OPPI-GRF: Optimizing Protein-Protein Interaction Prediction with Graph-Based Representation and Fusion

Anusuya Krishnan United Arab Emirates University Al Ain, UAE anusuyababy@uaeu.ac.ae Isaias Mehari Ghebrehiwet United Arab Emirates University Al Ain, UAE ighebrehiwet@uaeu.ac.ae

Abstract

Protein-protein interactions (PPIs) are essential to various biological processes, including cell signaling and metabolic regulation, making their accurate prediction vital for advancing drug discovery and therapeutic development. The exploration of PPIs has led to the development of several AI-based models that leverage recent advancements in artificial intelligence, primarily focusing on features extracted from diverse sources of protein data. In our study, we focus on PPIs by calculating the amino acid composition of the relevant proteins. We construct a PPI network represented as a graph, where each node signifies a protein pair, and an edge connects nodes if they share a common protein. Each node is linked to a feature vector encapsulating significant information about the proteins. We then employ the GraphFormer model to extract feature embeddings directly from the protein sequences, capturing intricate patterns within the data. To enhance our predictions, we implement a fusion technique that combines the amino acid composition features with the embeddings generated by the GraphFormer model using a Feature Attention Network. This network assigns weights to the features, allowing us to emphasize the most critical information from both the amino acid composition and embedding features. By integrating these two types of data, we leverage both sequence-level information, which describes the biochemical properties of proteins, and structural information from the graph embeddings, which illustrates how proteins interact. Finally, we evaluate our model's performance using key metrics such as accuracy, precision, recall, and F1 score. Our approach shows improved results compared to existing methods, effectively merging both sequence-level and structural information essential for understanding protein interactions.

1 Introduction

Proteins are essential building blocks of life, made up of amino acids that genes encode. They form peptides that fold into various proteins, playing a key role in building tissues and performing functions like catalyzing reactions, transporting molecules, responding to pathogens, and signaling between cells. Important processes such as DNA replication and transcription depend on specific proteins, which often interact through protein-protein interactions (PPIs) [1]. Identifying and targeting these interactions help researchers develop new therapies, including small-molecule inhibitors and biologics, for various diseases. PPI research is crucial for creating diagnostic tools and treatments in fields like drug design and medical diagnostics [2]. PPI identification methods can be traditional experimental techniques or computational approaches. However, traditional methods, like yeast two-hybrid assays and protein chips, can be time-consuming, labor-intensive, and costly, making them less practical for widespread use [1, 2].

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

To address these limitations, computational methods for predicting PPIs have rapidly evolved. However, these approaches encounter two significant challenges: the need to accurately identify various specific categories of PPIs and the necessity to achieve high predictive performance. Numerous computational techniques have been developed to explore protein-protein interaction (PPI) networks in various organisms, utilizing different types of protein information, including protein sequences, structures, gene co-expression, and gene ontology [3, 4]. Many studies focused on PPI prediction have utilized manually crafted features as the initial feature vectors, which are then fed into deep learning models to capture relevant patterns from the raw data [5].

Recently, there has been a notable rise in the use of graph neural networks (GNNs), establishing them as vital tools for graph-based applications. For instance, graph convolution has been utilized to uncover associations between miRNAs and drug resistance, framing the issue as a link prediction challenge to forecast these associations [6, 7]. A specific model for predicting protein-protein interactions (PPIs) represents the PPI network as a graph and employs the conjoint-triad (CT) method to derive node features. In this PPI network graph, each node symbolizes a protein, while an edge defines the relationship between protein pairs, indicating whether they interact or not. A signed variational graph auto-encoder (S-VGAE) that incorporates graph convolution layers is used to learn compact representations of the nodes. The learned representations of protein pairs are then concatenated and input into a neural network classifier to predict PPIs [7, 8].

All of these studies addressing different biomedical challenges have utilized graph convolutional networks (GCNs) [9]. However, existing variants of graph neural networks, including GCNs, may encounter issues such as suspended animation and over-smoothing, largely due to their heavy reliance on graph connections. To address the challenges associated with predicting protein-protein interactions (PPIs), we propose a comprehensive framework that utilizes attention-based fusion of sequence-derived and structural graph representations. The contributions of our proposed work are outlined as follows:

- Amino Acid Composition Calculation: We calculate the amino acid composition of the involved proteins to provide essential biochemical insights that serve as foundational features for PPI prediction [10]. This calculation helps in identifying the specific properties and behaviors of the proteins, which can influence their interactions.
- **Graph-Based PPI Representation:** We construct a PPI network as a graph representation, where each node corresponds to a protein pair and edges indicate relationships based on shared proteins, facilitating the exploration of interactions. This graph structure allows us to apply advanced analytical techniques to better capture the dynamics and intricacies of protein interactions.
- **Feature Fusion:** We apply a feature attention network (FA-Net) to fuse the amino acid composition features with the embeddings generated by the GraphFormer model [11], allowing us to emphasize the most relevant information and improve PPI prediction accuracy [12].

The paper is organized as follows: Section 2 reviews related work, Section 3 describes the methodology, Section 4 discusses the results, and Section 5 concludes with findings and future directions.

2 Related Works

Several computational approaches have been developed to categorize protein-protein interactions (PPIs) using various machine learning algorithms, such as SVM, naive Bayes, decision trees, and random forests [13, 16]. Among these, SVM has been widely utilized for PPI prediction. However, more recent studies have demonstrated that random forest and rotation forest-based approaches often outperform SVM-based methods in PPI classification [14, 15].

With advancements in technology, deep learning algorithms have been increasingly adopted for PPI prediction, showing superior performance compared to traditional machine learning methods [17]. Techniques like stacked auto-encoders have been employed to learn compact representations of protein sequence features. Other deep learning models, such as DeepInteract and EnsDNN, have leveraged deep neural networks (DNN) to predict protein interactions. For example, sequence-statistics-content (SSC) encoding formats, combined with convolutional neural networks (CNN), have been used to extract and utilize information from protein sequences for PPI prediction [17, 18, 22].



Figure 1: Proposed Architecture of Optimizing Protein-Protein Interaction Prediction with Graph-Based Representation and Fusion.

While earlier approaches relied heavily on hand-engineered features as input, the rapidly advancing field of deep learning has enabled automatic feature engineering [19, 20]. Recent studies have successfully applied auto-engineered features for PPI prediction, such as frameworks utilizing deep neural networks (DNN) that automatically learn features from protein sequences using CNNs and long-short-term memory (LSTM) networks [21, 22, 23, 24]. These techniques provide more effective and automated means of predicting PPIs, moving beyond manual feature extraction.

Recent studies have explored embedding techniques and recurrent neural network-based architectures to predict interactions directly from protein sequences [25, 26]. Other deep learning frameworks, such as models combining convolutional layers with bi-directional residual gated recurrent units, have been proposed to extract both local and global features from protein sequences [27]. These methods often rely on pre-trained embeddings of protein sequences as input.

Deep learning algorithms are highly effective at handling high-dimensional data and capturing hidden associations, even in datasets with complex, multi-modal distributions [28]. Recent approaches to PPI prediction have utilized multiple sources of protein information, such as protein sequences, 3D structures, and gene ontology, developing deep multi-modal models that combine these data sources. These models leverage the latest deep learning techniques to extract meaningful features from various protein data modalities [27, 28].

In addition to sequence-based approaches, protein-protein interactions can also be modeled graphically as PPI networks, where nodes represent proteins and edges denote interactions between them. Graph-based deep learning models (Gresnet, GraphBERT) have been used to predict PPIs by learning low-dimensional features from graph structures. For example, manually crafted sequence-based methods like the conjoint-triad (CT) method have been employed to generate node feature vectors for such graphs [29, 30, 31, 32].

The current work builds on these developments by utilizing graph-based neural networks to predict protein interactions. We formulate the PPI prediction task as a node classification problem, where each node represents a protein pair, and their interactions are classified. In this study, we employ the GraphFormer model to learn hidden representations of the nodes [11]. The feature vectors are derived directly from protein sequences, capturing intricate patterns and relationships within the data, which are then used for accurate PPI prediction.

3 Methodology

The proposed method, as depicted in Figure 1, consists of several essential stages: calculation of amino acid composition, construction of the graph, extraction of feature embeddings, integration of features through fusion, and classification of protein-protein interactions (PPIs).

3.1 Amino Acid Composition Analysis

The amino acid composition for each protein involved in the protein-protein interaction (PPI) network, denoted as D, is calculated to derive essential biochemical insights [10]. The amino acid composition C of a protein sequence S from the dataset D can be expressed as:

$$C_{aa} = \frac{N_{aa}}{N_{total}} \tag{1}$$

where C_{aa} represents the composition of amino acid aa, N_{aa} is the count of amino acid aa in the sequence S, and N_{total} is the total number of amino acids in the sequence S.

This composition yields a vector that encapsulates the normalized frequency of each of the 20 standard amino acids for proteins in the dataset D. Analyzing these amino acid compositions provides critical insights into the biochemical properties of the proteins, serving as foundational features for subsequent PPI prediction tasks.

3.2 Graph-Based PPI Representation

Euclidean distance was selected due to its ability to directly measure the absolute differences between feature values, making it well-suited for assessing the overall similarity between proteins in the amino acid composition matrix. Euclidean distance is computed for the amino acid composition matrix X to obtain a distance matrix D: The Euclidean distance between proteins i and j is calculated as:

$$D_{ij} = \sqrt{\sum_{k} (X_{ki} - X_{kj})^2}$$
(2)

where D_{ij} represents the Euclidean distance between the amino acid composition vectors of proteins i and j. Here, X_{ki} denotes the amino acid composition of the k-th amino acid in protein i, and X_{kj} denotes the amino acid composition of the k-th amino acid in protein j. The sum of squared differences $\sum_{k} (X_{ki} - X_{kj})^2$ measures the dissimilarity between the two proteins in terms of their amino acid composition. A binary distance matrix B is then created using a threshold τ :

$$\tau = k \times (\mu + 3\sigma) \tag{3}$$

where μ is the mean and σ is the standard deviation of the Euclidean distance values in the matrix D. The parameter k is optimized through hyperparameter tuning to determine the appropriate threshold for binarization. This binarization step, which converts the distance matrix into a binary form based on the threshold, plays a crucial role in identifying protein-protein interactions. The binary distance matrix B is defined by:

$$B_{ij} = \begin{cases} 1, & \text{if } D_{ij} < \tau \\ 0, & \text{if } D_{ij} \ge \tau \end{cases}$$

$$\tag{4}$$

In this matrix, B_{ij} represents the presence or absence of an edge between proteins *i* and *j*. If the Euclidean distance D_{ij} is less than the threshold τ , B_{ij} is set to 1, indicating a connection between the proteins. If D_{ij} is greater than or equal to τ , B_{ij} is set to 0, indicating no connection.

A graph G = (V, E) is constructed from the binary distance matrix B, where V represents the set of nodes, with each node corresponding to a protein and E represents the edges, where an edge exists between two proteins if they are connected in the binary matrix B.

This graph captures the protein-protein interaction network based on the biochemical similarities derived from their amino acid compositions, enabling further analysis of the interaction patterns.

3.3 GraphFormer Embeddings

The binary matrix B serves as the adjacency matrix A of the graph G. The feature matrix X, derived from the amino acid composition of the proteins, contains the corresponding features for each node

 $v_i \in V$. The GraphFormer model utilizes both the adjacency matrix A and the feature matrix X to generate node embeddings H [11]:

$$H = \text{GraphFormer}(A, X) \tag{5}$$

At each layer l in the GraphFormer, the node embeddings are updated through a graph attention mechanism and a feed-forward network:

$$H^{(l+1)} = \text{Attention}(H^{(l)}, A) + \text{FFN}(H^{(l)})$$
(6)

where $H^{(l)}$ is the node embedding at layer l, Attention $(H^{(l)}, A)$ is the multi-head attention mechanism using the adjacency matrix A, and $FFN(H^{(l)})$ is a feed-forward network applied to the node embeddings.

The final embeddings after L layers of the GraphFormer model, denoted as $H^{(L)}$, represent the learned graph embeddings that capture both the structural information from the graph and the biochemical properties of the proteins:

$$H^{(L)} = \text{GraphFormer}(A, X) \tag{7}$$

These embeddings $H^{(L)}$ can be used for downstream tasks, such as protein-protein interaction (PPI) prediction.

3.4 Feature Fusion Using Feature Attention Network

The amino acid composition features F_{aa} and the embeddings H generated by the GraphFormer model are combined using a Feature Attention Network (FA-Net) [12]. The fused feature vector F is calculated as follows:

$$F = \text{FA-Net}(F_{aa}, H) \tag{8}$$

The feature attention mechanism can be expressed as:

$$F = \sum_{i} \alpha_{i} F_{i} \tag{9}$$

where F_i represents the individual features from both the amino acid composition and embedding features, and α_i denotes the attention weights assigned to each feature. The attention weights are computed as:

$$\alpha_i = \frac{\exp\left(\operatorname{score}(F_i)\right)}{\sum_j \exp\left(\operatorname{score}(F_j)\right)} \tag{10}$$

Here, score(F_i) is a scoring function that evaluates the importance of the feature F_i . The resulting fused feature vector F is then used to predict protein-protein interactions (PPIs):

$$PPI_{prediction} = Model(F)$$
(11)

This approach enables the model to focus on the most relevant information from both the amino acid composition features and the graph embeddings, ultimately improving the accuracy of PPI predictions.

4 Experimental Results

This section presents the outcomes of our proposed OPPI-GRF model, including detailed performance metrics and comparative analyses against baseline methods.

4.1 Experimental Overview

Datasets The datasets employed in this study are publicly accessible: the Pan's Protein-Protein Interaction (PPI) dataset ¹ and the Database of Interacting Proteins (DIP) dataset.² The Pan's human Protein-Protein Interaction dataset, which includes both positive and negative samples, has been utilized as a benchmark to evaluate the proposed approach. The Human Protein Reference Database³ (HPRD) serves as the source of interacting pairs for the Pan's human PPI dataset. The non-interacting dataset consists of pairs derived from the Negatome database⁴, in addition to pairings of proteins from distinct subcellular localizations. Moreover, the DIP dataset offers experimentally validated interactions, emphasizing real-world PPI predictions. This dataset contains highly reliable interaction data obtained from experimental studies, rendering it an essential resource for evaluating the accuracy of computational models in predicting physical interactions between proteins. Table 1 summarizes the key characteristics of these two datasets.

Table 1	: Summary	of Interacting	and Non-Interacting	Proteins for	Selected Datasets
---------	-----------	----------------	---------------------	--------------	-------------------

Dataset	Number of Interacting Pairs	Number of Non-Interacting Pairs
Pan's PPI Dataset	6,380	3,173
DIP Dataset	2,700	2,300

Evaluation Metrics To assess the performance of our proposed OPPI-GRF model, we utilize several standard evaluation metrics. Accuracy is defined as the proportion of correctly predicted protein-protein interactions out of the total predictions made. Precision measures the ratio of true positive predictions to the total number of predicted positives, indicating how many of the predicted interactions are correct. Recall reflects the model's ability to identify relevant interactions by calculating the ratio of true positive predictions to the total number of actual positives. Finally, the F1-Score provides a balanced measure of prediction performance by computing the harmonic mean of precision and recall.

Experimental Setup In this work, we utilize the GraphFormer model to extract feature embeddings from our protein-protein interaction (PPI) datasets. The model is trained on a GeForce GTX 1080 GPU. For this experiment, we use deep learning libraries including PyTorch, scikit-learn, transformers, and networkx. The GraphFormer model is configured with 4 to 6 hidden layers, a hidden dimension of 256 or 512, 8 attention heads, a dropout rate ranging from 0.1 to 0.3, and a learning rate of 0.001. The batch size is set between 32 and 128, and training is performed for a maximum of 200 epochs.

To integrate additional protein features, we implement a Feature Attention Network (FA-Net), which merges the amino acid composition (AAC) features with the embeddings produced by GraphFormer. The FA-Net includes 1 to 3 hidden layers with a hidden dimension of 128 to 256 and a dropout rate of 0.2 to mitigate overfitting. The Adam optimizer is used to minimize the binary cross-entropy loss, and we employ an early stopping mechanism to prevent overfitting.

We follow a standard 80:20 train-test split for model training and evaluation. The dataset is split into 80% for training and 20% for testing to ensure the model learns from one subset while being evaluated on unseen data. Performance is measured using accuracy, precision, recall, and F1-score. Accuracy reflects overall correctness, precision evaluates true positives among predicted positives, and recall assesses true positives among actual positives. The F1-score balances precision and recall, particularly useful for imbalanced datasets.

4.2 Results on Pan's Human PPI Dataset

The performance of the proposed model on Pan's Human PPI dataset is summarized in Table 2. To validate the effectiveness of our approach, we designed several baselines. In baseline-1, we used the same graph structure but replaced the feature vectors with traditional conjoint-triad-based features. The accuracy, precision, recall, and F1-score for this baseline are 95.82%, 96.34%, 94.12%, and

¹http://www.csbio.sjtu.edu.cn/bioinf/LR_PPI/Data.htm

²http://dip.doe-mbi.ucla.edu/dip/Main.cgi

³https://www.hsls.pitt.edu/obrc/index.php?page=URL1055173331

⁴https://ngdc.cncb.ac.cn/databasecommons/database/id/207

93.56%, respectively (as reported in Table 2). These results demonstrate the limitations of manually crafted feature representations.

In contrast, our method achieves a significantly higher accuracy of 98.15% by using GraphFormer combined with amino acid composition features and sequence embeddings. These results highlight that the SeqVec-based embeddings, when combined with the attention mechanism of GraphFormer, outperform traditional feature engineering approaches, further underscoring the effectiveness of our proposed method for predicting protein-protein interactions on Pan's Human PPI dataset.

Method	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
Baseline-1 (Conjoint-Triad + AAC)	95.82	96.34	94.12	93.56
Baseline-2 (SeqVec + GraphFormer)	96.72	96.88	95.38	94.56
Proposed (GraphFormer + FA-Net)	98.89	99.42	98.92	99.34

Table 2: Performance comparison on Pan's Human PPI Dataset

The high precision and recall values indicate that the model is both accurate and sensitive in identifying true protein interactions, and the F1-score further demonstrates balanced performance. The experimental outcomes affirm the effectiveness of using feature attention network (FA-Net) for PPI prediction, combining both graph structural information and biochemical features.

4.3 Results on DIP Dataset

We have also validated the effectiveness of the proposed approach on the DIP dataset. The table 3 summarizes the performance comparison of various methods for predicting protein-protein interactions (PPIs) on the DIP Dataset. It includes two baseline methods and the proposed approach.

Method	Accuracy (%)	Precision (%)	Recall (%)	F1-score (%)
Baseline-1 (Conjoint-Triad + AAC)	88.82	89.34	88.12	87.56
Baseline-2 (SeqVec + GraphFormer)	89.72	89.88	90.38	88.56
Proposed (GraphFormer + FA-Net)	90.19	90.22	90.19	90.47

 Table 3: Performance comparison on DIP Dataset

From table 3, baseline-1 employs the Conjoint-Triad method combined with amino acid composition (AAC), achieving an accuracy of 88.82%, precision of 89.34%, recall of 88.12%, and F1-score of 87.56%. Baseline-2 utilizes SeqVec embeddings integrated with GraphFormer, resulting in improved performance with an accuracy of 89.72%, precision of 89.88%, recall of 90.38%, and F1-score of 88.56%. In contrast, the proposed method, which combines GraphFormer with AAC features, achieves the highest metrics, boasting an accuracy of 90.15%, precision of 90.42%, recall of 90.89%, and F1-score of 90.34%. These results demonstrate the effectiveness of the proposed approach in capturing complex protein interactions more accurately than traditional methods.

4.4 Performance Comparison of Existing Methods

Numerous studies have focused on predicting protein-protein interactions (PPI) using artificial intelligence-based approaches. Table 4 presents a comprehensive comparison of various existing methods used for protein-protein interaction (PPI) prediction on two datasets: Pan's Human PPI Dataset and DIP Dataset. It evaluates the performance of each method using five key metrics: accuracy, recall, specificity, precision, and F1-score.

In the Pan's Dataset section, the Long Short-Term Memory (LSTM) network achieves an accuracy of 97.20%, while the Convolutional Neural Network (CNN) slightly lags behind with 97.07%. A Multimodal PPI Model, which integrates multiple data sources, reaches an accuracy of 97.52%. Notably, GraphBERT, a model leveraging graph-based embeddings, demonstrates the best performance with an accuracy of 98.13%. The proposed method outshines all others with an accuracy of 98.89%, also achieving the highest values in precision, recall, and F1-score.

In the DIP Dataset section, the methods exhibit generally lower accuracy compared to the Pan's Dataset. The Bidirectional LSTM (BiLSTM) records an accuracy of 88.50%, followed by the CNN-

Performance Metrics on Pan's Dataset						
Method	Accuracy (%)	Recall (%)	Specificity (%)	Precision (%)	F1-score (%)	
LSTM	97.20	98.07	95.04	97.99	98.03	
CNN	97.07	98.19	93.46	97.98	98.09	
Multi-modal PPI Model	97.52	98.20	95.92	98.26	98.23	
GraphBERT	98.13	98.84	96.18	98.62	98.73	
Ours	98.89	98.92	98.03	99.42	99.34	
Performance Metrics on DIP Dataset						
(lr)1-6 BiLSTM	88.50	89.30	86.70	87.50	88.10	
CNN-RNN	87.20	88.50	85.90	86.80	87.20	
GAT	88.00	89.50	86.90	88.20	88.80	
GCN	89.30	89.90	87.20	89.50	89.10	
Ours	90.19	90.19	90.11	90.22	90.47	

Table 4: Comparison of existing methods against proposed method on two datasets

RNN combination at 87.20%. The Graph Attention Network (GAT) achieves 88.00%, while the Graph Convolutional Network (GCN) performs best in this section with an accuracy of 89.30%. The proposed method again leads with a remarkable accuracy of 90.19%, along with superior precision, recall, and F1-score values. These results highlight the effectiveness of the proposed method in PPI prediction, consistently outperforming existing models across both datasets and showcasing its potential advantages in accurately identifying protein interactions.

4.5 Discussion

The proposed method demonstrates a significant advancement in predicting protein-protein interactions (PPIs) through the integration of GraphFormer and Feature Attention Network (FA-Net) [11, 12]. By representing each protein pair as a node in a graph, our approach effectively captures the complex relationships inherent in biological data. This novel representation allows for a more comprehensive feature extraction process, leveraging both amino acid composition and the rich embeddings generated by GraphFormer. The results obtained indicate that our model achieves superior accuracy compared to traditional methods, highlighting the importance of advanced feature representations in enhancing PPI prediction performance.

In comparison to existing techniques, our proposed model not only surpasses traditional machine learning methods but also outperforms several state-of-the-art deep learning approaches. The improved accuracy and robustness of our method can be attributed to the synergistic effects of integrating feature attention mechanisms with graph-based embeddings. By emphasizing the most relevant features in the input data, FA-Net enhances the learning process, allowing the model to focus on critical interactions that drive the prediction of PPIs. This capability not only improves the accuracy of the predictions but also increases the interpretability of the model by revealing the underlying biological significance of the selected features.

5 Conclusion

In this study, we presented a novel approach for predicting protein-protein interactions (PPIs) that leverages a combination of advanced graph-based methodologies and attention mechanisms. Our proposed model demonstrated superior performance compared to traditional and recent deep learning methods in key metrics such as precision, recall, and F1-score. This highlights the effectiveness of integrating graph representations with feature attention networks to capture complex biological relationships more effectively. The promising results of our approach validate its potential as a reliable tool for PPI prediction and emphasize the importance of incorporating diverse features and sophisticated learning mechanisms in computational biology. Our findings suggest that further exploration of this methodology could lead to enhanced understanding and prediction of protein interactions, ultimately contributing to advancements in drug discovery and therapeutic interventions. Future work will focus on refining the model's scalability and exploring its integration with other biological data sources to enhance its predictive capabilities.

References

- [1] Alberts, B. (1998) The cell as a collection of protein machines: Preparing the next generation of molecular biologists. *Cell*, 92, 291–294.
- [2] Raman, K. (2010) Construction and analysis of protein–protein interaction networks. *Autom. Exp.*, 2, 2.
- [3] Skrabanek, L., Saini, H. K., Bader, G. D. & Enright, A. J. (2008) Computational prediction of protein-protein interactions. *Mol. Biotechnol.*, 38, 1–17.
- [4] Ding, Z. & Kihara, D. (2018) Computational methods for predicting protein-protein interactions using various protein features. *Curr. Protoc. Protein Sci.*, 93, e62.
- [5] Mathews, N., Tran, T., Rekabdar, B. & Ekenna, C. (2021) Predicting human-pathogen protein-protein interactions using Natural Language Processing methods. *Informatics in Medicine* Unlocked, 26, p.100738.
- [6] Huang, Y.A., Hu, P., Chan, K. C. & You, Z.H. (2020) Graph convolution for predicting associations between miRNA and drug resistance. *Bioinformatics*, 36, 851–858.
- [7] Gonzalez-Lopez, F., Morales-Cordovilla, J. A., Villegas-Morcillo, A., Gomez, A. M. & Sanchez, V. (2018) End-to-end prediction of protein-protein interaction based on embedding and recurrent neural networks. In 2018 IEEE International Conference on Bioinformatics and Biomedicine (BIBM), 2344–2350 (IEEE).
- [8] Yang, F., Fan, K., Song, D. & Lin, H. (2020) Graph-based prediction of protein-protein interactions with attributed signed graph embedding. *BMC Bioinform.*, 21, 1–16.
- [9] Jha, K., Saha, S. & Singh, H. (2022) Prediction of protein-protein interaction using graph neural networks. *Sci. Rep.*, 12, 1–12.
- [10] Zhai, Y., Chen, Y., Teng, Z. & Zhao, Y. (2020) Identifying antioxidant proteins by using amino acid composition and protein-protein interactions. *Frontiers in Cell and Developmental Biology*, 8, p.591487.
- [11] Yang, J., Liu, Z., Xiao, S., Li, C., Lian, D., Agrawal, S., Singh, A., Sun, G. & Xie, X. (2021) Graphformers: Gnn-nested transformers for representation learning on textual graph. Advances in Neural Information Processing Systems, 34, 28798–28810.
- [12] Qin, X., Wang, Z., Bai, Y., Xie, X. & Jia, H. (2020) FFA-Net: Feature fusion attention network for single image dehazing. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 34, No. 07, pp. 11908–11915).
- [13] Wang, L., You, Z.-H., Xia, S.X., Liu, F., Chen, X., Yan, X. & Zhou, Y. (2017) Advancing the prediction accuracy of protein-protein interactions by utilizing evolutionary information from position-specific scoring matrix and ensemble classifier. *Journal of Theoretical Biology*, 418, 105–110.
- [14] You, Z.H., Li, J., Gao, X., He, Z., Zhu, L., Lei, Y.K., & Ji, Z. (2015) Detecting protein-protein interactions with a novel matrix-based protein sequence representation and support vector machines. *BioMed Research International*, 2015, 867516.
- [15] You, Z.H., Yu, J.Z., Zhu, L., Li, S. & Wen, Z.K. (2014) A mapreduce based parallel SVM for large-scale predicting protein-protein interactions. *Neurocomputing*, 145, 37–43.
- [16] Zhou, C., Yu, H., Ding, Y., Guo, F. & Gong, X.-J. (2017) Multi-scale encoding of amino acid sequences for predicting protein interactions using gradient boosting decision tree. *PLoS ONE*, 12, e0181426.
- [17] Elnaggar, A., Heinzinger, M., Dallago, C., Rihawi, G., Wang, Y., Jones, L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M., & Bhowmik, D. (2020) ProtTrans: Towards cracking the language of Life's code through self-supervised deep learning and high performance computing. *arXiv* preprint arXiv:2007.06225.

- [18] Wang, Y.B., You, Z.H., Li, X., Jiang, T.H., Chen, X., Zhou, X., & Wang, L. (2017) Predicting protein–protein interactions from protein sequences by a stacked sparse autoencoder deep neural network. *Molecular BioSystems*, 13(7), 1336-1344.
- [19] Sun, T., Zhou, B., Lai, L. & Pei, J. (2017) Sequence-based prediction of protein-protein interaction using a deep-learning algorithm. *BMC Bioinform.*, 18, 1–8.
- [20] Patel, S., Tripathi, R., Kumari, V. & Varadwaj, P. (2017) Deepinteract: Deep neural network based protein-protein interaction prediction tool. *Curr. Bioinform.*, 12, 551–557.
- [21] Zhang, L., Yu, G., Xia, D. & Wang, J. (2019) Protein-protein interactions prediction based on ensemble deep neural networks. *Neurocomputing*, 324, 10–19.
- [22] Wang, Y., Li, Z., Zhang, Y., Ma, Y., Huang, Q., Chen, X., Dai, Z., & Zou, X. (2021) Performance improvement for a 2D convolutional neural network by using SSC encoding on protein–protein interaction tasks. *BMC Bioinformatics*, 22, 1-16.
- [23] Li, H., Gong, X.J., Yu, H. & Zhou, C. (2018) Deep neural network based predictions of protein interactions using primary sequences. *Molecules*, 23, 1923.
- [24] Chen, M., Ju, C.J.T., Zhou, G., Chen, X., Zhang, T., Chang, K.W., Zaniolo, C., & Wang, W. (2019) Multifaceted protein–protein interaction prediction based on Siamese residual RCNN. *Bioinformatics*, 35(14), i305-i314.
- [25] You, Z.-H., Lei, Y.K., Gui, J., Huang, D.S. & Zhou, X. (2010) Using manifold embedding for assessing and predicting protein interactions from high-throughput experimental data. *Bioinformatics*, 26, 2744–2751.
- [26] Li, X., Yan, X., Gu, Q., Zhou, H., Wu, D., & Xu, J. (2019) Deepchemstable: Chemical stability prediction with an attention-based graph convolution network. J. Chem. Inf. Model., 59, 1044–1049.
- [27] Fout, A.M. (2017) Protein Interface Prediction Using Graph Convolutional Networks. Ph.D. thesis, Colorado State University.
- [28] Xue, Y., Liu, Z., Fang, X. & Wang, F. (2022) Multimodal pre-training model for sequence-based prediction of protein-protein interaction. In *Machine Learning in Computational Biology* (pp. 34–46). PMLR.
- [29] Kipf, T. N. & Welling, M. (2016) Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907.
- [30] Zhang, J. & Meng, L. (2019) Gresnet: Graph residual network for reviving deep GNNs from suspended animation. arXiv preprint arXiv:1909.05729.
- [31] Zhang, J., Zhang, H., Xia, C. & Sun, L. (2020) Graph-bert: Only attention is needed for learning graph representations. arXiv preprint arXiv:2001.05140.
- [32] Jha, K., Karmakar, S., & Saha, S. (2023) Graph-BERT and language model-based framework for protein–protein interaction identification. *Scientific Reports*, 13(1), 5663.