

The supplement of scFTAT

Binhua Tang and Yiyao Chen

Contact: bh.tang@hhu.edu.cn

1. The main structure of scFTAT:

1.1 The improved Transformer layer:

After initial training in the frequency domain through the FFT encoding layer, the output gene expression matrix is fed into the improved Transformer layer. A complete Transformer layer generally consists of a multi-head attention mechanism module, a feedforward function module, and the necessary residual layers. This study has improved these three modules to achieve a more efficient and faster process for cell type identification.

Let X be the input gene expression data. Absolute positional encoding is performed before entering the multi-head attention layer to preserve the sequential information of the input gene expression data. Specifically, sine and cosine functions are used to represent the positions of the data. In the Transformer model, the general representation is as follows:

$$PE(p, 2i) = \sin(p / 1000^{2i/d_m}) \quad (S1)$$

$$PE(p, 2i + 1) = \cos(p / 1000^{2i/d_m}) \quad (S2)$$

where p represents the specific position of the input gene x in the gene expression matrix X , i represents the gene dimension, and d_m represents the dimension of the embedded positional encoding vector. Thus, the original input data X added with the result of absolute positional encoding, $PE(X)$, serves as the input X_{PE} for the subsequent multi-head attention layer.

The multi-head attention mechanism module is essentially composed of multiple parallel self-attention modules. The method for calculating the self-attention module is represented as below,

$$Attention(Q, K, V) = SoftMax\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (S3)$$

where $Q = W^q X_{PE}$, $K = W^k X_{PE}$, $V = W^v X_{PE}$, $W^{q,k,v}$, Q , K , and V represent the query vector, key vector, and value vector, respectively, and X_{PE} represents the output vector after positional encoding.

Equation (S3) calculates attention scores by directly computing the Q and K dot product, then normalizing it using SoftMax to obtain the final output. This method directly incorporates the absolute positional encoding information into the context representation. However, relative positional encoding has better generalization ability and scalability than absolute positional encoding, which is more advantageous when dealing with long sequences. Incorporating the relative information between input genes into model training for single-cell input data can improve the model's performance.

Some Transformer-based Computer Vision (CV) models have already explored relative positional encoding [4]. Thus, in this study, calculating attention scores is modified to implement relative positional encoding. In simple terms, after calculating Q and K , they are multiplied with the rotation encoding matrix. In any even dimension, the rotated positional encoding matrix R and the representation of the new Q are shown below:

$$R_{\theta,m}^d = \begin{pmatrix} \cos m\theta_1 & -\sin m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ \sin m\theta_1 & \cos m\theta_1 & 0 & 0 & \cdots & 0 & 0 \\ 0 & 0 & \cos m\theta_2 & -\sin m\theta_2 & \cdots & 0 & 0 \\ 0 & 0 & \sin m\theta_2 & \cos m\theta_2 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & \cdots & \cos m\theta_{d/2} & -\sin m\theta_{d/2} \\ 0 & 0 & 0 & 0 & \cdots & \sin m\theta_{d/2} & \cos m\theta_{d/2} \end{pmatrix} \quad (S4)$$

$$Q = R_{\theta,m}^d \begin{pmatrix} q^0 \\ q^1 \\ q^2 \\ \vdots \\ q^{d-1} \end{pmatrix} \quad (S5)$$

where d refers to the spatial dimension, m the position of the query vector Q , and q the

specific elements of Q in each dimension. Therefore, θ can be estimated on d , and its

vector is denoted as $\Theta = \{\theta_i = 10000^{\frac{-2(i-1)}{d}}, i \in [1, 2, \dots, \frac{d}{2}]\}$.

Because Equation (S4) is an orthogonal matrix, it does not change the vector magnitude during operation, thus further ensuring the model stability. In practical implementation, as $R_{\theta,m}^d$ is sparse, it can be estimated as follows to reduce the time complexity,

$$Q = R_{\theta,m}^d X = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \cos m\theta_1 \\ \cos m\theta_1 \\ \cos m\theta_2 \\ \vdots \\ \cos m\theta_{d/2} \end{pmatrix} + \begin{pmatrix} -x_2 \\ x_1 \\ -x_4 \\ \vdots \\ x_d \end{pmatrix} \otimes \begin{pmatrix} \sin m\theta_1 \\ \sin m\theta_1 \\ \sin m\theta_2 \\ \vdots \\ \sin m\theta_{d/2} \end{pmatrix} \quad (S6)$$

Thus, the newly updated Q and K will contain relative positional information, achieved by absolute positional encoding. Attention scores can be obtained after obtaining the new Q and K .

Here, the softmax operation is replaced with the kernel function approximation [5], depicted as,

$$\Phi_1 \Phi_2 \approx \text{SoftMax} \left(\frac{QK^T}{\sqrt{d}} \right) \quad (S7)$$

where Φ_1 and Φ_2 represent the updated Q and K . The specific form of Φ is denoted as,

$$\Phi = \frac{p}{\sqrt{m}} \exp(W^T x - \frac{\|x\|^2}{2}) \quad (S8)$$

where p is a positive constant, W the product of the original input data matrix and a random orthogonal matrix, m the dimension of matrix W , and x the input Q or K . The random orthogonal matrix here can reduce the dimensionality of the original input while retaining the corresponding features.

Next, we use the multi-head self-attention mechanism, with each head named as $head_h$, to calculate their respective attention weights in parallel, as shown in the

following equation. The specific number of heads is chosen based on the dataset's size and number of categories.

$$MulitHead(Q, K, V) = Concat(head_1, \dots, head_h) \quad (S9)$$

where *Concat* refers to the component-wise sum, $head_h = Attention(Q_h, K_h, V_h)$. After combining into a multi-head attention module, we adopt the Rezero method [6] to rescale the self-attention block. Specifically, the residual connection is represented in the following form,

$$X_i' = X_i + \alpha_i sublayer(X_i) \quad (S10)$$

where X_i is the input of the attention module, X_i' is the module's output, and α_i is a learnable residual weight shared by each multi-head attention module. This parameter is initialized to 0. This ensures that in the early training, the gradients of all parameters in the sublayer function in Equation (S10) will disappear and reach a suitable value during the training process, further speeding up the network's convergence rate.

After processing the attention module and the corresponding residual layer, a feedforward network is added to obtain nonlinear data features. The work employs a parallel feedforward function module based on global and local information enhancement [6]. In global information, the input global representation is obtained through average pooling, followed by a fully connected operation. In local information, features are directly extracted through a fully connected operation. The approach is relevant to the concept of channel attention [7]. Subsequently, the outputs of the two branches are interacted via concatenation. Finally, the two pieces of information are weighted by a gating unit to obtain the corresponding attention weights, and the output dimension is aligned with the original input dimension. This module can enhance the model's final expression ability without significantly increasing computational complexity.

2. The ablation experiments on the selected six scRNA-seq datasets

Table S1. The ablation experiment on Human_fetal_pancreasr

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.85	0.60	0.62	0.59	0.82
Transformer	0.88	0.67	0.64	0.67	0.85
Transformer+FFT	0.91	0.69	0.71	0.74	0.89
scFTAT	0.93	0.81	0.86	0.79	0.91

Table S2. The ablation experiment on Human_bladder

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.79	0.58	0.63	0.59	0.64
Transformer	0.88	0.62	0.72	0.57	0.78
Transformer+FFT	0.88	0.65	0.71	0.63	0.78
scFTAT	0.89	0.84	0.87	0.81	0.81

Table S3. The ablation experiment on Human_kidney

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.89	0.70	0.79	0.72	0.87
Transformer	0.92	0.79	0.82	0.77	0.89
Transformer+FFT	0.92	0.82	0.92	0.79	0.89
scFTAT	0.93	0.84	0.96	0.80	0.89

Table S4. The ablation experiment on Mouse_bladder

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.81	0.65	0.66	0.69	0.81
Transformer	0.81	0.70	0.71	0.74	0.77
Transformer+FFT	0.82	0.76	0.76	0.78	0.78
scFTAT	0.88	0.81	0.79	0.84	0.86

Table S5. The ablation experiment on Mouse_spleen

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.79	0.56	0.59	0.55	0.66
Transformer	0.79	0.60	0.67	0.61	0.66
Transformer+FFT	0.80	0.74	0.72	0.76	0.68
scFTAT	0.85	0.73	0.75	0.77	0.79

Table S6. The ablation experiment on Mouse_kidney

	ACC	F1	Precision	Recall	MCC
Transformer-LDA	0.75	0.34	0.38	0.34	0.61
Transformer	0.80	0.34	0.40	0.32	0.62
Transformer+FFT	0.88	0.66	0.71	0.64	0.78
scFTAT	0.90	0.90	0.91	0.91	0.89

References

1. Rao Y, Zhao W, Zhu Z et al. GFNet: Global Filter Networks for Visual Recognition, IEEE Trans Pattern Anal Mach Intell 2023;45:10960-10973.
2. Bachlechner TC, Majumder BP, Mao HH et al. ReZero is All You Need: Fast Convergence at Large Depth. In: Conference on Uncertainty in Artificial Intelligence. 2020.
3. Guo J, Gao J, Hu Y et al. Robust Adaptive Linear Discriminant Analysis with Bidirectional Reconstruction Constraint, ACM Transactions on Knowledge Discovery from Data 2020;14.
4. Dai Z, Yang Z, Yang Y et al. Transformer-XL: Attentive Language Models Beyond a Fixed-Length Context. 2019, arXiv:1901.02860.
5. Katharopoulos A, Vyas A, Pappas N, Fleuret F. Transformers are RNNs: Fast Autoregressive Transformers with Linear Attention. In: Hal D., III, Aarti S. eds). Proceedings of the 37th International Conference on Machine Learning. Proceedings of Machine Learning Research: PMLR, 2020, 5156--5165.
6. Huang T, Huang L, You S et al. LightViT: Towards Light-Weight Convolution-Free Vision Transformers. 2022, arXiv:2207.05557.
7. Hu J, Shen L, Albanie S et al. Squeeze-and-Excitation Networks, IEEE Trans Pattern Anal Mach Intell 2020;42:2011-2023.