

M2oE: Multimodal Collaborative Expert Peptide Model

Zengzhu Guo¹

School of Information Sciences
Guangdong University of Finance and Economics
Guangzhou, China
Email: gzz3383@163.com

Zhiqi Ma¹

School of Medicine
The Chinese University of Hong Kong
ShenZhen (CUHK-ShenZhen)
Email: zhiqima@link.cuhk.edu.cn

Abstract—Peptides are biomolecules comprised of amino acids that play an important role in our body. In recent years, peptides have received extensive attention in drug design and synthesis, and peptide prediction tasks help us better search for functional peptides. Typically, we use the primary sequence and structural information of peptides for model encoding. However, recent studies have focused more on single-modal information (structure or sequence) for prediction without multi-modal approaches. We found that single-modal models are not good at handling datasets with less information in that particular modality. Therefore, this paper proposes the M2oE multi-modal collaborative expert peptide model. Based on previous work, by integrating sequence and spatial structural information, employing expert model and Cross-Attention Mechanism, the model’s capabilities are balanced and improved. Experimental results indicate that the M2oE model performs excellently in complex task predictions. Code is available at: <https://github.com/goldzzmj/M2oE>

Index Terms—Antimicrobial peptides (AMP), MoE, Multi-modal

I. INTRODUCTION

Peptides, which are composed of amino acids, play pivotal roles in the modulation of physiological processes within the body. In contrast to proteins, peptides consist of shorter chains of amino acids[1]. The prediction of peptide properties entails forecasting their physicochemical characteristics, functions, and biological activities through advanced computational methods that have significantly evolved with the advent of deep learning techniques[2][3][4]. Recently, there has been a growing interest in peptides for drug design applications, particularly in the development of antimicrobial and anti-cancer agents due to the increasing prevalence of antibiotic resistance[5][6][7].

Typically, peptide encoding encompasses both the primary amino acid sequence and its spatial structure. Previous models, including RNN[8], LSTM[9], BiLSTM[10], and Transformer[11], indicate that the Transformer architecture is particularly effective in this context. Additionally, peptides can be represented as graph structures, rendering Graph Neural Networks (GNNs) instrumental for capturing molecular spatial information[12]. However, most studies predominantly focus on single-modality data, either sequence or structure, and even contrastive learning techniques often lack a genuine integration of these modalities[13].

Multimodal models have achieved significant advancements, especially within the AI4Science domain. For instance, GIT-Former[14] integrates graphical, imaging, and textual information to enhance prediction accuracy in molecular science; meanwhile, Mixture of Experts (MoE) models such as GMoE[15] and SwitchTransformer[16] optimize token allocation to improve adaptability. Despite these advancements, multimodal fusion continues to encounter challenges—particularly regarding the refinement of fusion methods for enhanced integration.

To enhance model performance, our M2oE model builds upon previous research by employing a mixed expert framework for embedding, which integrates multiple expert models to achieve more accurate task predictions[17][18][19]. This paper presents a multimodal collaborative expert peptide model with the following key contributions:

1. We propose a sequence-structure mixing expert model that addresses the challenge of expert allocation [20].
2. We leverage multimodal characteristics to improve mixed expert representation through interactive attention networks.
3. We utilize learnable weights α to evaluate the significance of sequence and spatial information across various data distribution scenarios.

II. METHODS

A. Benchmark dataset

The benchmark dataset utilized in this study is derived from Liu et al. [21]. According to the task division, the datasets encompass classification and regression tasks, which include antimicrobial peptides (AMP) [22] and aggregation propensity (AP) [13]. The processing of these two types of datasets aligns with previous work [21], having been partitioned into training, validation, and test sets at a ratio of 8:1:1. More detailed information is provided in Table I.

B. Sequence and Graph encodings

Peptide sequences $S \subseteq \mathbf{R}^M$ and sentence data are similar in that both require word-base embedding and positional identification combination as input. However, the difference is that the division of peptide sequences is based on amino acids and does not require complex tokenizer like natural language. Multi-head Self Attention (MSA) is the core in the Transformer

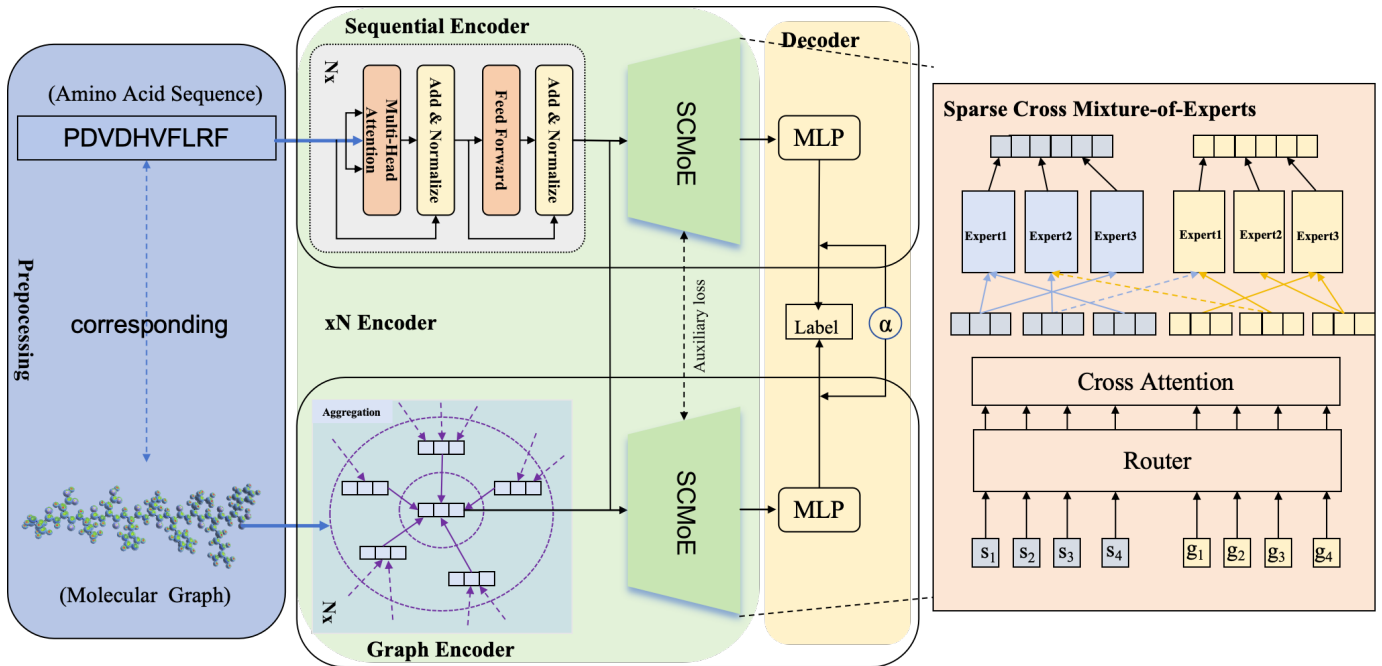


Fig. 1. The framework of the M2oE. The model is structured with an encoding module and a decoding module, incorporating the interactive attention mechanism in the SCMoE module and MoE token allocation to enhance the comprehensive ability of the M2oE encoding model. Additionally, MoE optimization is achieved through auxiliary loss. The decoding module utilizes MLP and learnable parameters α from both modes for making predictions.

TABLE I

THE CLASSIFICATION DATASET AND THE ASSOCIATION DATA SET WERE ANALYZED, WITH LABEL 1 IN AMPs REPRESENTING ANTIMICROBIAL PEPTIDES AND 0 REPRESENTING NON-ANTIMICROBIAL PEPTIDES.

Dataset	Property	Classification	Regression
		AMPs	AP
Train	AMP	5437	54159
	non-AMP	2019	
Validation	AMP	679	4000
	non-AMP	252	
Test	AMP	681	4000
	non-AMP	253	
Total		9321	62159

which scores the context and captures various dependencies within the sequence. Feed Forward (FFN) combines with non-linear activate function and additional trainable parameters, further capture non-linear relationships between amino acids and mapped to higher dimension. The sequence encoder output amino acids feature is represented as $s \in S^{M \times d}$, where d is feature hidden dimension.

The peptide molecule is defined as $\mathcal{G} = (\nu, \varepsilon)$, where $\nu = \{\nu_i\}_{i=1}^N$ represents the beads as nodes and $\varepsilon \subseteq \nu \times \nu$ represents the existence of chemical bonds between the beads as edges. Adjacent matrix $A \in \{0, 1\}^{N \times N}$ describes the relationship between nodes and is filled with 0 or 1 based on their corresponding edges, $A_{ij} = 1$, when it is existing connection $(i, j) \in \varepsilon$, otherwise $A_{ij} = 0$. Given feature adjacent X and join with adjacent A , GCN[12] leverages relative edges and nodes attribute to learn latent representation

of the node. One layer graph convolutional encoder represents as follow:

$$X^{(l+1)} = f_{GCN}(A, X^{(l)}; W^{(l)}) = \sigma(\tilde{A}X^{(l)}W^{(l)}), \quad (1)$$

where f_{GCN} is GCN encoder function, $\tilde{A} = A + I$ add diagonal matrix to keep and transmit the information of the node itself, $\hat{A} = D^{-\frac{1}{2}}\tilde{A}D^{-\frac{1}{2}}$ is to normalize the adjacency matrix. $W^{(l)}$ represent the learnable weight matrix of the l -layer of the model and σ is a non-linear activate function LeakyRelu. The Initial values $X^{(0)}$ are randomly initialized using a normal distribution and the final output by GCN is represented as $X \in \mathbb{R}^{N \times D}$ where D donates each node embedding dimension.

C. Sparse Cross Mixture of Experts

As shown in Figure 1, the parallel Transformer and SAGE-Graph capture the primary peptide sequence information and the secondary molecular structure information. However, sequence information and structural information can represent and complement each other. Therefore, we have designed a sparse interaction mixed expert system (SCMoE) fusion module.

The SCMoE model contains C sequence mixing experts and graph mixing experts, which can learn from tokens routed by different types of data through the expert network. In particular, the interactive attention network possesses the ability to focus on different modalities directly, endowing the mixing experts with stronger representational capabilities through this

multimodal alignment approach. Specifically, the routing network is controlled by a learnable matrix $W^r \in \mathbb{R}^{d \times C}$, which calculates the similarity between each token and the mixing experts, and assigns them to the *topk* most similar experts. The formula 2 shows this allocation method, where X_{ij} represents assigning the *i*-th token to the *j*-th expert with a coefficient of α .

However, using the Top *k* allocation method alone may result in some tokens never being assigned to experts, thus reducing the expressive power of the expert system[23]. To address this issue, a stochastic variable sampled from the standard normal distribution is added, allowing tokens ranked after *K* to also have a chance for allocation.

$$Router(X_i) = Topk(\alpha_j X_{ij} + N(0, 1) \cdot Softplus(X_{ij} W_{noise}^r)),$$

$$\alpha_j = \frac{X_{ij} W^r}{\sum_{j=0}^{topk} X_{ij} W^r} \quad (2)$$

Among them, $W_{noise} \in \mathbb{R}^{d \times C}$ are learnable parameters and $Softplus(\cdot)$ is a nonlinear activation function can prevent the problem of vanishing gradients.

The peptide sequence is composed of multiple amino acid symbols, so each character can be used as a local feature. The combination of local features assigned to the mixed experts implicitly expresses certain characteristics of the peptide sequence. However, relying solely on single-modal information makes it difficult to directly learn the implicit characteristics of peptides. Therefore, the Cross-Attention (CRA) is proposed to improve the MoE[14]. It can align similar characteristics between modalities while also drawing away different characteristics. Specifically, it can be represented as follows:

where F_{seq}, F_{gra} denote features from the sequence encoder and graph encoder, and d_k is the scaling factor respectively. Subsequently, we exchange the queries Q of the two modalities for spatial interaction:

$$F_{fgra} = \text{Softmax} \left(\frac{Q_{seq} K_{gra}^\top}{d_k} \right) V_{gra},$$

$$F_{fseq} = \text{Softmax} \left(\frac{Q_{gra} K_{seq}^\top}{d_k} \right) V_{seq}, \quad (3)$$

Subsequently, the cross-attention matrices of the two modalities are transformed and updated. The new sequence features are composed of graph node features and their corresponding attention coefficients. The updated interactive features also need to be allocated to different experts, similar to the formula2. Therefore, the updated sequence features can be integrated into the routing network through the operation $F_{seq}^{new} = \text{Concat}(F_{seq}, F_{fseq})$, as do the graph node features.

D. Fusion Module And Loss

The antimicrobial peptide prediction task is conducted based on the sequence and its spatial structure. Our designed SCMoE module ensures the expression of characteristics of each modality and enhances the expression of potential features

with the help of information from another modality. Therefore, the final fusion module only needs to utilize the nonlinear capability of MLP to capture the correlation between features and map them to the classification space of antimicrobial peptides. Traditional methods often involve combining multiple output results using fixed weights, but this approach is limited in that it is difficult to assess the importance of sequence and spatial information for prediction under different data distribution scenarios. As shown in formula 4, we employ learnable weights α to measure this importance.

$$\hat{y} = \sigma(\alpha MLP_1(Z_{seq}) + (1 - \alpha) MLP_2(Z_{gra})) \quad (4)$$

Among them, σ is Sigmoid function, mapping predictive data into the probability space. Z_{seq}, Z_{gra} are embeddings of the output from the sequence encoder and the graph encoder.

The routing network assigns tokens to experts based on the gating method, but this approach can sometimes lead to load imbalance issues, where one expert receives the majority of tokens, thereby degenerating into a single-expert model. Therefore, a strategy designed to ensure that each expert has an equal probability of being selected is formulated as shown in Equation 5. On the other hand, the capabilities of each expert are different, and the routing network tends to allocate tokens to the few experts with stronger capabilities, leaving the remaining experts idle, which similarly leads to load imbalance issues. As shown in Equation 6, the $CV(\cdot)$ function measures the degree of discreteness of expert importance, combined with fixed hyperparameter ω_{imp} to control the similar abilities of different experts.

$$L_{load} = \sum_{i=1}^C \left(\frac{n_i}{\sum_{j=1}^C n_j} - \frac{1}{C} \right)^2 \quad (5)$$

$$L_{importance} = \omega_{imp} \cdot CV(\sum_{x \in X} Router(x))$$

$$CV(X) = \frac{\sigma_x}{\mu_x} \quad (6)$$

Among them σ_x and μ_x are the variance and mean of data X .

Finally, the error between the predicted values and the true labels is calculated using the Binary Cross Entropy (BCE), and this is added to the balanced loss function of Mode of Expertise (MoE) regarding load and importance as the total optimization objective.

$$L = BCE(y, \hat{y}) + L_{Load} + L_{importance} \quad (7)$$

III. RESULTS AND DISCUSSION

We propose the M2oE model, which effectively balances and integrates sequence and structural features for downstream tasks. This model encompasses three types: sequence, graph, and hybrid models. The sequence and graph models are single-modality frameworks evaluated on classification (AP) and regression (AMP) datasets. For the sequence model, we employ Transformer and SwitchTransformer architectures, while the graph model utilizes GCN, GAT, GraphSAGE, and GMoE [16], [15]. The hybrid model incorporates Repcon, weighted

fusion M2oE methods, concatenation techniques, as well as our final approach.

Table III illustrates that the sequence model demonstrates superior performance on the AP dataset; notably, SwitchTransformer achieves an impressive R^2 of 95.1%. Conversely, on the AMP dataset, GraphSAGE leads with an accuracy of 84.7%. These findings suggest that single-modality models excel when a dataset is biased towards one modality but encounter challenges when it favors another.

The M2oE model synergistically combines the strengths of both sequence and graph models to achieve remarkable performance across both datasets: $R^2 = 0.951$ on AP with minimal MAE and MSE values (3.68E-2 and 2.21E-3), alongside an accuracy of 86.2% on AMP—surpassing baseline results.

While MoE enhances performance in single modalities independently without benefiting other modalities directly; therefore we implement Cross-Attention to ensure balanced improvements across modalities. Ablation experiments presented in Table II corroborate this assertion. Ultimately demonstrating that M2oE achieves optimal results by improving upon baseline metrics by 0.9%.

TABLE II
ABLATION EXPERIMENT RESULT ON THE AP DATASETS.

Variants	MAE	MSE	R^2
M2oE without CRA nor MoE	3.96E-2	2.57E-3	0.942
M2oE without CRA	3.74E-2	2.27E-3	0.949
M2oE without MoE	3.79E-2	2.38E-3	0.947
M2oE	3.68E-2	2.21E-3	0.951

TABLE III
COMPARISON WITH STATE-OF-THE-ART METHODS, INCLUDING SEQUENCE MODELS, GRAPH MODELS, AND MIXED MODELS.

Type	Model	AP			AMP
		MAE	MSE	R^2	ACC
Sequence	Transformer	3.81E-2	2.33E-3	0.947	0.813
	SwitchTransformer[16]	3.65E-2	2.15E-3	0.951	0.808
Graph	GCN	4.27E-2	3.02E-3	0.932	0.834
	GAT	4.40E-2	3.22E-3	0.928	0.843
	GraphSAGE	3.84E-2	2.36E-3	0.947	0.847
	GMoE [15]	3.82E-2	2.35E-3	0.947	0.837
Mixture	Repcon(Avg)	3.83E-2	2.24E-3	0.947	0.831
	M2oE(WS)	3.74E-2	2.29E-3	0.949	0.820
	M2oE(Concat)	3.73E-2	2.26E-3	0.949	0.824
	M2oE(Parallel)	3.68E-2	2.21E-3	0.951	0.862

IV. CONCLUSION

In this paper, we propose a multimodal collaborative expert peptide model, which integrates sequence and spatial structural information, utilizes a sparse mixed expert model, and takes into account the characteristics under different data distributions. Experimental results show that the M2oE model performs well in complex task prediction, and uses multimodal methods to solve problems that may arise in unimodal scenarios. Finally, we use ablation experiments to demonstrate the effectiveness of each module. In future work, we can also consider connecting the multimodal expert model to more complex tasks, such as peptide generation tasks, etc.

REFERENCES

- [1] I. W. Hamley, *Introduction to peptide science*. John Wiley & Sons, 2020.
- [2] G. Serrano, E. Guruceaga, and V. Segura, "Deepmspeptide: peptide detectability prediction using deep learning," *Bioinformatics*, vol. 36, no. 4, pp. 1279–1280, 2020.
- [3] X. Chen, C. Li, M. T. Bernards, Y. Shi, Q. Shao, and Y. He, "Sequence-based peptide identification, generation, and property prediction with deep learning: a review," *Molecular Systems Design & Engineering*, vol. 6, no. 6, pp. 406–428, 2021.
- [4] Y. Jiang, R. Wang, J. Feng, J. Jin, S. Liang, Z. Li, Y. Yu, A. Ma, R. Su, Q. Zou *et al.*, "Explainable deep hypergraph learning modeling the peptide secondary structure prediction," *Advanced Science*, vol. 10, no. 11, p. 2206151, 2023.
- [5] C. Wang, S. Garlick, and M. Zloh, "Deep learning for novel antimicrobial peptide design," *Biomolecules*, vol. 11, no. 3, p. 471, 2021.
- [6] C. Li, R. L. Warren, and I. Birol, "Models and data of amplify: a deep learning tool for antimicrobial peptide prediction," *BMC Research Notes*, vol. 16, no. 1, p. 11, 2023.
- [7] J. Chen, H. H. Cheong, and S. W. Siu, "xdeep-acpep: deep learning method for anticancer peptide activity prediction based on convolutional neural network and multitask learning," *Journal of chemical information and modeling*, vol. 61, no. 8, pp. 3789–3803, 2021.
- [8] J. L. Elman, "Finding structure in time," *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.
- [9] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [10] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE transactions on Signal Processing*, vol. 45, no. 11, pp. 2673–2681, 1997.
- [11] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [12] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.
- [13] Z. Liu, J. Wang, Y. Luo, S. Zhao, W. Li, and S. Z. Li, "Efficient prediction of peptide self-assembly through sequential and graphical encoding," *Briefings in Bioinformatics*, vol. 24, no. 6, p. bbad409, 2023.
- [14] P. Liu, Y. Ren, J. Tao, and Z. Ren, "Git-mol: A multi-modal large language model for molecular science with graph, image, and text," *Computers in Biology and Medicine*, vol. 171, p. 108073, 2024.
- [15] H. Wang, Z. Jiang, Y. You, Y. Han, G. Liu, J. Srinivasa, R. Kompella, Z. Wang *et al.*, "Graph mixture of experts: Learning on large-scale graphs with explicit diversity modeling," *Advances in Neural Information Processing Systems*, vol. 36, 2024.
- [16] W. Fedus, B. Zoph, and N. Shazeer, "Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity," *Journal of Machine Learning Research*, vol. 23, no. 120, pp. 1–39, 2022.
- [17] N. Jafarpour, D. Precup, M. Izadi, and D. Buckeridge, "Using hierarchical mixture of experts model for fusion of outbreak detection methods," in *AMIA Annual Symposium Proceedings*, vol. 2013. American Medical Informatics Association, 2013, p. 663.
- [18] A. Goyal, N. Kumar, T. Guha, and S. S. Narayanan, "A multimodal mixture-of-experts model for dynamic emotion prediction in movies," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 2822–2826.
- [19] D. Dai, W. Jiang, J. Zhang, Y. Lyu, Z. Sui, and B. Chang, "Mixture-of-experts for biomedical question answering," in *CCF International Conference on Natural Language Processing and Chinese Computing*. Springer, 2023, pp. 604–615.
- [20] P. Zhu, Y. Sun, B. Cao, and Q. Hu, "Task-customized mixture of adapters for general image fusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 7099–7108.
- [21] Z. Liu, G. Wang, J. Wang, J. Zheng, and S. Z. Li, "Co-modeling the sequential and graphical routes for peptide representation learning," *arXiv e-prints*, pp. arXiv–2310, 2023.
- [22] P. Bhadra, J. Yan, J. Li, S. Fong, and S. W. Siu, "Ampep: Sequence-based prediction of antimicrobial peptides using distribution patterns of amino acid properties and random forest," *Scientific reports*, vol. 8, no. 1, p. 1697, 2018.
- [23] N. Shazeer, A. Mirhoseini, K. Maziarz, A. Davis, Q. Le, G. Hinton, and J. Dean, "Outrageously large neural networks: The sparsely-gated mixture-of-experts layer," *arXiv preprint arXiv:1701.06538*, 2017.