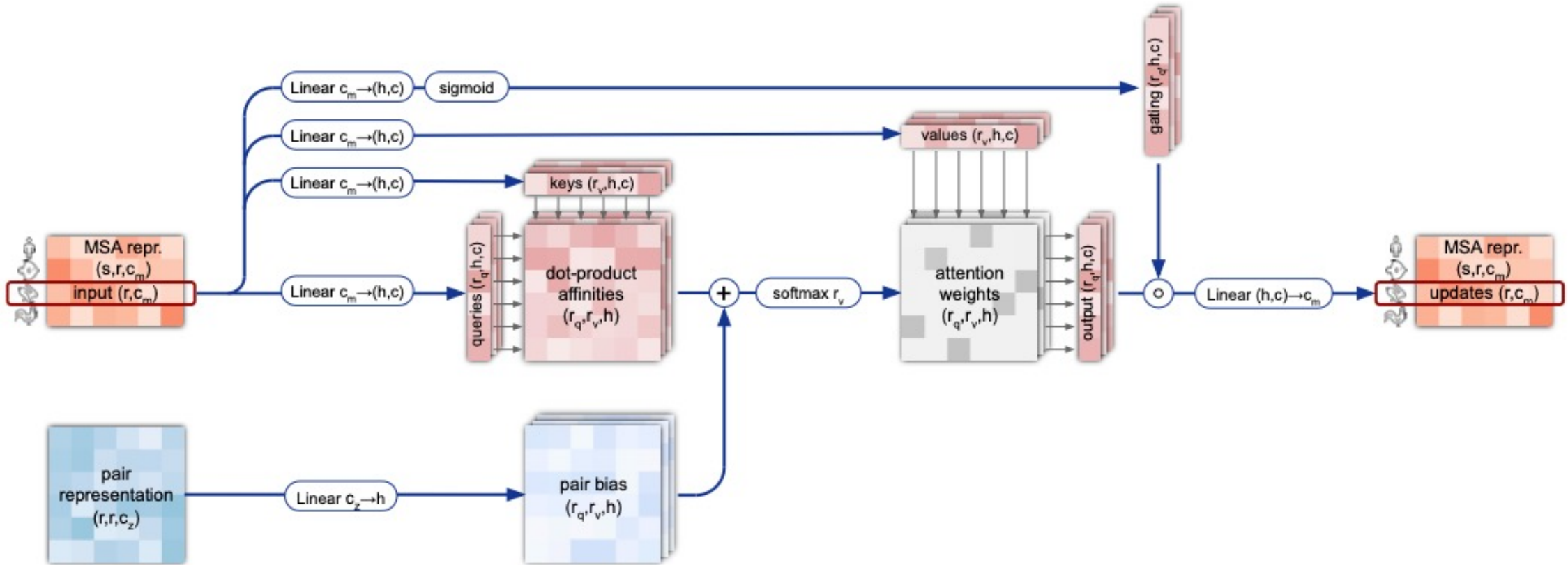**Algorithm 7** MSA row-wise gated self-attention with pair bias

---

**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: $\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h = \text{LinearNoBias}(\mathbf{m}_{si})$ $\qquad\qquad\qquad\qquad\qquad$ $\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h \in \mathbb{R}^c,\ h \in \{1, \dots, N_{\text{head}}\}$

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

4: $\mathbf{g}_{si}^h = \text{sigmoid}\left(\text{Linear}(\mathbf{m}_{si})\right)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\mathbf{g}_{si}^h \in \mathbb{R}^c$

# *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: $\tilde{\mathbf{m}}_{si} = \text{Linear}\left( \text{concat}_h(\mathbf{o}_{si}^h) \right)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\tilde{\mathbf{m}}_{si} \in \mathbb{R}^{c_m}$

8: **return** $\{\tilde{\mathbf{m}}_{si}\}$

---

# Column-wise gated self-attention



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

---

**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:  $\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h = \text{LinearNoBias}(\mathbf{m}_{si})$         $\mathbf{q}_{si}^h, \mathbf{k}_{si}^h, \mathbf{v}_{si}^h \in \mathbb{R}^c, \ h \in \{1, \ldots, N_{\text{head}}\}$

3:  $\mathbf{g}_{si}^h = \text{sigmoid}\left(\text{Linear}(\mathbf{m}_{si})\right)$         $\mathbf{g}_{si}^h \in \mathbb{R}^c$

   *# Attention*

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

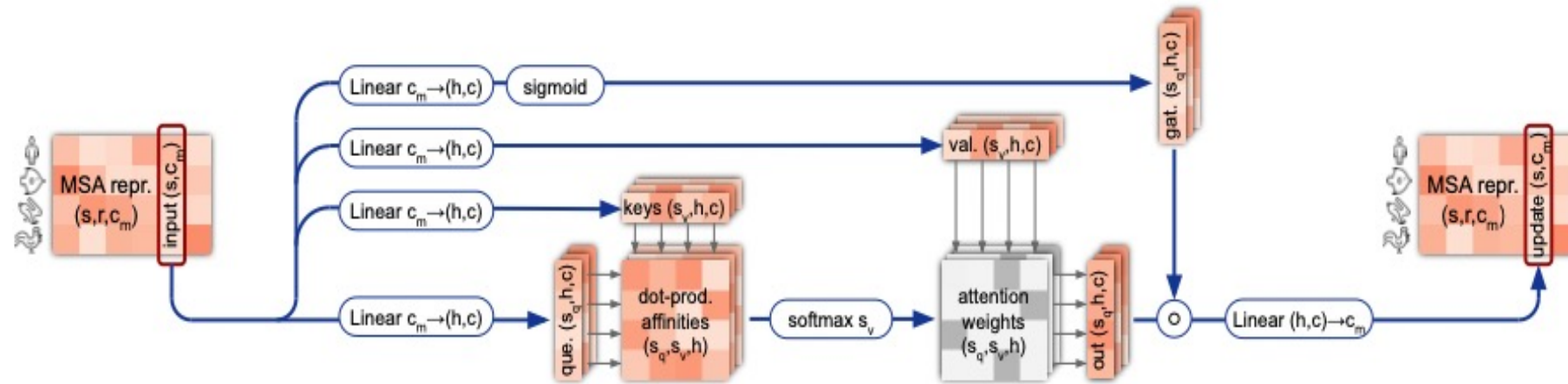5:  $\mathbf{o}_{si}^h = \mathbf{g}_{si}^h \odot \sum_t a_{sti}^h \mathbf{v}_{st}^h$

   *# Output projection*

6:  $\tilde{\mathbf{m}}_{si} = \text{Linear}\left(\text{concat}_h(\mathbf{o}_{si}^h)\right)$         $\tilde{\mathbf{m}}_{si} \in \mathbb{R}^{cm}$

7:  **return** $\{\tilde{\mathbf{m}}_{si}\}$

---

**Algorithm 13** Triangular gated self-attention around <mark>starting node</mark>

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

*# Input projections*

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

2: $\mathbf{q}_{ij}^h, \mathbf{k}_{ij}^h, \mathbf{v}_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$         $\mathbf{q}_{ij}^h, \mathbf{k}_{ij}^h, \mathbf{v}_{ij}^h \in \mathbb{R}^c, \ h \in \{1, \dots, N_{\text{head}}\}$

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

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

*# Attention*

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

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

*# Output projection*

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

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

**Algorithm 14** Triangular gated self-attention around ==ending== node

---

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

  *# Input projections*

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

2:   $\mathbf{q}_{ij}^h, \mathbf{k}_{ij}^h, \mathbf{v}_{ij}^h = \text{LinearNoBias}(\mathbf{z}_{ij})$          $\mathbf{q}_{ij}^h, \mathbf{k}_{ij}^h, \mathbf{v}_{ij}^h \in \mathbb{R}^c, \; h \in \{1, \ldots, N_{\text{head}}\}$

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

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

  *# Attention*

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

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

  *# Output projection*

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

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

---

---

**Algorithm 11** Triangular multiplicative update using "outgoing" edges

---

**def** TriangleMultiplicationOutgoing($\{\mathbf{z}_{ij}\}, c = 128$) :

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

2: $\mathbf{a}_{ij}, \mathbf{b}_{ij} = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right) \odot \text{Linear}(\mathbf{z}_{ij})$ $\qquad\qquad\qquad$ $\mathbf{a}_{ij}, \mathbf{b}_{ij} \in \mathbb{R}^c$

3: $\mathbf{g}_{ij} = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right)$ $\qquad\qquad\qquad$ $\mathbf{g}_{ij} \in \mathbb{R}^{c_z}$

4: $\tilde{\mathbf{z}}_{ij} = \mathbf{g}_{ij} \odot \text{Linear}(\text{LayerNorm}(\sum_k \mathbf{a}_{ik} \odot \mathbf{b}_{jk}))$ $\qquad\qquad\qquad$ $\tilde{\mathbf{z}}_{ij} \in \mathbb{R}^{c_z}$

5: **return** $\{\tilde{\mathbf{z}}_{ij}\}$

---

 

---

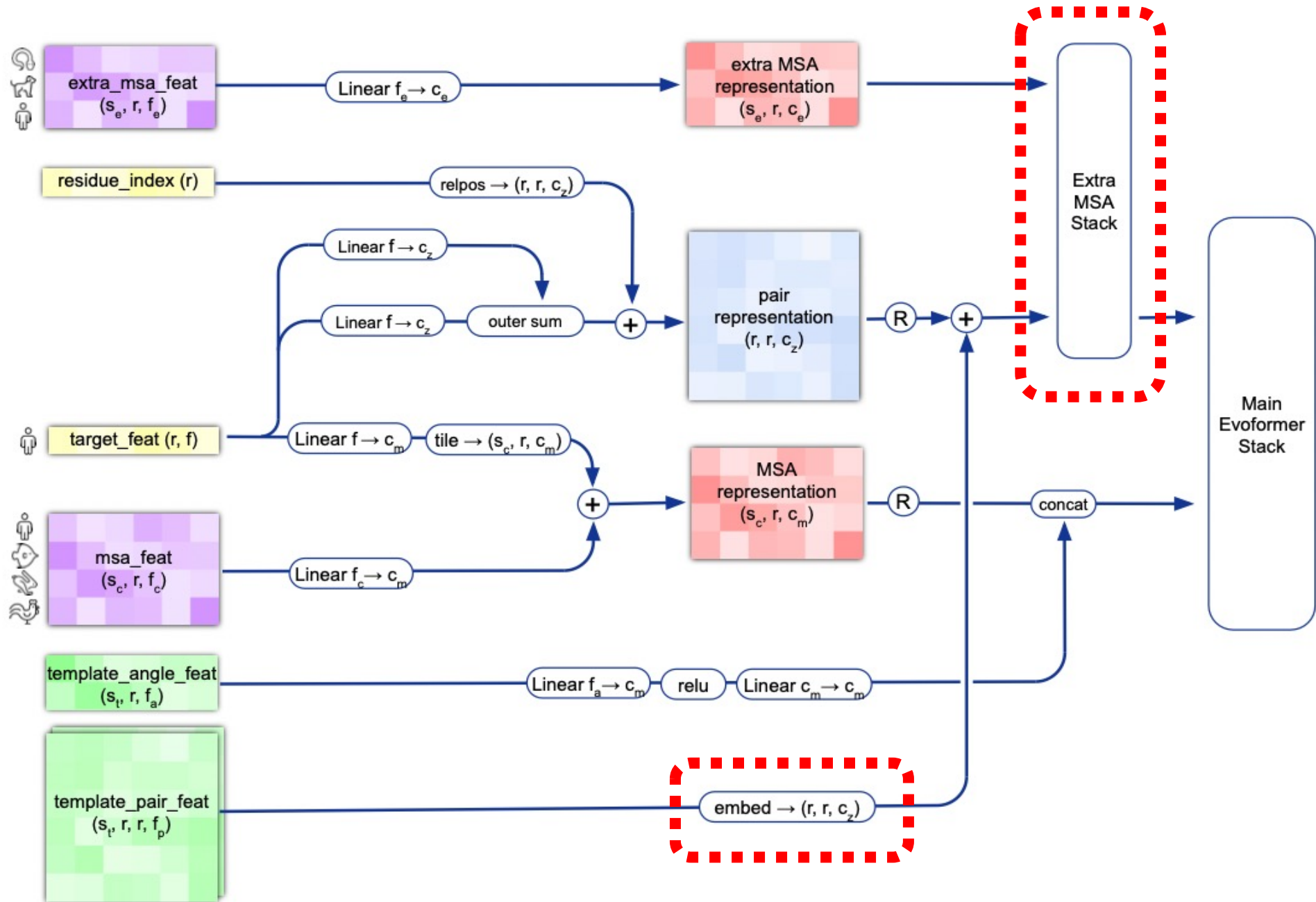**Algorithm 12** Triangular multiplicative update using "incoming" edges

---

**def** TriangleMultiplicationIncoming($\{\mathbf{z}_{ij}\}, c = 128$) :

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

2: $\mathbf{a}_{ij}, \mathbf{b}_{ij} = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right) \odot \text{Linear}(\mathbf{z}_{ij})$ $\qquad\qquad\qquad$ $\mathbf{a}_{ij}, \mathbf{b}_{ij} \in \mathbb{R}^c$

3: $\mathbf{g}_{ij} = \text{sigmoid}\left(\text{Linear}(\mathbf{z}_{ij})\right)$ $\qquad\qquad\qquad$ $\mathbf{g}_{ij} \in \mathbb{R}^{c_z}$

4: $\tilde{\mathbf{z}}_{ij} = \mathbf{g}_{ij} \odot \text{Linear}(\text{LayerNorm}(\sum_k \mathbf{a}_{ki} \odot \mathbf{b}_{kj}))$ $\qquad\qquad\qquad$ $\tilde{\mathbf{z}}_{ij} \in \mathbb{R}^{c_z}$

5: **return** $\{\tilde{\mathbf{z}}_{ij}\}$

---

**Algorithm 16** Template pair stack

---

**def** $\text{TemplatePairStack}(\{\mathbf{t}_{ij}\}, N_{\text{block}} = 2)$ :

1: **for all** $l \in [1, \ldots, N_{\text{block}}]$ **do**

2: $\quad \{\mathbf{t}_{ij}\} \mathrel{+}= \text{DropoutRowwise}_{0.25}(\text{TriangleAttentionStartingNode}(\{\mathbf{t}_{ij}\}, c = 64, N_{\text{head}} = 4))$

3: $\quad \{\mathbf{t}_{ij}\} \mathrel{+}= \text{DropoutColumnwise}_{0.25}(\text{TriangleAttentionEndingNode}(\{\mathbf{t}_{ij}\}, c = 64, N_{\text{head}} = 4))$

4: $\quad \{\mathbf{t}_{ij}\} \mathrel{+}= \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationOutgoing}(\{\mathbf{t}_{ij}\}, c = 64))$

5: $\quad \{\mathbf{t}_{ij}\} \mathrel{+}= \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationIncoming}(\{\mathbf{t}_{ij}\}, c = 64))$

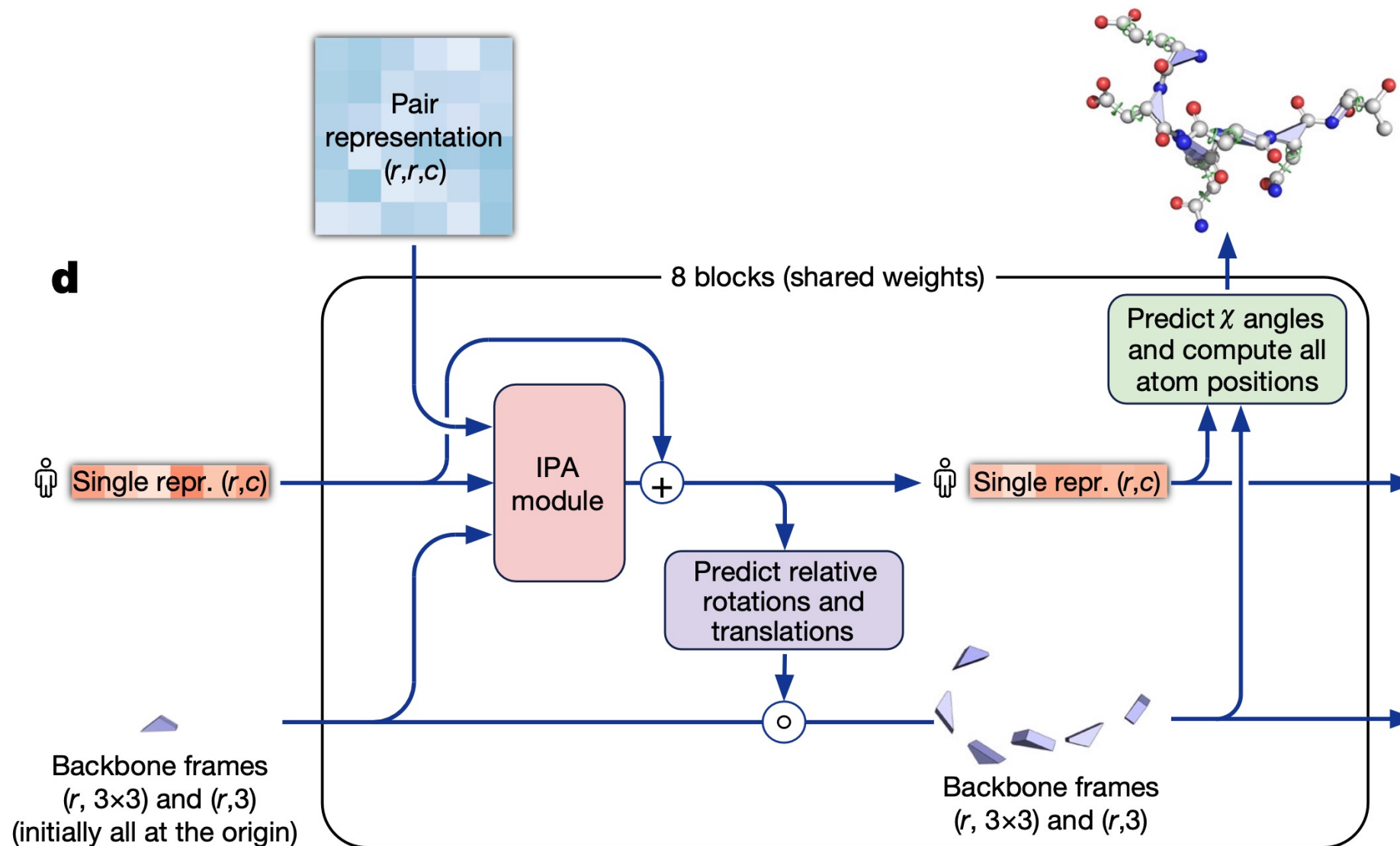6: $\quad \{\mathbf{t}_{ij}\} \mathrel{+}= \text{PairTransition}(\{\mathbf{t}_{ij}\}, n = 2)$

7: **end for**

8: **return** $\text{LayerNorm}\left(\{\mathbf{t}_{ij}\}\right)$

---

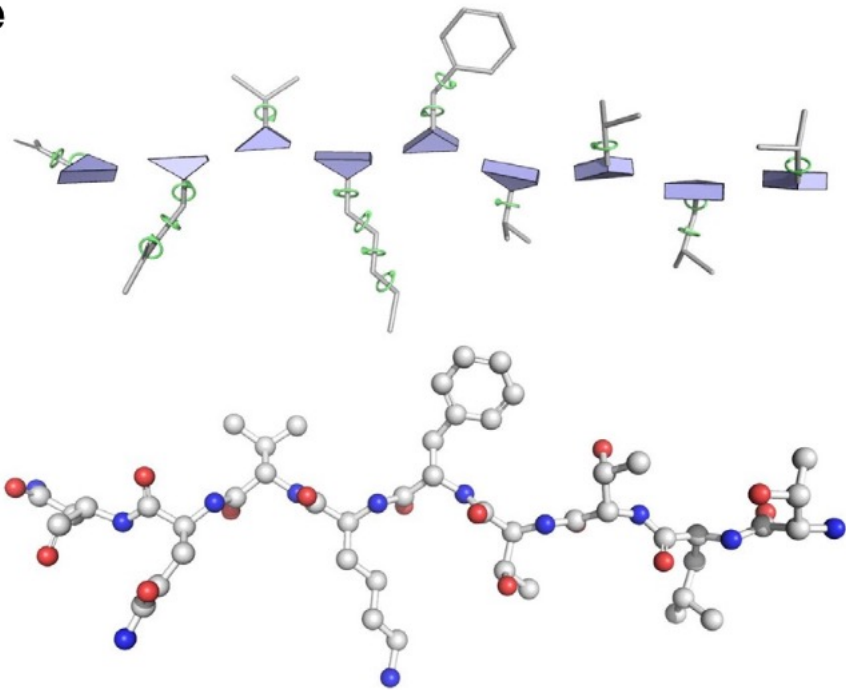**Algorithm 18** Extra MSA stack

$\textbf{def}\ \text{ExtraMsaStack}(\{\mathbf{e}_{s_e i}\}, \{\mathbf{z}_{ij}\}, N_{\text{block}} = 4):$

1: **for all** $l \in [1, \dots, N_{\text{block}}]$ **do**

   #    *MSA stack*

2:     $\{\mathbf{e}_{s_e i}\}\ \mathrel{+}=\ \text{DropoutRowwise}_{0.15}(\text{MSARowAttentionWithPairBias}(\{\mathbf{e}_{s_e i}\}, \{\mathbf{z}_{ij}\}, c = 8))$

3:     $\{\mathbf{e}_{s_e i}\}\ \mathrel{+}=\ \text{MSAColumn Global Attention}(\{\mathbf{e}_{s_e i}\}$

4:     $\{\mathbf{e}_{s_e i}\}\ \mathrel{+}=\ \text{MSATransition}(\{\mathbf{e}_{s_e i}\})$

   #    *Communication*

5:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{OuterProductMean}(\{\mathbf{e}_{s_e i}\})$

   #    *Pair stack*

6:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationOutgoing}(\{\mathbf{z}_{ij}\}))$

7:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{DropoutRowwise}_{0.25}(\text{TriangleMultiplicationIncoming}(\{\mathbf{z}_{ij}\}))$

8:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{DropoutRowwise}_{0.25}(\text{TriangleAttentionStartingNode}(\{\mathbf{z}_{ij}\}))$

9:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{DropoutColumnwise}_{0.25}(\text{TriangleAttentionEndingNode}(\{\mathbf{z}_{ij}\}))$

10:     $\{\mathbf{z}_{ij}\}\ \mathrel{+}=\ \text{PairTransition}(\{\mathbf{z}_{ij}\})$

11: **end for**

12: **return** $\{\mathbf{z}_{ij}\}$

# Structure Module

# Frame and torsion angle



- Backbone: $\{T_i = (R_i, t_i)\}$ map from local frame to global frame

- Side chain: the torsion angles are the only degrees of freedom, while all bond angles and bond lengths are fully rigid.

# 3D Equivariance

---

**Algorithm 22** Invariant point attention (IPA)

---

$\textbf{def}$ InvariantPointAttention($\{\mathbf{s}_i\}, \{\mathbf{z}_{ij}\}, \{T_i\}, N_{\text{head}} = 12, c = 16, N_{\text{query points}} = 4, N_{\text{point values}} = 8$) :

1: $\mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h = \text{LinearNoBias}(\mathbf{s}_i)$ $\hspace{3em} \mathbf{q}_i^h, \mathbf{k}_i^h, \mathbf{v}_i^h \in \mathbb{R}^c, \ h \in \{1, \ldots, N_{\text{head}}\}$

2: $\vec{\mathbf{q}}_i^{hp}, \vec{\mathbf{k}}_i^{hp} = \text{LinearNoBias}(\mathbf{s}_i)$ $\hspace{2em} \vec{\mathbf{q}}_i^{hp}, \vec{\mathbf{k}}_i^{hp}, \in \mathbb{R}^3, \ p \in \{1, \ldots, N_{\text{query points}}\}, \ \text{units: nanometres}$

3: $\vec{\mathbf{v}}_i^{hp} = \text{LinearNoBias}(\mathbf{s}_i)$ $\hspace{3em} \vec{\mathbf{v}}_i^{hp} \in \mathbb{R}^3, \ p \in \{1, \ldots, N_{\text{point values}}\}, \ \text{units: nanometres}$

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

5: $w_C = \sqrt{\dfrac{2}{9 N_{\text{query points}}}},$

6: $w_L = \sqrt{\dfrac{1}{3}}$

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

8: $\tilde{\mathbf{o}}_i^h = \sum_j a_{ij}^h \mathbf{z}_{ij}$

9: $\mathbf{o}_i^h = \sum_j a_{ij}^h \mathbf{v}_j^h$

10: $\vec{\mathbf{o}}_i^{hp} = T_i^{-1} \circ \sum_j a_{ij}^h \left( T_j \circ \vec{\mathbf{v}}_j^{hp} \right)$

11: $\tilde{\mathbf{s}}_i = \text{Linear} \left( \text{concat}_{h,p}(\tilde{\mathbf{o}}_i^h, \mathbf{o}_i^h, \vec{\mathbf{o}}_i^{hp}, \left\| \vec{\mathbf{o}}_i^{hp} \right\|) \right)$

12: $\textbf{return} \ \{\tilde{\mathbf{s}}_i\}$

---

---

**Algorithm 23** Backbone update

---

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

1: $b_i, c_i, d_i, \vec{\mathbf{t}}_i = \text{Linear}(\mathbf{s}_i)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad b_i, c_i, d_i \in \mathbb{R}, \ \vec{\mathbf{t}}_i \in \mathbb{R}^3$

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

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

3: $R_i = \begin{pmatrix} a_i^2 + b_i^2 - c_i^2 - d_i^2 & 2b_ic_i - 2a_id_i & 2b_id_i + 2a_ic_i \\ 2b_ic_i + 2a_id_i & a_i^2 - b_i^2 + c_i^2 - d_i^2 & 2c_id_i - 2a_ib_i \\ 2b_id_i - 2a_ic_i & 2c_id_i + 2a_ib_i & a_i^2 - b_i^2 - c_i^2 + d_i^2 \end{pmatrix}$

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

5: **return** $T_i$

---

R: orthogonal matrix, norm=1

**Algorithm 24** Compute all atom coordinates

---

**def** computeAllAtomCoordinates($T_i, \vec{\alpha}_i^f, F_i^{\text{aatype}}$) :

1: $\hat{\vec{\alpha}}_i^f = \vec{\alpha}_i^f / \|\vec{\alpha}_i^f\|$

2: $(\vec{\omega}_i, \vec{\phi}_i, \vec{\psi}_i, \vec{\chi}_{1_i}, \vec{\chi}_{2_i}, \vec{\chi}_{3_i}, \vec{\chi}_{4_i}) = \hat{\vec{\alpha}}_i^f$

<span style="color:darkred">*# Make extra backbone frames.*</span>

3: $r_i = F_i^{\text{aatype}}$

4: $T_{i1} = T_i \circ T^{\text{lit}}_{r_i,(\omega \rightarrow \text{bb})} \circ \text{makeRotX}(\vec{\omega}_i)$

5: $T_{i2} = T_i \circ T^{\text{lit}}_{r_i,(\phi \rightarrow \text{bb})} \circ \text{makeRotX}(\vec{\phi}_i)$

6: $T_{i3} = T_i \circ T^{\text{lit}}_{r_i,(\psi \rightarrow \text{bb})} \circ \text{makeRotX}(\vec{\psi}_i)$

<span style="color:darkred">*# Make side chain frames (chain them up along the side chain).*</span>

7: $T_{i4} = T_i \circ T^{\text{lit}}_{r_i,(\chi 1 \rightarrow \text{bb})} \circ \text{makeRotX}(\vec{\chi}_{1_i})$

8: $T_{i5} = T_{i4} \circ T^{\text{lit}}_{r_i,(\chi 2 \rightarrow \chi 1)} \circ \text{makeRotX}(\vec{\chi}_{2_i})$

9: $T_{i6} = T_{i5} \circ T^{\text{lit}}_{r_i,(\chi 3 \rightarrow \chi 2)} \circ \text{makeRotX}(\vec{\chi}_{3_i})$

10: $T_{i7} = T_{i6} \circ T^{\text{lit}}_{r_i,(\chi 4 \rightarrow \chi 3)} \circ \text{makeRotX}(\vec{\chi}_{4_i})$

<span style="color:darkred">*# Map atom literature positions to the global frame.*</span>

11: $\vec{\mathbf{x}}_i^a = \text{concat}_{f,a'}\left(\{T_i^f \circ \vec{\mathbf{x}}^{\text{lit}}_{r_i,f,a'}\}\right)$

12: **return** $T_i^f, \vec{\mathbf{x}}_i^a$

---

# Loss

- FAPE: scores a set of predicted atom coordinates under a set of predicted local frames against the corresponding ground truth atom coordinates and ground truth local frames

- Auxiliary loss: FAPE + torsion Angle loss

- pLDDT(confidence loss): the per-residue lDDT-Cα scores.

- TM score: assessing global structure of predicted protein

# Recycle

---

**Algorithm 32** Embedding of Evoformer and Structure module outputs for recycling

---

**def** RecyclingEmbedder($\{\mathbf{m}_{1i}\}, \{\mathbf{z}_{ij}\}, \{\vec{\mathbf{x}}_i^{C^\beta}\}$) :

    *# Embed pair distances of backbone atoms:*

1: $d_{ij} = \left\| \vec{\mathbf{x}}_i^{C^\beta} - \vec{\mathbf{x}}_j^{C^\beta} \right\|$                                            $C^\alpha$ used for glycin

2: $\mathbf{d}_{ij} = \text{Linear}(\text{one\_hot}(d_{ij}, \mathbf{v}_{\text{bins}} = [3\tfrac{3}{8}\,\text{Å},\ 5\tfrac{1}{8}\,\text{Å},\ \ldots, 21\tfrac{3}{8}\,\text{Å}]))$             $\mathbf{d}_{ij} \in \mathbb{R}^{c_z}$

    *# Embed output Evoformer representations:*

3: $\tilde{\mathbf{z}}_{ij} = \mathbf{d}_{ij} + \text{LayerNorm}(\mathbf{z}_{ij})$

4: $\tilde{\mathbf{m}}_{1i} = \text{LayerNorm}(\mathbf{m}_{1i})$

5: **return** $\{\tilde{\mathbf{m}}_{1i}\}, \{\tilde{\mathbf{z}}_{ij}\}$

---

# Inference time

-V100 GPU

-seq_len= 256    4.8min

-seq_len =384    9.2min

-seq_len =2500  18h

# Ablation Study