

Supplemental File for “Integrating Unsupervised Language  
Model with Multi-View Multiple Sequence Alignments for High-  
Accuracy Inter-Chain Contact Prediction”

Zi Liu <sup>1,2, #</sup>, Yi-Heng Zhu <sup>3, #</sup>, Long-Chen Shen <sup>1</sup>, Xuan Xiao <sup>2</sup>, Wang-Ren Qiu <sup>2</sup>,  
\*, Dong-Jun Yu <sup>1, \*</sup>

<sup>1</sup>School of Computer Science and Engineering, Nanjing University of Science  
and Technology, Xiaolingwei 200, Nanjing, China, 210094

<sup>2</sup>Computer Department, Jingdezhen Ceramic University,  
Jingdezhen, China, 333403

<sup>3</sup> College of Artificial Intelligence, Nanjing Agricultural University,  
Nanjing, China, 210095

\*All correspondence should be addressed to E-mail: [njyudj@njjust.edu.cn](mailto:njyudj@njjust.edu.cn)

#These authors contributed equally

## Table of Content

### Supporting Texts

**Text S1.** The details of ESM-MSA transformer

### Supporting Tables

**Table S1.** Average precision in top  $N$  predicted contacts on TS630 dataset.

**Table S2.** Average precision in top  $L/K$  predicted contacts, AUPR, and AUROC on TS630 dataset. ALL represents the number of native contacts on the target.

**Table S3.** Average contact precision of different features on the TS630 dataset.

### Supporting Figures

**Figure S1.** The workflow of ESM-MSA

**Figure S2.** Performance comparison between nine feature embeddings regarding AUROC on the TS630 dataset.

# Supporting Texts

## Text S1. The details of ESM-MSA transformer

### A. Masking

For an input MSA, the masking strategy is performed. Specifically, for each sequence in MSA, we randomly sample 15% tokens (amino acids), each of which is changed as a special “masking” token with 80% probability, a randomly-chosen alternate amino acid with 10% probability, and the original input token (i.e., no change) with 10% probability.

### B. One-hot encoding

The masked MSA is encoded as three matrices using one-hot encoding from three different views. Specifically, for the  $j$ -th position of the  $i$ -th sequence in the masked MSA, we encode it as three one-hot vectors, i.e.,  $\mathbf{x}_{ij}$ ,  $\mathbf{y}_{ij}$ , and  $\mathbf{z}_{ij}$ , from the views of token type, row position, and column position, respectively.

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ijC_{max}}) \in R^{C_{max}}, x_{ijk} = \begin{cases} 1, & k = c_{ij} \\ 0, & k \neq c_{ij} \end{cases} \quad (1)$$

$$\mathbf{y}_{ij} = (y_{ij1}, y_{ij2}, \dots, y_{ijM_{max}}) \in R^{M_{max}}, y_{ijk} = \begin{cases} 1, & k = i \\ 0, & k \neq i \end{cases} \quad (2)$$

$$\mathbf{z}_{ij} = (z_{ij1}, z_{ij2}, \dots, z_{ijL_{max}}) \in R^{L_{max}}, z_{ijk} = \begin{cases} 1, & k = j \\ 0, & k \neq j \end{cases} \quad (3)$$

where  $c_{ij}$  is the index of token type for the  $j$ -th position of the  $i$ -th sequence,  $C_{max}$  is the number of types of tokens,  $L_{max}$  and  $M_{max}$  are preset maximum values for sequence length and alignments, respectively. In this work,  $C_{max} = 28$  and  $L_{max} = M_{max} = 1024$ , where 28 types of tokens include 20 common amino acids, 6 non-common amino acids (B, J, O, U, X and Z), 1 gap token, and 1 “masking” token.

According to Eqs. 1-3, the masked MSA can be encoded as three matrices, i.e.,  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{Z}$ , through one-hot encoding from the view of token type, row position, and column position, respectively, where  $\mathbf{X} \in R^{M \times L \times C_{max}}$ ,  $\mathbf{Y} \in R^{M \times L \times M_{max}}$  and  $\mathbf{Z} \in R^{M \times L \times L_{max}}$ ,  $M$  is the number of alignments, and  $L$  is the length of individual sequence in the masked MSA.

### C. Initial embedding

Each one-hot coding matrix is multiplied by a weight matrix to generate the corresponding embedding matrix:

$$\mathbf{H}_{token} = \mathbf{X}\mathbf{W}_{token} = \begin{bmatrix} \mathbf{X}[1] \\ \mathbf{X}[2] \\ \dots \\ \mathbf{X}[M] \end{bmatrix} \mathbf{W}_{token} = \begin{bmatrix} \mathbf{X}[1]\mathbf{W}_{token} \\ \mathbf{X}[2]\mathbf{W}_{token} \\ \dots \\ \mathbf{X}[M]\mathbf{W}_{token} \end{bmatrix} \in R^{M \times L \times D} \quad (4)$$

$$\mathbf{X}[i] \in R^{L \times C_{max}}, \mathbf{W}_{token} \in R^{C_{max} \times D}$$

$$\mathbf{H}_{row} = \mathbf{Y}\mathbf{W}_{row} = \begin{bmatrix} \mathbf{Y}[1] \\ \mathbf{Y}[2] \\ \dots \\ \mathbf{Y}[M] \end{bmatrix} \mathbf{W}_{row} = \begin{bmatrix} \mathbf{Y}[1]\mathbf{W}_{row} \\ \mathbf{Y}[2]\mathbf{W}_{row} \\ \dots \\ \mathbf{Y}[M]\mathbf{W}_{row} \end{bmatrix} \in R^{M \times L \times D} \quad (5)$$

$$\mathbf{Y}[i] \in R^{L \times M_{max}}, \mathbf{W}_{row} \in R^{M_{max} \times D}$$

$$\mathbf{H}_{col} = \mathbf{Z}\mathbf{W}_{col} = \begin{bmatrix} \mathbf{Z}[1] \\ \mathbf{Z}[2] \\ \dots \\ \mathbf{Z}[M] \end{bmatrix} \mathbf{W}_{col} = \begin{bmatrix} \mathbf{Z}[1]\mathbf{W}_{col} \\ \mathbf{Z}[2]\mathbf{W}_{col} \\ \dots \\ \mathbf{Z}[M]\mathbf{W}_{col} \end{bmatrix} \in R^{M \times L \times D} \quad (6)$$

$$\mathbf{Z}[i] \in R^{L \times L_{max}}, \mathbf{W}_{col} \in R^{L_{max} \times D}$$

where  $\mathbf{X}[i]$ ,  $\mathbf{Y}[i]$  and  $\mathbf{Z}[i]$  are the one-hot coding matrices for the  $i$ -th sequence in the masked MSA from the view of token type, row position, and column position, respectively,  $\mathbf{H}_{token}$ ,  $\mathbf{H}_{row}$ , and  $\mathbf{H}_{col}$  are token type-based, row position-based, and column position-based embedding matrices for the masked MSA, respectively, and  $D$  is the embedding dimension. In this work,  $D = 768$ .

Three embedding matrices are added as an initial embedding matrix  $\mathbf{H}_{init}$ :

$$\mathbf{H}_{init} = \mathbf{H}_{token} + \mathbf{H}_{row} + \mathbf{H}_{col}, \mathbf{H}_{init} \in R^{M \times L \times D} \quad (7)$$

#### D. Batch normalization and dropout

The initial embedding matrix  $\mathbf{H}_{init}$  is fed to the batch normalization layer to generate the corresponding normalized matrix  $\mathbf{H}_1$ :

$$\mathbf{H}_1 = BN(\mathbf{H}_{init}) = \begin{bmatrix} BN(\mathbf{h}_{11}) & \dots & BN(\mathbf{h}_{1L}) \\ \vdots & \ddots & \vdots \\ BN(\mathbf{h}_{M1}) & \dots & BN(\mathbf{h}_{ML}) \end{bmatrix} \quad (8)$$

$$BN(\mathbf{h}_{ij}) = \gamma \cdot \frac{\mathbf{h}_{ij} - u_{ij}}{\sqrt{\sigma_{ij}^2 + \epsilon}} + \beta, \mathbf{h}_{ij} \in R^D \quad (9)$$

where  $\mathbf{h}_{ij}$  is the initial embedding vector for the  $j$ -th position of the  $i$ -th sequence in the masked MSA,  $u_{ij}$  and  $\sigma_{ij}^2$  are mean and variance for  $\mathbf{h}_{ij}$ , respectively, and  $\gamma$ ,  $\beta$ , and  $\epsilon$  are normalized factors.

The normalized matrix  $\mathbf{H}_1$  is fed to dropout layer:

$$\mathbf{H}_1 \leftarrow \text{dropout}(\mathbf{H}_1, r) \quad (10)$$

where  $r$  is the rate of neurons which are randomly dropped in each training step, indicating that the corresponding weight vectors will be not optimized.

### E. Self-attention

The initial embedding matrix  $\mathbf{H}_1$  is fed to the self-attention network with  $N$  blocks, each of which consists of three sub-blocks. In this work,  $N = 12$ .

The first sub-block consists of a batch normalization layer, a row attention layer, a dropout layer, and a short connection, as follows.

$$\mathbf{H}_k^B = \text{BN}(\mathbf{H}_k) \quad (11)$$

$$\mathbf{H}_k^R = \text{RA}(\mathbf{H}_k^B) \quad (12)$$

$$\mathbf{H}_k^R \leftarrow \text{dropout}(\mathbf{H}_k^R, r) \quad (13)$$

$$\mathbf{F}_k = \text{SC}(\mathbf{H}_k, \mathbf{H}_k^R) = \mathbf{H}_k + \mathbf{H}_k^R \quad (14)$$

where  $\mathbf{H}_k$  and  $\mathbf{F}_k$  are the input and output matrices in the first sub-block of the  $k$ -th self-attention block, respectively,  $\text{BN}(\cdot)$  is the batch normalization function (see Eqs. 8-9),  $\text{SC}(\cdot)$  is the short connection, and  $\text{RA}(\cdot)$  is the row attention layer (see Eqs. 23-30),  $\mathbf{H}_k, \mathbf{H}_k^B, \mathbf{H}_k^R, \mathbf{F}_k \in R^{M \times L \times D}$ .

The second sub-block consists of a batch normalization layer, a column attention layer, a dropout layer, and a short connection, as follows.

$$\mathbf{F}_k^B = \text{BN}(\mathbf{F}_k) \quad (15)$$

$$\mathbf{F}_k^C = \text{CA}(\mathbf{F}_k^B) \quad (16)$$

$$\mathbf{F}_k^C \leftarrow \text{dropout}(\mathbf{F}_k^C, r) \quad (17)$$

$$\mathbf{U}_k = \text{SC}(\mathbf{F}_k, \mathbf{F}_k^C) = \mathbf{F}_k + \mathbf{F}_k^C \quad (18)$$

where  $\mathbf{F}_k$  and  $\mathbf{U}_k$  are the input and output matrices in the second sub-block of the  $k$ -th self-attention block, respectively,  $\text{CA}(\cdot)$  is the column attention layer (see Eqs. 31-39), and  $\mathbf{F}_k^B, \mathbf{F}_k^C, \mathbf{U}_k \in R^{M \times L \times D}$ .

The last sub-block consists of a batch normalization layer, a feed-forward network, a dropout layer, and a short connection, as follows.

$$\mathbf{U}_k^B = \text{BN}(\mathbf{U}_k) \quad (19)$$

$$\mathbf{U}_k^F = \text{FFN}(\mathbf{U}_k^B) \quad (20)$$

$$\mathbf{U}_k^F \leftarrow \text{dropout}(\mathbf{U}_k^F, r) \quad (21)$$

$$\mathbf{H}_{k+1} = \text{SC}(\mathbf{U}_k, \mathbf{U}_k^F) = \mathbf{U}_k + \mathbf{U}_k^F \quad (22)$$

where  $\mathbf{U}_k$  and  $\mathbf{H}_{k+1}$  are the input and output matrices in the third sub-block of the

$k$ -th self-attention block, respectively,  $FFN(\cdot)$  is the feed-forward network (see Eqs. 40-45), and  $\mathbf{U}_k^B, \mathbf{U}_k^F, \mathbf{H}_{k+1} \in R^{M \times L \times D}$ .

### (A) Row attention

Each row attention layer consists of  $m$  attention heads and a linear unit, where  $m = 12$ . In each attention head, the input matrix is multiplied by three weight matrices to generate the corresponding Query, Key, and Value matrices.

$$\mathbf{Q}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{QR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{QR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{QR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{QR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{QR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (23)$$

$$\mathbf{K}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{KR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{KR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{KR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{KR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{KR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (24)$$

$$\mathbf{V}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{VR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{VR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{VR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{VR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{VR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (25)$$

$$\mathbf{H}_k^B [i] \in R^{L \times D}, \mathbf{W}_{kt}^{QR}, \mathbf{W}_{kt}^{KR}, \mathbf{W}_{kt}^{VR} \in R^{D \times (\frac{D}{m})}$$

where  $\mathbf{H}_k^B$  is the input matrix of row attention layer in the  $k$ -th self-attention block (See Eq. 12),  $\mathbf{Q}_{kt}^R$ ,  $\mathbf{K}_{kt}^R$ , and  $\mathbf{V}_{kt}^R$  are Query, Key, and Value matrices in the  $t$ -th head of the row attention layer in the  $k$ -th block, respectively,  $\mathbf{W}_{kt}^{QR}$ ,  $\mathbf{W}_{kt}^{KR}$ , and  $\mathbf{W}_{kt}^{VR}$  are corresponding weight matrices.

Then, the dot-product between  $\mathbf{Q}_{kt}^R$  and  $\mathbf{K}_{kt}^R$  is performed and then normalized by SoftMax function to generate a row attention weight matrix:

$$\mathbf{W}_{kt}^{AR} = \text{SoftMax}\left(\frac{\sum_{i=1}^M \mathbf{Q}_{kt}^R [i] \cdot (\mathbf{K}_{kt}^R [i])^T}{\sqrt{MD/m}}\right) \in R^{L \times L}, \mathbf{Q}_{kt}^R [i], \mathbf{K}_{kt}^R [i] \in R^{L \times (D/m)} \quad (26)$$

$$\mathbf{W}_{kt}^{AR} \leftarrow \text{dropout}(\mathbf{W}_{kt}^{AR}, r) \quad (27)$$

where  $\mathbf{W}_{kt}^{AR}$  is the attention weight matrix in the  $t$ -th head of the row attention layer in the  $k$ -th block and measures the correlation for each pair of columns in the masked MSA.

Next, the row attention weight matrix  $\mathbf{W}_{kt}^{AR}$  is multiplied by Value matrix  $\mathbf{V}_{kt}^R$  to generate the corresponding row attention matrix:

$$\mathbf{A}_{kt}^R = \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R = \mathbf{W}_{kt}^{AR} \begin{bmatrix} \mathbf{V}_{kt}^R[1] \\ \mathbf{V}_{kt}^R[2] \\ \dots \\ \mathbf{V}_{kt}^R[M] \end{bmatrix} = \begin{bmatrix} \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[1] \\ \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[2] \\ \dots \\ \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[M] \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})}, \mathbf{V}_{kt}^R[i] \in R^{L \times (\frac{D}{m})} \quad (28)$$

where  $\mathbf{A}_{kt}^R$  is the attention matrix in the  $t$ -th head of the row attention layer in the  $k$ -th block.

Finally, the outputs of all attention heads are concatenated as a new matrix, which is further fed to a linear unit:

$$\mathbf{A}_k^R = \mathbf{A}_{k1}^R \mathbf{A}_{k2}^R \dots \mathbf{A}_{km}^R \in R^{M \times L \times D} \quad (29)$$

$$\mathbf{H}_k^R = \mathbf{A}_k^R \mathbf{W}_k^R + \mathbf{b}_k^R = \begin{bmatrix} \mathbf{A}_k^R[1] \\ \mathbf{A}_k^R[2] \\ \dots \\ \mathbf{A}_k^R[M] \end{bmatrix} \mathbf{W}_k^R + \mathbf{b}_k^R = \begin{bmatrix} \mathbf{A}_k^R[1] \mathbf{W}_k^R \\ \mathbf{A}_k^R[2] \mathbf{W}_k^R \\ \dots \\ \mathbf{A}_k^R[M] \mathbf{W}_k^R \end{bmatrix} + \mathbf{b}_k^R \in R^{M \times L \times D} \quad (30)$$

$$\mathbf{W}_k^R \in R^{D \times D}, \mathbf{A}_k^R[i] \in R^{L \times D}$$

where  $\mathbf{H}_k^R$  is the output matrix of row attention layer in the  $k$ -th attention block (See Eq. 12), and  $\mathbf{W}_k^R$  and  $\mathbf{b}_k^R$  are weight matrix and bias in the linear unit, respectively.

## (B) Column attention

Each column attention layer consists of  $m$  attention heads and a linear unit. In each attention head, the input matrix is multiplied by three weight matrices to generate the corresponding Query, Key, and Value matrices.

$$\mathbf{Q}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{QC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{QC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{QC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{QC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{QC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (31)$$

$$\mathbf{K}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{KC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{KC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{KC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{KC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{KC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (32)$$

$$\mathbf{V}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{VC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{VC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{VC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{VC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{VC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (33)$$

$$\mathbf{F}_k^B[i] \in R^{L \times D}, \mathbf{W}_{kt}^{QC}, \mathbf{W}_{kt}^{KC}, \mathbf{W}_{kt}^{VC} \in R^{D \times (\frac{D}{m})}$$

where  $\mathbf{F}_k^B$  is the input matrix of column attention layer in the  $k$ -th self-attention block (see Eq. 16),  $\mathbf{Q}_{kt}^C$ ,  $\mathbf{K}_{kt}^C$ , and  $\mathbf{V}_{kt}^C$  are Query, Key, and Value matrices in the  $t$ -th head of column attention layer in the  $k$ -th block, respectively,  $\mathbf{W}_{kt}^{QC}$ ,  $\mathbf{W}_{kt}^{KC}$ , and  $\mathbf{W}_{kt}^{VC}$  are

corresponding weight metrics.

Then, the dot-product between  $\mathbf{Q}_{kt}^C$  and  $\mathbf{K}_{kt}^C$  is performed and then normalized by SoftMax function to generate an attention weight matrix:

$$\mathbf{W}_{kt}^{AC} = \text{SoftMax}\left(\frac{\mathbf{Q}_{kt}^C(\mathbf{K}_{kt}^C)^T}{\sqrt{D/m}}\right) \in R^{M \times L \times M} \quad (34)$$

$$\mathbf{W}_{kt}^{AC} \leftarrow \text{dropout}(\mathbf{W}_{kt}^{AC}, r) \quad (35)$$

$$\begin{aligned} \mathbf{Q}_{kt}^C(\mathbf{K}_{kt}^C)^T &= [\mathbf{Q}_{kt}^C[:, 1, :] \mathbf{Q}_{kt}^C[:, 2, :] \dots \mathbf{Q}_{kt}^C[:, L, :]] \cdot [\mathbf{K}_{kt}^C[:, 1, :] \mathbf{K}_{kt}^C[:, 2, :] \dots \mathbf{K}_{kt}^C[:, L, :]]^T = \\ &[\mathbf{Q}_{kt}^C[:, 1, :] \cdot \mathbf{K}_{kt}^C[:, 1, :]}^T \mathbf{Q}_{kt}^C[:, 2, :] \cdot \mathbf{K}_{kt}^C[:, 2, :]}^T \dots \mathbf{Q}_{kt}^C[:, L, :] \cdot \mathbf{K}_{kt}^C[:, L, :]}^T] \in R^{M \times L \times M} \end{aligned} \quad (36)$$

$$\mathbf{Q}_{kt}^C[:, j, :], \mathbf{K}_{kt}^C[:, j, :] \in R^{M \times (\frac{D}{m})}, \mathbf{Q}_{kt}^C[:, j, :] \cdot \mathbf{K}_{kt}^C[:, j, :]}^T \in R^{M \times M}$$

where  $\mathbf{W}_{kt}^{AC}$  is the attention weight matrix in the  $t$ -th head of column attention layer in the  $k$ -th block, and  $\mathbf{W}_{kt}^{AC}[:, j, :]$  measures the correlation for each pair of alignments at the  $j$ -th position.

Next, the column attention weight matrix  $\mathbf{W}_{kt}^{AC}$  is multiplied by Value matrix  $\mathbf{V}_{kt}^C$  to generate the corresponding column attention matrix:

$$\begin{aligned} \mathbf{A}_{kt}^C = \mathbf{W}_{kt}^{AC} \mathbf{V}_{kt}^C &= [\mathbf{W}_{kt}^{AC}[:, 1, :] \mathbf{W}_{kt}^{AC}[:, 2, :] \dots \mathbf{W}_{kt}^{AC}[:, L, :]] \cdot [\mathbf{V}_{kt}^C[:, 1, :] \mathbf{V}_{kt}^C[:, 2, :] \dots \mathbf{V}_{kt}^C[:, L, :]] = [\mathbf{W}_{kt}^{AC}[:, 1, :]} \cdot \\ &\mathbf{V}_{kt}^C[:, 1, :]} \mathbf{W}_{kt}^{AC}[:, 2, :]} \cdot \mathbf{V}_{kt}^C[:, 2, :]} \dots \mathbf{W}_{kt}^{AC}[:, L, :]} \cdot \mathbf{V}_{kt}^C[:, L, :]}] \in R^{M \times L \times (\frac{D}{m})} \end{aligned} \quad (37)$$

$$\mathbf{W}_{kt}^{AC}[:, j, :] \in R^{M \times M}, \mathbf{V}_{kt}^C[:, j, :] \in R^{M \times (\frac{D}{m})}, \mathbf{W}_{kt}^{AC}[:, j, :]} \cdot \mathbf{V}_{kt}^C[:, j, :]} \in R^{M \times (\frac{D}{m})}$$

where  $\mathbf{A}_{kt}^C$  is the attention matrix in the  $t$ -th head of column attention layer in the  $k$ -th block.

Finally, the outputs of all attention heads are concatenated as a new matrix, which is further fed to a linear unit:

$$\mathbf{A}_k^C = \mathbf{A}_{k1}^C \mathbf{A}_{k2}^C \dots \mathbf{A}_{km}^C \in R^{M \times L \times D} \quad (38)$$

$$\mathbf{F}_k^C = \mathbf{A}_k^C \mathbf{W}_k^C + \mathbf{b}_k^C = \begin{bmatrix} \mathbf{A}_k^C[1] \\ \mathbf{A}_k^C[2] \\ \dots \\ \mathbf{A}_k^C[M] \end{bmatrix} \mathbf{W}_k^C = \begin{bmatrix} \mathbf{A}_1^C[1] \mathbf{W}_k^C \\ \mathbf{A}_2^C[2] \mathbf{W}_k^C \\ \dots \\ \mathbf{A}_k^C[M] \mathbf{W}_k^C \end{bmatrix} + \mathbf{b}_k^C \in R^{M \times L \times D} \quad (39)$$

$$\mathbf{W}_k^C \in R^{D \times D}, \mathbf{A}_k^C[i] \in R^{L \times D}$$

where  $\mathbf{F}_k^C$  in the output matrix of column attention layer in the  $k$ -th attention block, (See Eq. 16), and  $\mathbf{W}_k^C$  and  $\mathbf{b}_k^C$  are weight matrix and bias in the linear unit, respectively.



### (C) Feed-forward network

$$\mathbf{T}_k^F = \text{gelu}(\mathbf{U}_k^B \mathbf{W}_k^1 + \mathbf{b}_k^1) \in R^{M \times L \times D_1} \quad (40)$$

$$\mathbf{T}_k^F \leftarrow \text{dropout}(\mathbf{T}_k^F, r) \quad (41)$$

$$\mathbf{U}_k^F = \mathbf{T}_k^F \mathbf{W}_k^2 + \mathbf{b}_k^2 \in R^{M \times L \times D} \quad (42)$$

$$\text{gelu}(x) = x\phi(x) \quad (43)$$

$$\mathbf{U}_k^B \mathbf{W}_k^1 = \begin{bmatrix} \mathbf{U}_k^B[1] \\ \mathbf{U}_k^B[2] \\ \dots \\ \mathbf{U}_k^B[M] \end{bmatrix} \mathbf{W}_k^1 = \begin{bmatrix} \mathbf{U}_k^B[1] \mathbf{W}_k^1 \\ \mathbf{U}_k^B[2] \mathbf{W}_k^1 \\ \dots \\ \mathbf{U}_k^B[M] \mathbf{W}_k^1 \end{bmatrix} \in R^{M \times L \times D_1} \quad (44)$$

$$\mathbf{T}_k^F \mathbf{W}_k^2 = \begin{bmatrix} \mathbf{T}_k^F[1] \\ \mathbf{T}_k^F[2] \\ \dots \\ \mathbf{T}_k^F[M] \end{bmatrix} \mathbf{W}_k^2 = \begin{bmatrix} \mathbf{T}_k^F[1] \mathbf{W}_k^2 \\ \mathbf{T}_k^F[2] \mathbf{W}_k^2 \\ \dots \\ \mathbf{T}_k^F[M] \mathbf{W}_k^2 \end{bmatrix} \in R^{M \times L \times D} \quad (45)$$

$$\mathbf{U}_k^B[i] \in R^{L \times D}, \mathbf{W}_k^1 \in R^{D \times D_1}, \mathbf{T}_k^F[i] \in R^{L \times D_1}, \mathbf{W}_k^2 \in R^{D_1 \times D}, D_1=3072$$

where  $\mathbf{U}_k^B$  and  $\mathbf{U}_k^F$  are the input and output matrices of feed-forward network in the  $k$ -th self-attention block, respectively, (see Eq. 20),  $\mathbf{W}_k^1$  and  $\mathbf{W}_k^2$  are weight matrices,  $\mathbf{b}_k^1$  and  $\mathbf{b}_k^2$  are bias, and  $\phi(x)$  is the integral of Gaussian Distribution for  $x$ .

### G. Output layer

The output of the last self-attention block is fed to a fully connected layer with SoftMax function to generate a probability matrix:

$$\mathbf{P} = \text{SoftMax}(\mathbf{H}_{N+1} \mathbf{W}^O + \mathbf{b}^O) \in R^{M \times L \times C_{max}} \quad (46)$$

$$\mathbf{H}_{N+1} \mathbf{W}^O = \begin{bmatrix} \mathbf{H}_{N+1}[1] \mathbf{W}^O \\ \mathbf{H}_{N+1}[2] \mathbf{W}^O \\ \dots \\ \mathbf{H}_{N+1}[M] \mathbf{W}^O \end{bmatrix}, \mathbf{H}_{N+1}[i] \in R^{L \times D}, \mathbf{W}^O \in R^{D \times C_{max}} \quad (47)$$

where  $\mathbf{H}_{N+1}$  is the outputted embedding matrix in the  $N$ -th self-attention block,  $\mathbf{W}^O$  and  $\mathbf{b}^O$  are weight matrix and bias, respectively, and the  $\mathbf{P}(i, j, c)$  indicates the probability that the  $j$ -th position of the  $i$ -th sequence in the masked MSA is predicted as the  $c$ -th type of amino acid.

### F. Loss function

For an individual MSA, the loss function is designed as:

$$\text{Loss}_{msa} = \frac{1}{M} \cdot \sum_{i=1}^M \left\{ \frac{1}{|\text{mask}(i)|} \cdot \sum_{j \in \text{mask}(i)} -\log \mathbf{P}_{i,j,c(i,j)} \right\} \quad (48)$$

where  $M$  is the number of alignments,  $mask(i)$  is a set of masking position in the  $i$ -th sequence,  $|mask(i)|$  is the number of elements in  $mask(i)$ ,  $c(i, j)$  is the type index of amino acid for the  $j$ -th position in the  $i$ -th sequence before masking, and  $-\log P_{i,j,c(i,j)}$  is negative log likelihood of the true amino acid at the  $j$ -th position in the  $i$ -th sequence under condition of masking.

## Supporting Tables

Table S1. Average precision in top  $N$  predicted contacts on TS630 dataset.

Methods	Top 1	Top 5	Top 10	Top 20	Top 50	Top 100
ICCPred	<b>0.267</b>	<b>0.245</b>	<b>0.240</b>	<b>0.231</b>	<b>0.217</b>	<b>0.198</b>
GLINTER	0.153	0.154	0.160	0.154	0.146	0.139
HDIContact	0.041	0.038	0.035	0.040	0.081	0.103

Table S2. Average precision in top  $L/K$  predicted contacts, AUPR, and AUROC on TS630 dataset. ALL represents the number of native contacts on the target.

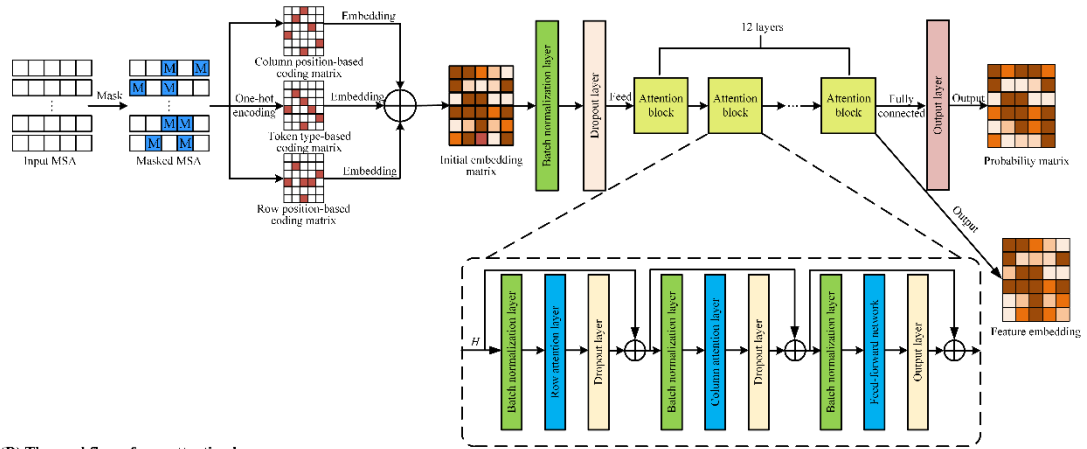
Methods	AUROC	AUPR	Top $L/30$	Top $L/20$	Top $L/10$	Top $L/5$	Top $L/2$	Top ALL
ICCPred	<b>0.725</b>	<b>0.115</b>	<b>0.238</b>	<b>0.233</b>	<b>0.225</b>	<b>0.210</b>	<b>0.181</b>	<b>0.133</b>
GLINTER	0.411	0.111	0.156	0.152	0.149	0.143	0.137	0.116
HDIContact	0.659	0.070	0.035	0.035	0.055	0.091	0.107	0.091

Table S3. Average contact precision of different features on the TS630 dataset.

Feature	Top										
	1	5	10	20	50	100	<i>L</i> /30	<i>L</i> /20	<i>L</i> /10	<i>L</i> /5	<i>L</i> /2
ESM2	0.046	0.039	0.036	0.038	0.039	0.038	0.039	0.038	0.038	0.038	0.037
CPX	0.251	0.233	0.218	0.209	0.199	0.181	0.218	0.214	0.205	0.193	0.167
GDS	0.190	0.177	0.170	0.167	0.157	0.145	0.171	0.168	0.161	0.152	0.133
PPIS	0.157	0.160	0.157	0.151	0.142	0.131	0.159	0.154	0.148	0.138	0.122
PIS	0.237	0.223	0.214	0.201	0.186	0.171	0.214	0.207	0.194	0.181	0.158
PP	0.208	0.203	0.200	0.194	0.179	0.166	0.202	0.198	0.187	0.175	0.154
GP	0.235	0.227	0.223	0.218	0.201	0.186	0.219	0.218	0.208	0.195	0.172
GI	0.195	0.190	0.188	0.181	0.172	0.160	0.186	0.184	0.178	0.167	0.149
GSP	<b>0.267</b>	<b>0.245</b>	<b>0.240</b>	<b>0.231</b>	<b>0.217</b>	<b>0.198</b>	<b>0.238</b>	<b>0.233</b>	<b>0.225</b>	<b>0.210</b>	<b>0.181</b>

# Supporting Figures

(A) The framework of ESM-MSA



(B) The workflow of row attention layer

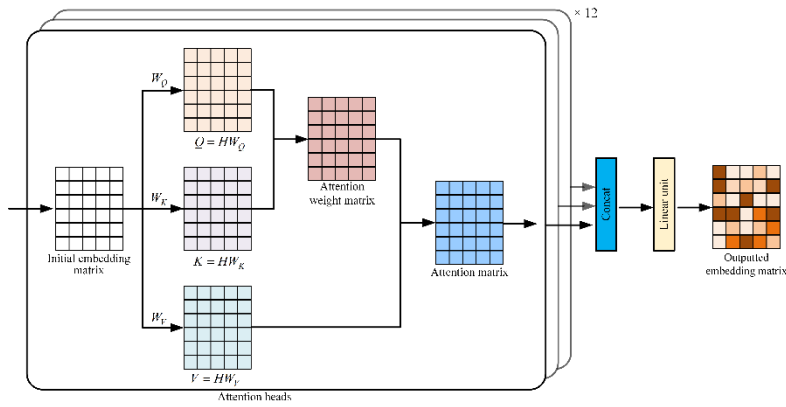
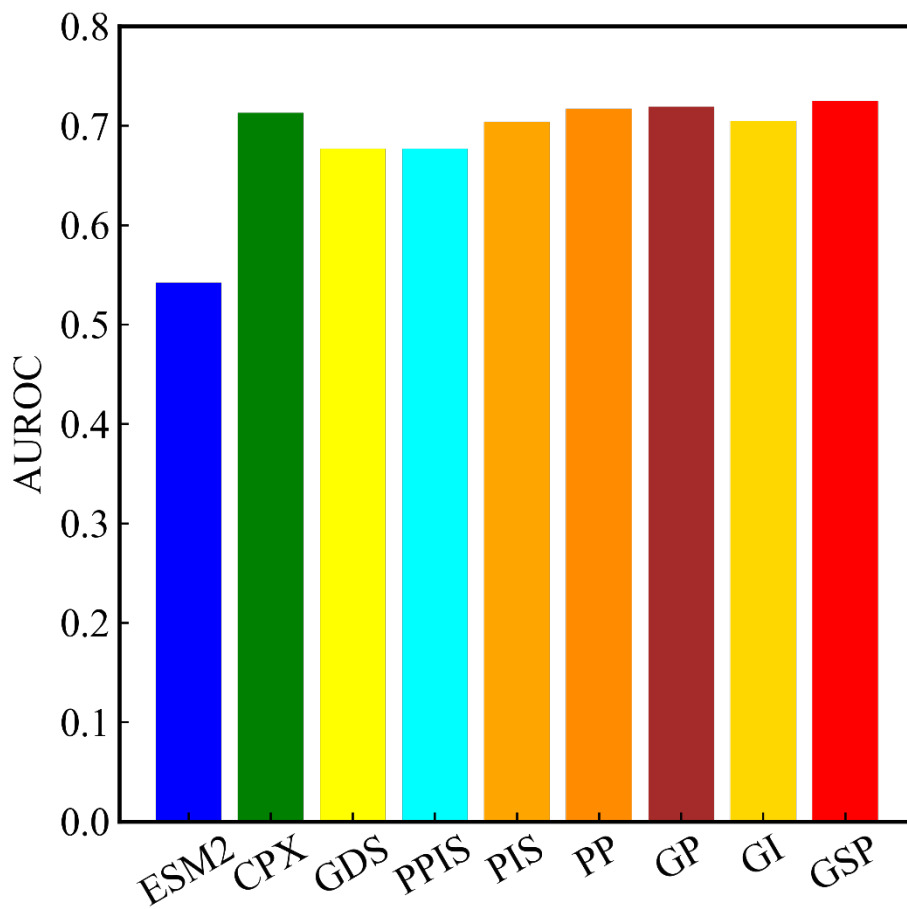


Figure S1. The workflow of ESM-MSA



**Figure S2.** Performance comparison between nine feature embeddings regarding AUROC on the TS630 dataset.