

Supplementary Materials

Table of content

Supporting Text

- Text S1. Definitions of DNA-binding sites.
- Text S2. Procedures for the ESM2 transformer.
- Text S3. Procedures for the ESM-MSA transformer.
- Text S4. p -value calculation between EPE and the other six feature embeddings.
- Text S5. Procedures for constructing PDNA-960 and PDNA-136 datasets.
- Text S6. Formulas for calculating average precision.

Supporting Figures

- Figure S1. The frequency distributions of 20 native amino acids among DNA-binding and non-DNA-binding sites for five benchmark datasets.
- Figure S2. The workflow of the ESM2 transformer.
- Figure S3. The workflow of the ESM-MSA transformer.
- Figure S4. The architectures of three ablation models.

Supporting Tables

- Table S1. The performance of five DNA-binding site predictors on PDNA-543 over ten-fold cross-validation.
- Table S2. The performance of five DNA-binding site predictors on PDNA-335 over ten-fold cross-validation.
- Table S3. The performance of eleven DNA-binding site predictors on PDNA-316 over ten-fold cross-validation.
- Table S4. The p -values of MCC values between EPE and the other six feature embeddings on five benchmark datasets, where the base model is the designed LSTM-attention network.
- Table S5. The p -values of AUROC values between EPE and the other six feature embeddings on five benchmark datasets, where the base model is the designed LSTM-attention network.
- Table S6. The predicted and native DNA-binding sites of two representative proteins for five DNA-binding prediction methods.

Supporting Text

Text S1. Definitions of DNA-binding sites

In the PDNA-335 and PDNA-52 datasets, the DNA-binding sites are defined using the criterion of the protein-ligand binding database BioLip [1], consistent with that used in the Critical Assessment of Structure Prediction (CASP) [2, 3]. Specifically, a protein residue with at least an inter-molecular atomic contact to a DNA molecule is labeled as a DNA binding site, where the inter-molecular atomic contact is a non-hydrogen protein-DNA atom pair whose Euclidian distance is less than the sum of van der Waals radii plus 0.5 Å. It is noted that this criterion has been tightened up in the newest BioLip database (i.e., BioLip2) [4], in which a protein residue with at least two inter-molecular atomic contacts to a DNA is defined as a DNA binding site.

In the PDNA-316 dataset, the DNA-binding sites are defined by the criteria of Ahmad's work [5], which is the first work for protein-DNA binding site prediction (to our best knowledge). Specifically, in a protein-DNA complex, an amino acid residue in the protein is defined as a DNA-binding site if the distance between any atoms of this residue and any atoms of the DNA molecule is less than a cut-off value. In the PDNA-316, this cut-off is set to be 3.5 Å.

PDNA-543 and PDNA-41 were constructed by the TargetDNA paper [6], which does not provide any details for defining DNA-binding sites. However, we infer that the definition of DNA-binding sites for these two datasets is the same as that for the PDNA-316 dataset via the criteria of Ahmad's work due to the following observations. Specifically, there are 101 overlap proteins between the PDNA-543 and PDNA-316 datasets, where 98.0% overlap proteins have consistent DNA-binding site annotations between these two datasets. Meanwhile, there are 91 overlap proteins between the PDNA-543 and PDNA-335 datasets, but only 12.1% overlap proteins have consistent DNA-binding site annotations.

To further demonstrate our inference, we designed the following computational experiments. Specifically, we separately used different criteria to re-define the DNA-binding sites for the protein chains in the PDNA-543 and PDNA-41 datasets and then calculated the consistency ratio between re-defined DNA-binding sites and originally annotated DNA-binding sites. The higher consistency ratio means that the re-implemented criterion is more similar to the original criterion used in the TargetDNA paper. It is noted that there are 86 protein chains whose DNA-binding sites cannot be

re-defined in the PDNA-543 and PDNA-41. The underlying reason is that these proteins have been obsoleted or updated in the PDB database from the year 2016 to now. As a result, their PDB structures are unavailable, or the sequences extracted from their structures are inconsistent with the original sequences provided by the TargetDNA paper. Therefore, we only collected 498 protein chains from the PDNA-543 and PDNA-41 whose consistency ratio between re-defined and originally annotated DNA-binding sites could be calculated. Here, the consistency ratio (denoted as CR) is defined as $CR=N1/N2$, where N1 is the number of protein chains whose re-defined DNA-binding sites are consistent with the originally annotated DNA-binding sites in the TargetDNA paper; N2 is the total number of available protein chains (i.e., $N2=498$). In our experiments, we used four different criteria to re-define DNA-binding sites, including the criterion of CASP and the criterion of Ahmad’s work with the cutoff values of 3.0 Å, 3.5 Å, and 4.0 Å, where the corresponding CR values are 9.8%, 0.0%, 89.6%, and 4.2%, respectively. These experiment data further demonstrated that the TargetDNA paper probably used the criterion of Ahmad’s work with the cut-off=3.5 Å to define DNA-binding sites. It could not escape our notice that there is nearly a 10% difference between the re-defined and originally annotated DNA-binding sites. This difference may be due to that the PDB structures for part of proteins have been updated from 2016 to now. The source codes and experiment data for re-defining DNA-binding sites in 498 protein chains from the PDNA-543 and PDNA-41 datasets are available at https://github.com/yiheng-zhu/ULDNA/tree/main/Check_TargetDNA.

Text S2. Procedures for the ESM2 transformer

A. Masking

For an input sequence, the masking strategy [7] is performed on the corresponding tokens (i.e., amino acids). Specifically, we randomly sample 15% tokens, 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 sequence is represented as a $L \times 28$ matrix using one-hot encoding [8], where 28 is the types of tokens, including 20 common amino acids, 6 non-common amino acids (B, J, O, U, X, and Z), 1 gap token, and 1 “masking” token.

C. Embedding with position information

The one-hot coding matrix X of the masked sequence is multiplied by an embedding weight matrix W_E to generate an embedding matrix H_E :

$$H_E = XW_E, X \in R^{L \times 28}, W_E \in R^{28 \times D}, H_E \in R^{L \times D} \quad (S1)$$

where L is the length of the masked sequence, 28 is the types of tokens in the masked sequence, and D is the embedding dimension.

Then, the position embedding strategy is used to record the position of each token in the masked sequence to generate a position embedding matrix H_P :

$$H_P = \begin{bmatrix} h_1 \\ h_2 \\ \dots \\ h_L \end{bmatrix}, h_i = (v_{i,1}, v_{i,2}, \dots, v_{i,D}), H_P \in R^{L \times D}, \text{ and } h_i \in R^D \quad (S2)$$

$$v_{i,2k} = \sin\left(\frac{i}{10000^{2k/D}}\right), v_{i,2k+1} = \cos\left(\frac{i}{10000^{(2k+1)/D}}\right), k = 0, 1, \dots, (D-1)/2 \quad (S3)$$

where h_i is the embedding vector for the i -th position in the masked sequence.

Finally, two embedding matrices are added as a combination embedding matrix H_1 :

$$H_1 = H_E + H_P, H_1 \in R^{L \times D} \quad (S4)$$

D. Self-attention

The embedding matrix H_1 is fed to a self-attention block with n layers, each of which consists of m attention heads, a linear unit, and a feed-forward network (FFN). In each attention head, the scale dot-product attention is performed as follows:

$$A_{i,j} = \text{Softmax}(M_{i,j}^Q M_{i,j}^{K^T} / \sqrt{d_{ij}}) M_{i,j}^V \quad (S5)$$

$$M_{i,j}^Q = H_i W_{i,j}^Q, M_{i,j}^K = H_i W_{i,j}^K, M_{i,j}^V = H_i W_{i,j}^V \quad (S6)$$

$$d_{ij} = D/m, W_{i,j}^Q, W_{i,j}^K, W_{i,j}^V \in R^{D \times (\frac{D}{m})}, M_{i,j}^Q, M_{i,j}^K, M_{i,j}^V, A_{i,j} \in R^{L \times (\frac{D}{m})} \quad (S7)$$

where $A_{i,j}$ is the attention matrix in the (i -th layer, j -th head) and measures the evolution correlation for each amino acid pair in the sequence, $M_{i,j}^Q$, $M_{i,j}^K$, and $M_{i,j}^V$ are Query, Key, and Value matrices in the (i -th layer, j -th head), H_i is the input matrix in the i -th layer, $W_{i,j}^Q$, $W_{i,j}^K$, and $W_{i,j}^V$ are weight matrices, and d_{ij} is the scale parameter.

The outputs of all attention heads in the i -th layer are concatenated as a new matrix A_i , which is further fed to a linear unit to output the matrix U_i :

$$A_i = A_{i,1} A_{i,2} \dots A_{i,m} \quad (S8)$$

$$U_i = A_i W_i^1 + b_i^1, W_i^1 \in R^{D \times D}, A_i, b_i^1, U_i \in R^{L \times D} \quad (\text{S9})$$

where W_i^1 and b_i^1 are the weight matrix and bias, respectively, in the linear unit.

E. Feed-forward network with shortcut connections

The U_i is added by H_i to generate a new matrix F_i , which is further fed to the FFN to output the matrix T_i :

$$F_i = H_i + U_i \quad (\text{S10})$$

$$T_i = \text{gelu}(F_i W_i^2 + b_i^2) W_i^3 + b_i^3, W_i^2, W_i^3 \in R^{D \times D}, b_i^2, b_i^3, T_i \in R^{L \times D} \quad (\text{S11})$$

$$\text{gelu}(x) = x \Phi(x) \quad (\text{S12})$$

where W_i^2 and W_i^3 are weight matrices in the FFN, b_i^2 and b_i^3 are bias in the FFN, and $\Phi(x)$ is the integral of Gaussian Distribution for x

The F_i is added by T_i as the output of the i -th attention layer:

$$H_{i+1} = F_i + T_i, H_{i+1} \in R^{L \times D} \quad (\text{S13})$$

where H_{i+1} is the evolution diversity-based embedding matrix in the i -th attention layer.

The output of the last attention layer is fed to a fully connected layer with SoftMax function to generate a $L \times 28$ probability matrix:

$$P = \text{SoftMax}(H^n W^n + b^n), P \in R^{L \times 28} \quad (\text{S14})$$

where the (l -th, c -th) value in P indicates the probability that the l -th token in the masked sequence is predicted as the c -th type of amino acid, W^n and b^n are weight matrix and bias, respectively.

F. Loss function

The loss function is designed as a negative log-likelihood function between inputted one-hot and outputted probability matrices, to ensure that the prediction model correctly predicts the true amino acids in the masked position as much as possible:

$$\text{Loss}_{esm} = E_{x \sim X} \sum_{l \in x(M)} \left(-\frac{\log P_{l,c(l)}}{|x(M)|} \right) \quad (\text{S15})$$

where x is a sequence in training protein set X , $x(M)$ is a set of masking positions in x , $|x(M)|$ is the number of elements in $x(M)$, $c(l)$ is the type index of amino acid for the l -th token in x before masking, and $-\log P_{l,c(l)}$ is the negative log-likelihood of the true amino acid x_l under the condition of masking.

The ESM2 transformer is optimized by minimizing the loss function via Adam optimization algorithm [9]. Then, the output of the last attention layer is represented as

a $L \times D$ matrix, as the evolution diversity-based embedding for DNA-binding site prediction, where D is the number of neurons of FFN. The current ESM2 model with 3 billion parameters was trained over 60 million proteins from the UniRef50 database and can be freely downloaded at <https://github.com/facebookresearch/esm>, where $n = 36$, $m = 20$, and $D = 2560$.

Text S3. Procedures for the ESM-MSA transformer

A. Masking

For an input multiple sequence alignment (MSA), the masking strategy is performed. Specifically, for each individual 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.

$$\mathbf{x}_{ij} = (x_{ij1}, x_{ij2}, \dots, x_{ijC_{max}}) \in R^{C_{max}}, x_{ijk} = \begin{cases} 1, & k = c_{ij} \\ 0, & k \neq c_{ij} \end{cases} \quad (\text{S16})$$

$$\mathbf{y}_{ij} = (y_{ij1}, y_{ij2}, \dots, y_{ijM_{max}}) \in R^{M_{max}}, y_{ijk} = \begin{cases} 1, & k = i \\ 0, & k \neq i \end{cases} \quad (\text{S17})$$

$$\mathbf{z}_{ij} = (z_{ij1}, z_{ij2}, \dots, z_{ijL_{max}}) \in R^{L_{max}}, z_{ijk} = \begin{cases} 1, & k = j \\ 0, & k \neq j \end{cases} \quad (\text{S18})$$

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. S16-S18, 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:

$$\mathbf{H}_{token} = \mathbf{X}\mathbf{W}_{token} = \begin{bmatrix} \mathbf{X}[1] \\ \mathbf{X}[2] \\ \dots \\ \mathbf{X}[M] \end{bmatrix} \mathbf{W}_{token} = \begin{bmatrix} \mathbf{X}[1]\mathbf{W}_{token} \\ \mathbf{X}[2]\mathbf{W}_{token} \\ \dots \\ \mathbf{X}[M]\mathbf{W}_{token} \end{bmatrix} \in R^{M \times L \times D} \quad (\text{S19})$$

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

$$\mathbf{H}_{row} = \mathbf{Y}\mathbf{W}_{row} = \begin{bmatrix} \mathbf{Y}[1] \\ \mathbf{Y}[2] \\ \dots \\ \mathbf{Y}[M] \end{bmatrix} \mathbf{W}_{row} = \begin{bmatrix} \mathbf{Y}[1]\mathbf{W}_{row} \\ \mathbf{Y}[2]\mathbf{W}_{row} \\ \dots \\ \mathbf{Y}[M]\mathbf{W}_{row} \end{bmatrix} \in R^{M \times L \times D} \quad (\text{S20})$$

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

$$\mathbf{H}_{col} = \mathbf{Z}\mathbf{W}_{col} = \begin{bmatrix} \mathbf{Z}[1] \\ \mathbf{Z}[2] \\ \dots \\ \mathbf{Z}[M] \end{bmatrix} \mathbf{W}_{col} = \begin{bmatrix} \mathbf{Z}[1]\mathbf{W}_{col} \\ \mathbf{Z}[2]\mathbf{W}_{col} \\ \dots \\ \mathbf{Z}[M]\mathbf{W}_{col} \end{bmatrix} \in R^{M \times L \times D} \quad (\text{S21})$$

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

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} :

$$\mathbf{H}_{init} = \mathbf{H}_{token} + \mathbf{H}_{row} + \mathbf{H}_{col}, \mathbf{H}_{init} \in R^{M \times L \times D} \quad (\text{S22})$$

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 :

$$\mathbf{H}_1 = \text{BN}(\mathbf{H}_{init}) = \begin{bmatrix} \text{BN}(\mathbf{h}_{11}) & \dots & \text{BN}(\mathbf{h}_{1L}) \\ \vdots & \ddots & \vdots \\ \text{BN}(\mathbf{h}_{M1}) & \dots & \text{BN}(\mathbf{h}_{ML}) \end{bmatrix} \quad (\text{S23})$$

$$\text{BN}(\mathbf{h}_{ij}) = \gamma \cdot \frac{\mathbf{h}_{ij} - u_{ij}}{\sqrt{\sigma_{ij}^2 + \epsilon}} + \beta, \mathbf{h}_{ij} \in R^D \quad (\text{S24})$$

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 σ_{ij}^2 are mean and variance for \mathbf{h}_{ij} , respectively, and γ , β ,

and ϵ are normalized factors.

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

$$\mathbf{H}_1 \leftarrow \text{dropout}(\mathbf{H}_1, r) \quad (\text{S25})$$

where r is the rate of neurons that 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.

$$\mathbf{H}_k^B = \text{BN}(\mathbf{H}_k) \quad (\text{S26})$$

$$\mathbf{H}_k^R = \text{RA}(\mathbf{H}_k^B) \quad (\text{S27})$$

$$\mathbf{H}_k^R \leftarrow \text{dropout}(\mathbf{H}_k^R, r) \quad (\text{S28})$$

$$\mathbf{F}_k = \text{SC}(\mathbf{H}_k, \mathbf{H}_k^R) = \mathbf{H}_k + \mathbf{H}_k^R \quad (\text{S29})$$

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. S23-S24), $\text{SC}(\cdot)$ is the short connection, and $\text{RA}(\cdot)$ is the row attention layer (see Eqs. S38-S45), $\mathbf{H}_k, \mathbf{H}_k^B, \mathbf{H}_k^R, \mathbf{F}_k \in \mathbb{R}^{M \times L \times D}$.

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

$$\mathbf{F}_k^B = \text{BN}(\mathbf{F}_k) \quad (\text{S30})$$

$$\mathbf{F}_k^C = \text{CA}(\mathbf{F}_k^B) \quad (\text{S31})$$

$$\mathbf{F}_k^C \leftarrow \text{dropout}(\mathbf{F}_k^C, r) \quad (\text{S32})$$

$$\mathbf{U}_k = \text{SC}(\mathbf{F}_k, \mathbf{F}_k^C) = \mathbf{F}_k + \mathbf{F}_k^C \quad (\text{S33})$$

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. S46-S54), and $\mathbf{F}_k^B, \mathbf{F}_k^C, \mathbf{U}_k \in \mathbb{R}^{M \times L \times D}$.

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

$$\mathbf{U}_k^B = \text{BN}(\mathbf{U}_k) \quad (\text{S34})$$

$$\mathbf{U}_k^F = \text{FFN}(\mathbf{U}_k^B) \quad (\text{S35})$$

$$\mathbf{U}_k^F \leftarrow \text{dropout}(\mathbf{U}_k^F, r) \quad (\text{S36})$$

$$\mathbf{H}_{k+1} = SC(\mathbf{U}_k, \mathbf{U}_k^F) = \mathbf{U}_k + \mathbf{U}_k^F \quad (\text{S37})$$

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. S55-S60), 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.

$$\mathbf{Q}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{QR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{QR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{QR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{QR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{QR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S38})$$

$$\mathbf{K}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{KR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{KR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{KR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{KR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{KR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S39})$$

$$\mathbf{V}_{kt}^R = \mathbf{H}_k^B \mathbf{W}_{kt}^{VR} = \begin{bmatrix} \mathbf{H}_k^B [1] \\ \mathbf{H}_k^B [2] \\ \dots \\ \mathbf{H}_k^B [M] \end{bmatrix} \mathbf{W}_{kt}^{VR} = \begin{bmatrix} \mathbf{H}_k^B [1] \mathbf{W}_{kt}^{VR} \\ \mathbf{H}_k^B [2] \mathbf{W}_{kt}^{VR} \\ \dots \\ \mathbf{H}_k^B [M] \mathbf{W}_{kt}^{VR} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S40})$$

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

where \mathbf{H}_k^B is the input matrix of row attention layer in the k -th self-attention block (See Eq. S27), \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:

$$\mathbf{W}_{kt}^{AR} = \text{SoftMax}\left(\frac{\sum_{i=1}^M \mathbf{Q}_{kt}^R [i] \cdot (\mathbf{K}_{kt}^R [i])^T}{\sqrt{MD/m}}\right) \in R^{L \times L}, \mathbf{Q}_{kt}^R [i], \mathbf{K}_{kt}^R [i] \in R^{L \times (D/m)} \quad (\text{S41})$$

$$\mathbf{W}_{kt}^{AR} \leftarrow \text{dropout}(\mathbf{W}_{kt}^{AR}, r) \quad (\text{S42})$$

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 the Value matrix \mathbf{V}_{kt}^R

to generate the corresponding row attention matrix:

$$\mathbf{A}_{kt}^R = \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R = \mathbf{W}_{kt}^{AR} \begin{bmatrix} \mathbf{V}_{kt}^R[1] \\ \mathbf{V}_{kt}^R[2] \\ \dots \\ \mathbf{V}_{kt}^R[M] \end{bmatrix} = \begin{bmatrix} \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[1] \\ \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[2] \\ \dots \\ \mathbf{W}_{kt}^{AR} \mathbf{V}_{kt}^R[M] \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})}, \mathbf{V}_{kt}^R[i] \in R^{L \times (\frac{D}{m})} \quad (\text{S43})$$

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:

$$\mathbf{A}_k^R = \mathbf{A}_{k1}^R \mathbf{A}_{k2}^R \dots \mathbf{A}_{km}^R \in R^{M \times L \times D} \quad (\text{S44})$$

$$\mathbf{H}_k^R = \mathbf{A}_k^R \mathbf{W}_k^R + \mathbf{b}_k^R = \begin{bmatrix} \mathbf{A}_k^R[1] \\ \mathbf{A}_k^R[2] \\ \dots \\ \mathbf{A}_k^R[M] \end{bmatrix} \mathbf{W}_k^R + \mathbf{b}_k^R = \begin{bmatrix} \mathbf{A}_k^R[1] \mathbf{W}_k^R \\ \mathbf{A}_k^R[2] \mathbf{W}_k^R \\ \dots \\ \mathbf{A}_k^R[M] \mathbf{W}_k^R \end{bmatrix} + \mathbf{b}_k^R \in R^{M \times L \times D} \quad (\text{S45})$$

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

where \mathbf{H}_k^R is the output matrix of row attention layer in the k -th attention block (See Eq. S27), 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.

$$\mathbf{Q}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{QC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{QC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{QC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{QC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{QC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S46})$$

$$\mathbf{K}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{KC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{KC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{KC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{KC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{KC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S47})$$

$$\mathbf{V}_{kt}^C = \mathbf{F}_k^B \mathbf{W}_{kt}^{VC} = \begin{bmatrix} \mathbf{F}_k^B[1] \\ \mathbf{F}_k^B[2] \\ \dots \\ \mathbf{F}_k^B[M] \end{bmatrix} \mathbf{W}_{kt}^{VC} = \begin{bmatrix} \mathbf{F}_k^B[1] \mathbf{W}_{kt}^{VC} \\ \mathbf{F}_k^B[2] \mathbf{W}_{kt}^{VC} \\ \dots \\ \mathbf{F}_k^B[M] \mathbf{W}_{kt}^{VC} \end{bmatrix} \in R^{M \times L \times (\frac{D}{m})} \quad (\text{S48})$$

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

where \mathbf{F}_k^B is the input matrix of column attention layer in the k -th self-attention block (see Eq. S31), \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 matrices.

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:

$$\mathbf{W}_{kt}^{AC} = \text{SoftMax} \left(\frac{\mathbf{Q}_{kt}^C (\mathbf{K}_{kt}^C)^T}{\sqrt{D/m}} \right) \in R^{M \times L \times M} \quad (\text{S49})$$

$$\mathbf{W}_{kt}^{AC} \leftarrow \text{dropout}(\mathbf{W}_{kt}^{AC}, r) \quad (\text{S50})$$

$$\begin{aligned} \mathbf{Q}_{kt}^C (\mathbf{K}_{kt}^C)^T &= [\mathbf{Q}_{kt}^C[:, 1, :] \mathbf{Q}_{kt}^C[:, 2, :] \dots \mathbf{Q}_{kt}^C[:, L, :]] \cdot [\mathbf{K}_{kt}^C[:, 1, :] \mathbf{K}_{kt}^C[:, 2, :] \dots \mathbf{K}_{kt}^C[:, L, :]]^T = \\ &[\mathbf{Q}_{kt}^C[:, 1, :] \cdot \mathbf{K}_{kt}^C[:, 1, :]}^T \mathbf{Q}_{kt}^C[:, 2, :] \cdot \mathbf{K}_{kt}^C[:, 2, :]}^T \dots \mathbf{Q}_{kt}^C[:, L, :] \cdot \mathbf{K}_{kt}^C[:, L, :]}^T] \in R^{M \times L \times M} \end{aligned} \quad (\text{S51})$$

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

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:

$$\begin{aligned} \mathbf{A}_{kt}^C &= \mathbf{W}_{kt}^{AC} \mathbf{V}_{kt}^C = [\mathbf{W}_{kt}^{AC}[:, 1, :] \mathbf{W}_{kt}^{AC}[:, 2, :] \dots \mathbf{W}_{kt}^{AC}[:, L, :]] \cdot [\mathbf{V}_{kt}^C[:, 1, :] \mathbf{V}_{kt}^C[:, 2, :] \dots \mathbf{V}_{kt}^C[:, L, :]] = [\mathbf{W}_{kt}^{AC}[:, 1, :]} \cdot \\ &\mathbf{V}_{kt}^C[:, 1, :]} \mathbf{W}_{kt}^{AC}[:, 2, :]} \cdot \mathbf{V}_{kt}^C[:, 2, :]} \dots \mathbf{W}_{kt}^{AC}[:, L, :]} \cdot \mathbf{V}_{kt}^C[:, L, :]}] \in R^{M \times L \times (\frac{D}{m})} \end{aligned} \quad (\text{S52})$$

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

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:

$$\mathbf{A}_k^C = \mathbf{A}_{k1}^C \mathbf{A}_{k2}^C \dots \mathbf{A}_{km}^C \in R^{M \times L \times D} \quad (\text{S53})$$

$$\mathbf{F}_k^C = \mathbf{A}_k^C \mathbf{W}_k^C + \mathbf{b}_k^C = \begin{bmatrix} \mathbf{A}_k^C[1] \\ \mathbf{A}_k^C[2] \\ \dots \\ \mathbf{A}_k^C[M] \end{bmatrix} \mathbf{W}_k^C = \begin{bmatrix} \mathbf{A}_1^C[1] \mathbf{W}_k^C \\ \mathbf{A}_2^C[2] \mathbf{W}_k^C \\ \dots \\ \mathbf{A}_k^C[M] \mathbf{W}_k^C \end{bmatrix} + \mathbf{b}_k^C \in R^{M \times L \times D} \quad (\text{S54})$$

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

where \mathbf{F}_k^C is the output matrix of column attention layer in the k -th attention block, (See Eq. S31), and \mathbf{W}_k^C and \mathbf{b}_k^C are weight matrix and bias in the linear unit, respectively.

(c) Feed-forward network

$$\mathbf{T}_k^F = \text{gelu}(\mathbf{U}_k^B \mathbf{W}_k^1 + \mathbf{b}_k^1) \in R^{M \times L \times D_1} \quad (\text{S55})$$

$$\mathbf{T}_k^F \leftarrow \text{dropout}(\mathbf{T}_k^F, r) \quad (\text{S56})$$

$$\mathbf{U}_k^F = \mathbf{T}_k^F \mathbf{W}_k^2 + \mathbf{b}_k^2 \in R^{M \times L \times D} \quad (\text{S57})$$

$$\text{gelu}(x) = x\phi(x) \quad (\text{S58})$$

$$\mathbf{U}_k^B \mathbf{W}_k^1 = \begin{bmatrix} \mathbf{U}_k^B[1] \\ \mathbf{U}_k^B[2] \\ \dots \\ \mathbf{U}_k^B[M] \end{bmatrix} \mathbf{W}_k^1 = \begin{bmatrix} \mathbf{U}_k^B[1] \mathbf{W}_k^1 \\ \mathbf{U}_k^B[2] \mathbf{W}_k^1 \\ \dots \\ \mathbf{U}_k^B[M] \mathbf{W}_k^1 \end{bmatrix} \in R^{M \times L \times D_1} \quad (\text{S59})$$

$$\mathbf{T}_k^F \mathbf{W}_k^2 = \begin{bmatrix} \mathbf{T}_k^F[1] \\ \mathbf{T}_k^F[2] \\ \dots \\ \mathbf{T}_k^F[M] \end{bmatrix} \mathbf{W}_k^2 = \begin{bmatrix} \mathbf{T}_k^F[1] \mathbf{W}_k^2 \\ \mathbf{T}_k^F[2] \mathbf{W}_k^2 \\ \dots \\ \mathbf{T}_k^F[M] \mathbf{W}_k^2 \end{bmatrix} \in R^{M \times L \times D} \quad (\text{S60})$$

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

where \mathbf{U}_k^B and \mathbf{U}_k^F are the input and output matrices of the feed-forward network in the k -th self-attention block, respectively, (see Eq. S35), \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 .

F. 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:

$$\mathbf{P} = \text{SoftMax}(\mathbf{H}_{N+1} \mathbf{W}^O + \mathbf{b}^O) \in R^{M \times L \times C_{max}} \quad (\text{S61})$$

$$\mathbf{H}_{N+1} \mathbf{W}^O = \begin{bmatrix} \mathbf{H}_{N+1}[1] \mathbf{W}^O \\ \mathbf{H}_{N+1}[2] \mathbf{W}^O \\ \dots \\ \mathbf{H}_{N+1}[M] \mathbf{W}^O \end{bmatrix}, \mathbf{H}_{N+1}[i] \in R^{L \times D}, \mathbf{W}^O \in R^{D \times C_{max}} \quad (\text{S62})$$

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.

G. Loss function

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

$$\text{Loss}_{msa} = \frac{1}{M} \cdot \sum_{i=1}^M \left\{ \frac{1}{|\text{mask}(i)|} \cdot \sum_{j \in \text{mask}(i)} -\log \mathbf{P}_{i,j,c(i,j)} \right\} \quad (\text{S63})$$

where M is the number of alignments, $\text{mask}(i)$ is a set of masking positions 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 the negative log-likelihood of the true amino acid at the j -th position in the i -th sequence under the condition of masking.

Text S4. p -value calculation between EPE and the other six feature embeddings

We select the two-sided Student's t-test [10] to calculate the p -values between EPE and the other six feature embeddings. Specifically, for each feature embedding, we fed it to the designed LSTM-attention network to train a DNA-binding site prediction model, which is then evaluated on the test dataset to calculate the corresponding evaluation index, such as MCC and AUROC values. To reduce the influence of randomness, we repeat this procedure 10 times to generate a group of 10 evaluation indices. Finally, the p -value between two feature embeddings is calculated on their groups of evaluation indices under the two-sided Student's t-test. In this work, we use the Python package "scipy" to implement the Student's t-test to calculate the p -values.

Text S5. Procedures for constructing PDNA-960 and PDNA-136 datasets

Firstly, we downloaded 7345 protein-DNA complex structures which were released in the PDB database before October 15, 2023. In each complex structure, we removed the protein chains whose lengths are more than 1000 or less than 30. Then, the CD-HIT software [11] with a cut-off of 30% sequence identity was performed on all protein chains to remove the redundant chains. After this, we collected 1096 non-redundant protein chains, each of which was labeled with DNA-binding sites using the criteria of Critical Assessment of Structure Prediction (CASP) [2, 3] (see details in Text S1). Finally, the 136 chains released in the PDB after January 1, 2023, were used as the test dataset (i.e., PDNA-136, 2193 DNA-binding sites, and 47287 non-DNA-binding sites), while the remaining 960 chains were used as the training dataset (i.e., PDNA-960, 18336 DNA-binding sites, and 271988 non-DNA-binding sites).

Text S6. Formulas for calculating average precision

In the test dataset, the average precision is calculated as follows:

$$AP = \frac{1}{n} \cdot \sum_{i=1}^n \frac{TP_i}{TP_i + FP_i} \quad (\text{S64})$$

where TP_i and FP_i are the numbers of true positives and false positives, respectively, in the i -th test protein, and n is the number of test proteins.

Supporting Figures

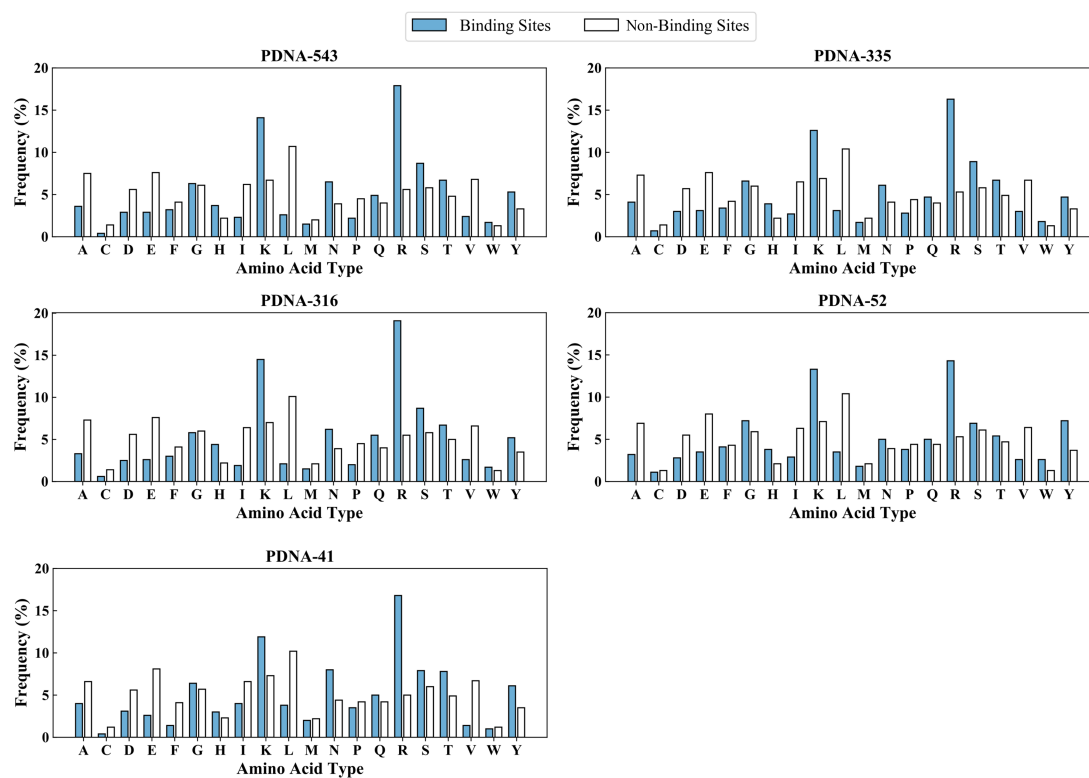


Figure S1. The frequency distributions of 20 native amino acids among DNA-binding and non-DNA-binding sites for five benchmark datasets.

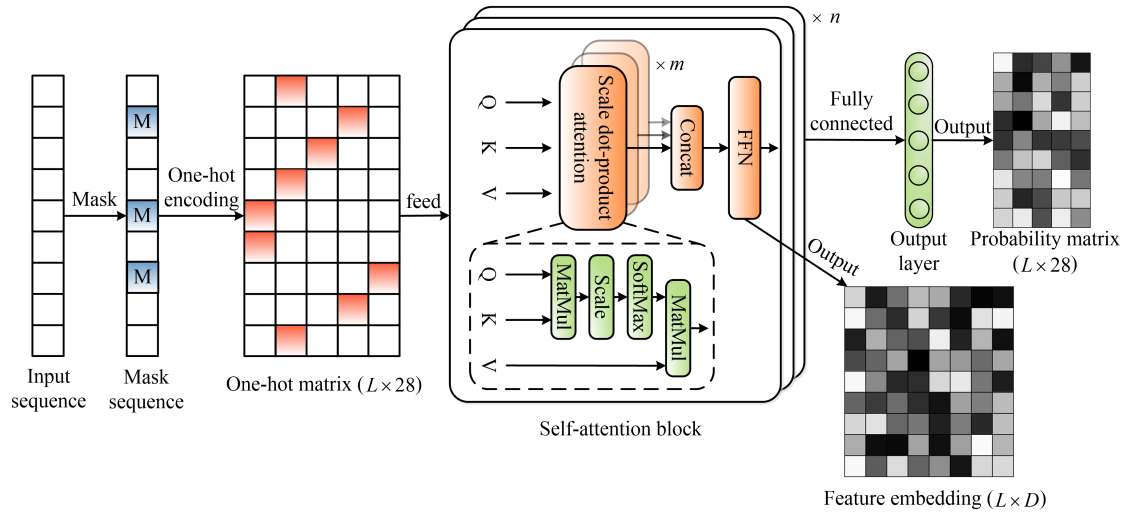
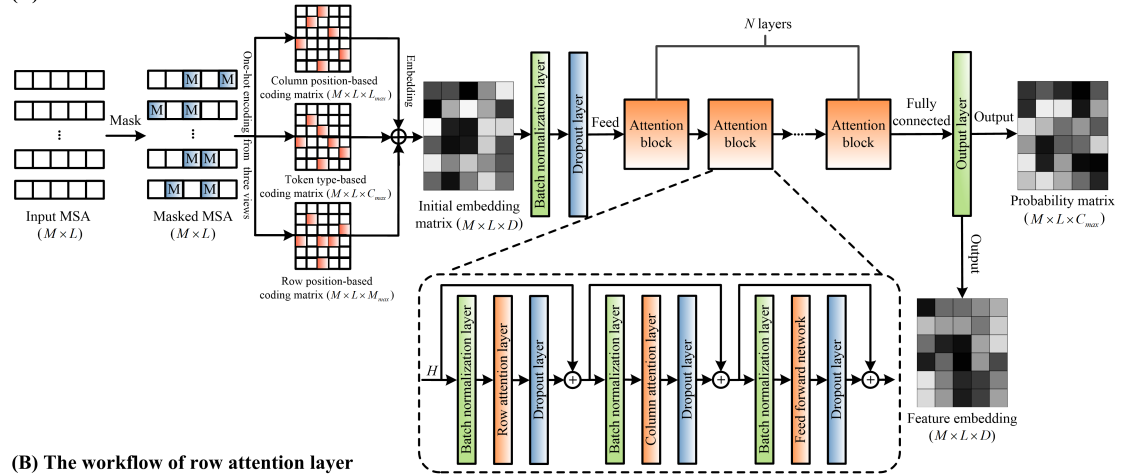
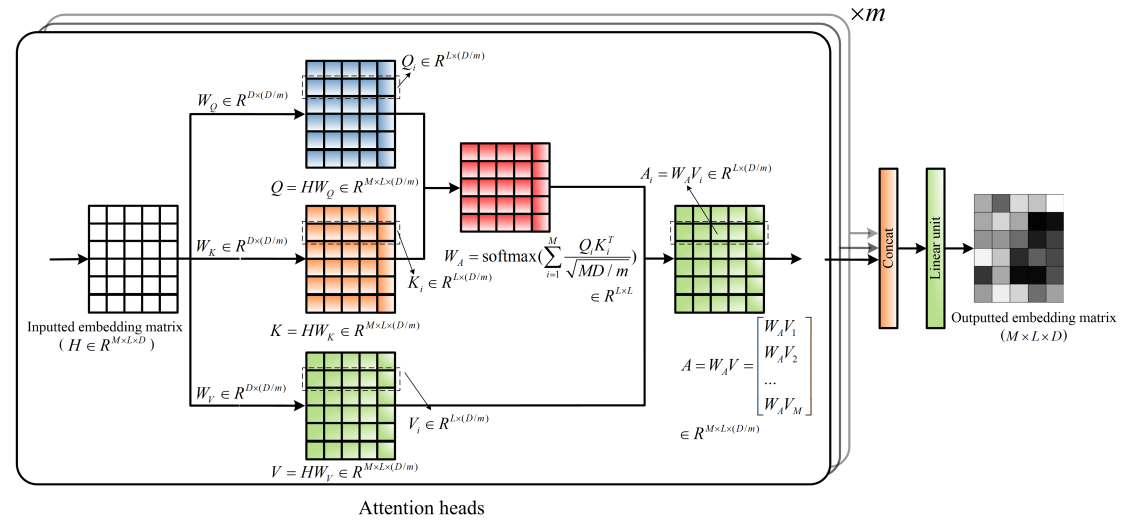


Figure S2. The workflow of the ESM2 transformer.

(A) The framework of ESM-MSA



(B) The workflow of row attention layer



(C) The workflow of column attention layer

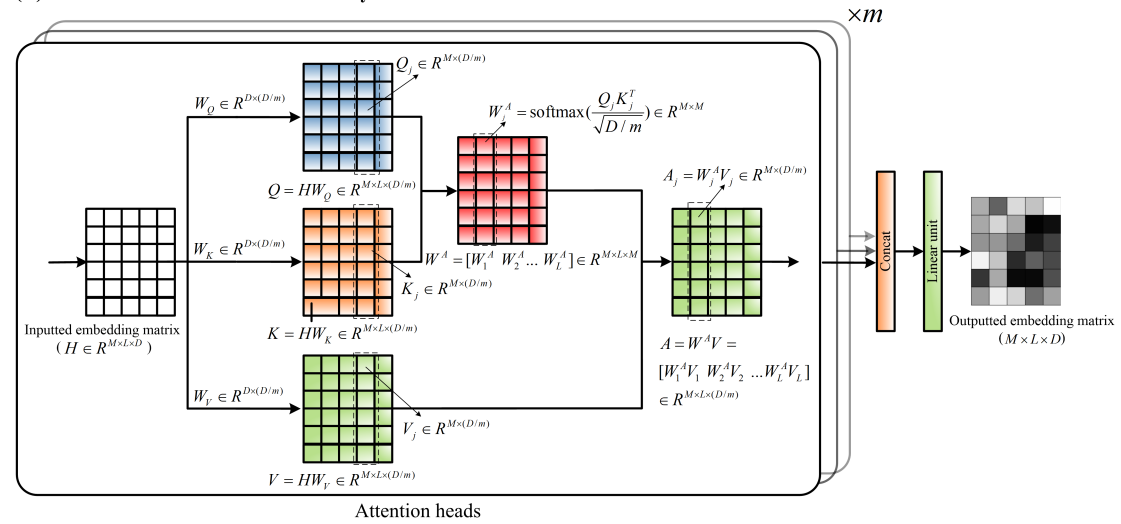


Figure S3. The workflow of the ESM-MSA transformer.

Supporting Tables

Table S1. The performance of 5 DNA-binding site predictors on the PDNA-543 dataset over ten-fold cross-validation.

Method	Sen	Spe	Acc	MCC	AUROC
TargetDNA (<i>Sen</i> \approx <i>Spe</i>) ^a	0.770	0.771	0.770	0.304	0.845
DNAPred (<i>Sen</i> \approx <i>Spe</i>) ^b	0.771	0.785	0.784	0.318	0.861
PredDBR (<i>Sen</i> \approx <i>Spe</i>) ^c	0.776	0.774	0.774	0.338	-
ULDNA (<i>Sen</i> \approx <i>Spe</i>)	0.864	0.861	0.861	0.462	0.933
TargetDNA (<i>Spe</i> \approx 0.95) ^a	0.406	0.950	0.914	0.339	0.845
DNAPred (<i>Spe</i> \approx 0.95) ^b	0.449	0.950	0.917	0.373	0.861
PredDBR (<i>Spe</i> \approx 0.95) ^c	0.465	0.950	0.911	0.409	-
Guan's method (<i>Spe</i> \approx 0.95) ^d	0.452	0.954	0.928	0.352	-
ULDNA (<i>Spe</i> \approx 0.95)	0.668	0.950	0.931	0.534	0.933

^{a, b, c, d} Results excerpted from TargetDNA [6], DNAPred [12], PredDBR [13], and Guan et al [14], respectively. "*Sen* \approx *Spe*" and "*Spe* \approx 0.95" mean that the thresholds make *Sen* \approx *Spe* and "*Spe* \approx 0.95", respectively, on the PDNA-543 dataset over ten-fold cross-validation. '-' means the value is not available. Bold fonts highlight the best performer in each evaluation index.

Table S2. The performance of 5 DNA-binding site predictors on the PDNA-335 dataset over ten-fold cross-validation.

Method	Sen	Spe	Acc	MCC	AUROC
EC-RUS ^a	0.487	0.951	0.926	0.378	0.852
TargetS ^b	0.417	0.945	0.899	0.362	0.824
DNAPred ^c	0.543	0.917	0.886	0.390	0.856
PredDBR ^d	0.426	0.953	0.910	0.390	-
ULDNA	0.676	0.948	0.925	0.565	0.940

^{a, b, c, d} Results excerpted from EC-RUS [15], TargetS [16], DNAPred [12], and PredDBR [13], respectively. '-' means the value is not available. Bold fonts highlight the best performer in each evaluation index.

Table S3. The performance of 11 DNA-binding site predictors on the PDNA-316 dataset over ten-fold cross-validation.

Method	Sen	Spe	Acc	MCC
DBS-PRED ^a	0.530	0.760	0.750	0.170
BindN ^a	0.540	0.800	0.780	0.210
DNABindR ^a	0.660	0.740	0.730	0.230
DISIS ^a	0.190	0.980	0.920	0.250
DP-Bind ^a	0.690	0.790	0.780	0.290
BindN-rf ^a	0.670	0.830	0.820	0.320
MetaDBSite ^a	0.770	0.770	0.770	0.320
TargetDNA (<i>Sen</i> \approx <i>Spe</i>) ^a	0.780	0.780	0.780	0.339
TargetDNA (<i>Spe</i> \approx 0.95) ^a	0.430	0.950	0.910	0.375
DNAPred (<i>Sen</i> \approx <i>Spe</i>) ^b	0.800	0.799	0.799	0.370
DNAPred (<i>Spe</i> \approx 0.95) ^b	0.521	0.951	0.918	0.452
PredDBR (<i>Sen</i> \approx <i>Spe</i>) ^c	0.815	0.807	0.808	0.398
PredDBR (<i>Spe</i> \approx 0.95) ^c	0.561	0.953	0.921	0.497
PredDBR (threshold = 0.5) ^c	0.531	0.958	0.923	0.489
ULDNA (<i>Sen</i> \approx <i>Spe</i>)	0.871	0.867	0.867	0.502
ULDNA (<i>Spe</i> \approx 0.95)	0.676	0.950	0.929	0.561
ULDNA (threshold = 0.5)	0.449	0.983	0.942	0.526

^{a, b, c} Results excerpted from TargetDNA [6], DNAPred [12], and PredDBR [13], respectively. “*Sen* \approx *Spe*” and “*Spe* \approx 0.95” mean that the thresholds make *Sen* \approx *Spe* and “*Spe* \approx 0.95”, respectively, on the PDNA-316 dataset over ten-fold cross-validation. Bold fonts highlight the best performer in each evaluation index.

Table S4. The *p*-values of MCC values between EPE and the other six feature embeddings on five benchmark datasets, where the base model is the designed LSTM-attention network.

Dataset	Feature embeddings from different protein language models					
	(EPE, ESM2)	(EPE, ProtTrans)	(EPE, ESM-MSA)	(EPE, PE)	(EPE, EE)	(EPE, EP)
PDNA-543	1.96×10^{-06}	2.96×10^{-12}	3.43×10^{-12}	1.01×10^{-08}	4.25×10^{-04}	9.55×10^{-04}
PDNA-335	3.89×10^{-10}	7.74×10^{-16}	1.40×10^{-11}	6.08×10^{-08}	1.25×10^{-03}	6.64×10^{-06}
PDNA-316	1.37×10^{-09}	1.13×10^{-10}	1.15×10^{-14}	1.57×10^{-07}	4.34×10^{-06}	9.42×10^{-03}
PDNA-41	8.22×10^{-08}	9.25×10^{-07}	5.21×10^{-09}	2.75×10^{-04}	1.97×10^{-03}	5.99×10^{-02}
PDNA-52	5.82×10^{-05}	6.39×10^{-04}	8.78×10^{-09}	1.12×10^{-03}	1.91×10^{-02}	3.14×10^{-02}

ESM2: the feature embedding from the ESM2 transformer; ProtTrans: the feature embedding from the ProtTrans transformer; ESM-MSA: the feature embedding from the ESM-MSA transformer; PE: the concatenation of two feature embeddings from the ProtTrans and ESM-MSA transformers; EE: the concatenation of two feature embeddings from the ESM2 and ESM-MSA transformers; EP: the concatenation of two feature embeddings from the ESM2 and ProtTrans transformers; EPE: the concatenation of three feature embeddings from the ESM2, ProtTrans, and ESM-MSA transformers.

Table S5. The p -values of AUROC values between EPE and the other six feature embeddings on five benchmark datasets, where the base model is the designed LSTM-attention network.

Dataset	Feature embeddings from different protein language models					
	(EPE, ESM2)	(EPE, ProtTrans)	(EPE, ESM-MSA)	(EPE, PE)	(EPE, EE)	(EPE, EP)
PDNA-543	1.14×10^{-08}	5.02×10^{-10}	5.84×10^{-08}	1.06×10^{-05}	3.40×10^{-02}	9.89×10^{-05}
PDNA-335	1.76×10^{-10}	8.12×10^{-13}	1.45×10^{-12}	1.38×10^{-11}	2.67×10^{-05}	4.52×10^{-07}
PDNA-316	1.77×10^{-08}	8.99×10^{-11}	1.77×10^{-12}	1.31×10^{-08}	6.42×10^{-05}	1.85×10^{-07}
PDNA-41	2.28×10^{-07}	1.37×10^{-03}	5.84×10^{-05}	7.17×10^{-02}	3.48×10^{-01}	2.37×10^{-02}
PDNA-52	3.77×10^{-03}	2.40×10^{-04}	4.47×10^{-06}	9.24×10^{-04}	6.44×10^{-03}	2.07×10^{-01}

ESM2: the feature embedding from the ESM2 transformer; ProtTrans: the feature embedding from the ProtTrans transformer; ESM-MSA: the feature embedding from the ESM-MSA transformer; PE: the concatenation of two feature embeddings from the ProtTrans and ESM-MSA transformers; EE: the concatenation of two feature embeddings from the ESM2 and ESM-MSA transformers; EP: the concatenation of two feature embeddings from the ESM2 and ProtTrans transformers; EPE: the concatenation of three feature embeddings from the ESM2, ProtTrans, and ESM-MSA transformers.

Table S6. The predicted and native DNA-binding sites of two representative proteins for five DNA-binding prediction methods.

Protein	Method	Predicted DNA-binding sites
2MXF_A	LA-ESM2	2R 5K 7Y 10P 11H 18T 19K 20G 21G 22N 23H 24K 27K 30K
	LA-ProtTrans	1A 2R 3K 4V 5K 16I 17E 18T 19K 20G 21G 22N 23H 24K 25T 27K
	LA-ESM-MSA	2R 5K 7Y 18T 19K 20G 21G 22N 23H 24K 27K 30K 41W
	ULDNA	2R 3K 5K 7Y 18T 19K 20G 21G 22N 23H 24K 27K 30K
	PredDBR	2R 5K 7Y 8K 9N 18T 19K 20G 21G 23H
	Native DNA-binding sites	1A 2R 3K 5K 7Y 18T 19K 20G 21G 22N 23H 24K 26L 27K 30K 39E
3ZQL_A	LA-ESM2	12R 13R 14S 15A 16R 17S 18H 19R 20T 43S 44M 45R 54G 55T 56M 57S 59Y 60Y 61Y 180R 183 M
	LA-ProtTrans	13R 14S 15A 16R 17S 18H 19R 20T 21L 43S 44M 45R 55T 56M 57S 59Y 60Y 61Y 65K
	LA-ESM-MSA	12R 13R 14S 15A 16R 17S 18H 19R 20T 23R 43S 44M 45R 46R 53A 54G 55T 56M 57S 59Y 60Y 61Y 64T 65K
	ULDNA	12R 13R 14S 15A 16R 17S 18H 19R 20T 43S 44M 45R 53A 54G 55T 56M 57S 59Y 60Y 61Y 64T 65K
	PredDBR	1V 4W 6H 7P 12R 15A 18H 19R 22S 23R 43S 44M 45R 53A 54G 55T 56M 57S 59Y 60Y 61Y 64T 65K 115W 117N 119H 124P 125N 126S 182W 187G 236D
	Native DNA-binding sites	12R 19R 43S 44M 45R 54G 55T 56M 57S 59Y 60Y 61Y 64T 65K

Bold font means a DNA-binding site can be correctly predicted by a DNA-binding prediction method.

Reference

- [1] J. Yang, A. Roy, and Y. Zhang, "BioLiP: a semi-manually curated database for biologically relevant ligand–protein interactions," *Nucleic acids research*, vol. 41, no. D1, pp. D1096-D1103, 2012.
- [2] T. Gallo Cassarino, L. Bordoli, and T. Schwede, "Assessment of ligand binding site predictions in CASP10," *Proteins: Structure, Function, and Bioinformatics*, vol. 82, pp. 154-163, 2014.
- [3] T. Schmidt, J. Haas, T. G. Cassarino, and T. Schwede, "Assessment of ligand-binding residue predictions in CASP9," *Proteins: Structure, Function, and Bioinformatics*, vol. 79, no. S10, pp. 126-136, 2011.
- [4] C. Zhang, X. Zhang, P. L. Freddolino, and Y. Zhang, "BioLiP2: an updated structure database for biologically relevant ligand–protein interactions," *Nucleic Acids Research*, p. gkad630, 2023.
- [5] S. Ahmad, M. M. Gromiha, and A. Sarai, "Analysis and prediction of DNA-binding proteins and their binding residues based on composition, sequence and structural information," *Bioinformatics*, vol. 20, no. 4, pp. 477-486, 2004.
- [6] J. Hu, Y. Li, M. Zhang, X. Yang, H.-B. Shen, and D.-J. Yu, "Predicting protein-DNA binding residues by weightedly combining sequence-based features and boosting multiple SVMs," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 14, no. 6, pp. 1389-1398, 2016.
- [7] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [8] J. Buckman, A. Roy, C. Raffel, and I. Goodfellow, "Thermometer encoding: One hot way to resist adversarial examples," in *International Conference on Learning Representations*, 2018.
- [9] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [10] J. C. De Winter, "Using the Student's t-test with extremely small sample sizes," *Practical Assessment, Research, and Evaluation*, vol. 18, no. 1, p. 10, 2019.
- [11] W. Li and A. Godzik, "CD-HIT: a fast program for clustering and comparing large sets of protein or nucleotide sequences," *Bioinformatics*, vol. 22, no. 13, pp. 1658-1659, 2006.
- [12] Y.-H. Zhu, J. Hu, X.-N. Song, and D.-J. Yu, "DNAPred: accurate identification of DNA-binding sites from protein sequence by ensembled hyperplane-distance-based support vector machines," *Journal of Chemical Information and Modeling*, vol. 59, no. 6, pp. 3057-3071, 2019.
- [13] J. Hu, Y.-S. Bai, L.-L. Zheng, N.-X. Jia, D.-J. Yu, and G.-J. Zhang, "Protein-DNA binding residue prediction via bagging strategy and sequence-based cube-format feature," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 19, no. 6, pp. 3635-3645, 2021.
- [14] S. Guan, Q. Zou, H. Wu, and Y. Ding, "Protein-DNA binding residues prediction using a deep learning model with hierarchical feature extraction," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, p. DOI:

10.1109/TCBB.2022.3190933, 2022.

- [15] Y. Ding, J. Tang, and F. Guo, "Identification of protein-ligand binding sites by sequence information and ensemble classifier," *Journal of Chemical Information and Modeling*, vol. 57, no. 12, pp. 3149-3161, 2017.
- [16] D.-J. Yu, J. Hu, J. Yang, H.-B. Shen, J. Tang, and J.-Y. Yang, "Designing template-free predictor for targeting protein-ligand binding sites with classifier ensemble and spatial clustering," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 10, no. 4, pp. 994-1008, 2013.