

# SemiGNN-PPI: Self-Ensembling Multi-Graph Neural Network for Efficient and Generalizable Protein-Protein Interaction Prediction

Ziyuan Zhao<sup>1,2</sup>, Peisheng Qian<sup>1</sup>, Xulei Yang<sup>1\*</sup>, Zeng Zeng<sup>3</sup>, Cuntai Guan<sup>2</sup>,  
Wai Leong Tam<sup>4</sup> and Xiaoli Li<sup>1,2</sup>

<sup>1</sup>Institute for Infocomm Research (I<sup>2</sup>R), A\*STAR, Singapore

<sup>2</sup>School of Computer Science and Engineering (SCSE), Nanyang Technological University, Singapore

<sup>3</sup>School of Microelectronics, Shanghai University, China

<sup>4</sup>Genome Institute of Singapore (GIS), A\*STAR, Singapore

{zhaoz, qianp, yangx}@i2r.a-star.edu.sg, zengz@shu.edu.cn, ctguan@ntu.edu.sg,  
tamwl@gis.a-star.edu.sg, xlli@i2r.a-star.edu.sg

## Abstract

Protein-protein interactions (PPIs) are crucial in various biological processes and their study has significant implications for drug development and disease diagnosis. Existing deep learning methods suffer from significant performance degradation under complex real-world scenarios due to various factors, *e.g.*, label scarcity and domain shift. In this paper, we propose a self-ensembling multi-graph neural network (SemiGNN-PPI) that can effectively predict PPIs while being both efficient and generalizable. In SemiGNN-PPI, we not only model the protein correlations but explore the label dependencies by constructing and processing multiple graphs from the perspectives of both features and labels in the graph learning process. We further marry GNN with Mean Teacher to effectively leverage unlabeled graph-structured PPI data for self-ensemble graph learning. We also design multiple graph consistency constraints to align the student and teacher graphs in the feature embedding space, enabling the student model to better learn from the teacher model by incorporating more relationships. Extensive experiments on PPI datasets of different scales with different evaluation settings demonstrate that SemiGNN-PPI outperforms state-of-the-art PPI prediction methods, particularly in challenging scenarios such as training with limited annotations and testing on unseen data.

## 1 Introduction

Protein-protein Interactions (PPIs) are central to various cellular functions and processes, such as signal transduction, cell-cycle progression, and metabolic pathways [Acuner Ozbabacan *et al.*, 2011]. Therefore, the identification and characterization of PPIs are of great importance for understanding protein functions and disease occurrence, which can potentially facilitate therapeutic target identification [Petta *et al.*, 2016] and the novel drug design [Skrabaneck

*et al.*, 2008]. In past decades, high-throughput experimental methods, *e.g.*, yeast two-hybrid screens (Y2H) [Fields and Song, 1989], and mass spectrometric protein complex identification (MS-PCI) [Ho *et al.*, 2002] have been developed to identify PPIs. Nevertheless, genome-scale experiments are expensive, tedious, and time-consuming while suffering from high error rates and low coverage [Luo *et al.*, 2015]. As such, there is an urgent need to establish reliable computational methods to identify PPIs with high quality and accuracy.

In recent years, a large variety of high-throughput computational approaches for PPI prediction have been proposed, which can be broadly divided into two groups: classic machine learning (ML)-based methods [Browne *et al.*, 2007; Lin and Chen, 2013; Guo *et al.*, 2008; Wong *et al.*, 2015; Chen and Liu, 2005] and deep learning (DL)-based methods [Sun *et al.*, 2017; Du *et al.*, 2017; Hashemifar *et al.*, 2018; Chen *et al.*, 2019a; Lv *et al.*, 2021]. Compared to classic ML methods, DL algorithms are capable of processing complicated and large-scale data and extracting useful features automatically, achieving significant success in a diverse range of bioinformatics applications [Min *et al.*, 2017; Soleymani *et al.*, 2022], including PPI prediction [Soleymani *et al.*, 2022]. Most existing DL-based methods treat interactions as independent instances, ignoring protein correlations. PPI can be naturally formulated as graph networks with proteins and interactions represented as nodes and edges, respectively [Margolin *et al.*, 2006; Pio *et al.*, 2020]. To improve PPI prediction performance, recent works [Yang *et al.*, 2020; Lv *et al.*, 2021] have been proposed to investigate the correlations between PPIs using various graph neural network (GNN) architectures [Kipf and Welling, 2016; Xu *et al.*, 2019]. However, they are limited by ignoring learning label dependencies for multi-type PPI prediction. It has recently become common practice to employ Graph Convolutional Networks (GCNs) to capture label correlation in a wide range of multi-label tasks [Chen *et al.*, 2019b; Wang *et al.*, 2020]. Nevertheless, multi-label learning utilizing label graphs predominantly works in the visual domain and has yet to be extended to PPI prediction tasks.

In general, a desired PPI prediction framework should be efficient, transferable, and generalizable, whereas two ma-

\*Corresponding author

major bottlenecks deriving from imperfect datasets have hindered the development of such models. **Label scarcity:** Despite the tremendous progress in PPI research using various computational and experimental methods, many interactions still need to be annotated from experimental data. Consequently, only a small portion of labeled samples can be used for model training. It can be a significant bottleneck in obtaining robust and accurate PPI prediction models. **Domain shift:** Most existing methods are only developed and validated using in-distribution data (*i.e.*, trainset-homologous testsets), receiving severe performance degradation when being deployed to unseen data with different distributions (*i.e.*, trainset-heterologous testsets). Although [Lv *et al.*, 2021] design new evaluations to better reflect model generalization, giving instructive and consistent assessment across datasets, the domain shift issue still needs to be fully explored for PPI prediction. Therefore, how to deal with imperfect data for improving model efficiency and generalization remains a vital issue in PPI prediction. Recent studies [Zhang *et al.*, 2021; Zhao *et al.*, 2022] show that self-ensemble methods with semi-supervised learning (SSL) [Laine and Aila, 2017; Tarvainen and Valpola, 2017] have demonstrated effectiveness in addressing both label scarcity and domain shift.

In this work, to tackle the above challenges and limitations, we propose an efficient and generalizable PPI prediction framework, referred to as **Self-ensembling multi-Graph Neural Network (SemiGNN-PPI)**. Firstly, we propose leveraging graph structure to model protein correlations and label dependencies for multi-graph learning. Specifically, we learn inter-dependent classifiers to extract information from the label graph, which are then applied to the protein representations aggregated by neighbors in the protein graph for multi-type PPI prediction. Secondly, we propose combining GNN with Mean Teacher [Tarvainen and Valpola, 2017], a powerful SSL model, to explore unlabeled data for self-ensemble graph learning. In our framework, the student model learns to classify the labeled data accurately and also distills the knowledge beneath unlabeled data from the teacher model with multiple graph consistency constraints for improving the model performance under complex scenarios. To the best of our knowledge, this is the first study to explore efficient and generalizable multi-type PPI prediction. Precisely, the main contributions of the work can be summarized as follows:

- For multi-type PPI prediction, we first investigate the limitations and challenges of existing methods under complex but realistic scenarios, and then propose an effective **Self-ensembling multi-Graph Neural Network-based PPI prediction (SemiGNN-PPI)** framework for improving model efficiency and generalization.
- In SemiGNN-PPI, we construct multiple graphs to learn correlations between proteins and label dependencies simultaneously. We further advance GNN with Mean Teacher to effectively utilize unlabeled data by consistency regularization with multiple constraints.
- Extensive experiments on three PPI datasets with different settings demonstrate that SemiGNN-PPI outperforms other state-of-the-art methods for multi-label PPI prediction under various challenging scenarios.

## 2 Related Work

**Protein-Protein Interaction Prediction.** Amino acid sequence-based methods have received considerable attention in PPI prediction. Early works leverage machine learning (ML) techniques [Browne *et al.*, 2007; Chen and Liu, 2005; Lin and Chen, 2013; Guo *et al.*, 2008; Wong *et al.*, 2015] to map pairs of handcrafted sequence features of proteins to interaction types. With the advent of deep learning (DL), more recent works have utilized deep neural networks [Sun *et al.*, 2017; Hashemifar *et al.*, 2018; Du *et al.*, 2017; Chen *et al.*, 2019a; Lv *et al.*, 2021] to automatically extract features from protein sequences for enhancing feature representation. Furthermore, the latest works consider protein correlations and utilize graph neural networks (GNN) to model graph-structured PPI data [Yang *et al.*, 2020; Kipf and Welling, 2016; Lv *et al.*, 2021]. However, it is essential to explore label dependencies for improving the model performance, which has long been ignored for multi-type PPI prediction. Moreover, the generalization and efficiency problems for PPI prediction are still under-explored under complex scenarios, such as data scarcity and distribution shift.

**Multi-Label Learning.** MLL addresses the problem of assigning multiple labels to a single instance. It has been utilized successfully in numerous fields, *e.g.*, computer vision [Liu *et al.*, 2021; Xu *et al.*, 2022]. Traditional MLL methods typically train independent classifiers for all labels but fail to consider the potential label interdependence, leading to suboptimal performance. Recent trends in MLL incorporate deep learning to capture the label dependencies [Wang *et al.*, 2016; Guo *et al.*, 2019]. For example, CNN-RNN [Wang *et al.*, 2016] leverages recurrent neural networks (RNNs) to transform the label vectors into an embedded space to learn label correlations implicitly. More recently, graph-based MLL methods have aroused great attention from researchers [Chen *et al.*, 2019b; Wang *et al.*, 2020]. Especially, ML-GCN [Chen *et al.*, 2019b] successfully applies Graph Convolutional Network (GCN) by constructing a directed graph over object labels to explicitly model the label dependencies adaptively. In this regard, we propose to explore correlations between PPI types with GCN on the structured label graph for more accurate PPI prediction.

**Learning from Imperfect Data.** In recent years, deep learning has made tremendous progress in numerous domains, *e.g.*, computer vision and bioinformatics. However, the applicability of deep learning is limited by heavy reliance on training data. We rarely have a perfect dataset for model training [Tajbakhsh *et al.*, 2020; Bekker and Davis, 2020], especially in biomedical imaging and bioinformatics [Zhao *et al.*, 2021; Lu *et al.*, 2022; Qu and Hickey, 2022; Pio *et al.*, 2022]. The commonly encountered challenges in PPI prediction include label scarcity, where only limited annotations are available for training (semi-supervised learning, SSL), and domain shift, where unseen data (target domain) with different distributions from training data (source domain) is used for testing (unsupervised domain adaptation, UDA). In this regard, model efficiency and generalization would be heavily constrained, limiting the wide real-world applications. Self-ensemble learning [Laine and Aila, 2017]

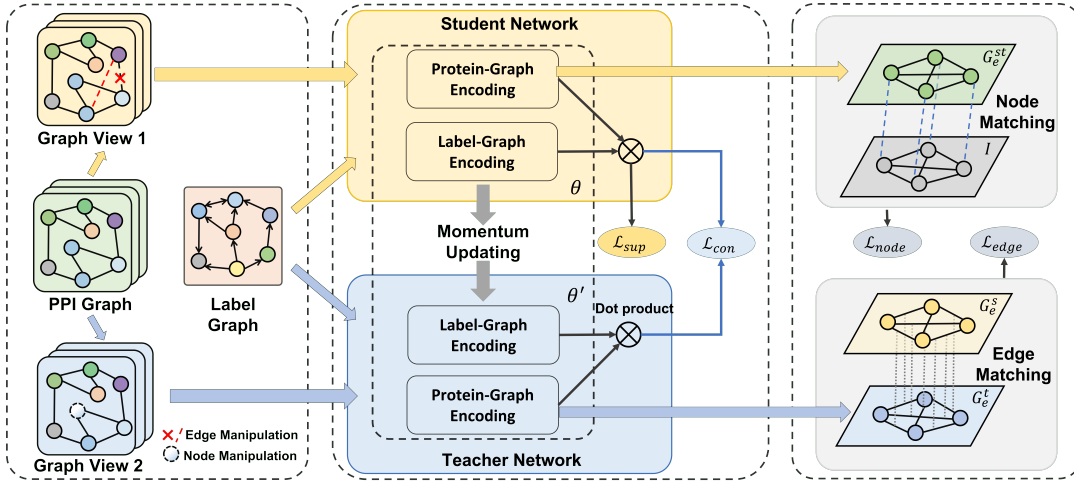


Figure 1: The overall framework of SemiGNN-PPI. First, we generate two augmented graph views with node and edge manipulations. Then, protein graphs and label graphs are fed into the multi-graph teacher-student network, which models both protein relations and label dependencies for self-ensemble learning. Simultaneously, to better capture fine-grained structural information, we align student and teacher feature embeddings by jointly optimizing multiple graph consistency constraints (node matching and edge matching).

is one of the most prevalent methods for SSL, which works by enforcing consistency in model predictions from different epochs with the network parameter average [Tarvainen and Valpola, 2017]. Recently, self-ensemble learning has been extended to visual domain adaptation tasks [Choi *et al.*, 2019; Zhang *et al.*, 2021; Zhao *et al.*, 2022], achieving promising UDA performance. Inspired by these observations, we advance GNN with self-ensemble learning to handle imperfect data for efficient and generalizable PPI prediction.

## 3 Methodology

### 3.1 Task Definition

Given a set of proteins  $P = \{p_0, p_1, \dots, p_n\}$  and a set of PPIs  $E = \{e_{ij} = \{p_i, p_j\} | i \neq j, p_i, p_j \in P, I(e_{ij}) \in \{0, 1\}\}$ , where  $I(e_{ij})$  is a binary PPI indicator function that is 1 if the PPI between proteins  $p_i$  and  $p_j$  has been confirmed, and 0 otherwise, the types of PPI can be represented by the label space  $C = \{c_0, c_1, \dots, c_{t-1}\}$  with  $t$  different types of interactions, and the labels for a confirmed PPI  $e_{ij}$  can be represented as  $y_{ij} \subseteq C$ . The goal of multi-type PPI learning is to learn a function  $f : e_{ij} \rightarrow \hat{y}_{ij}$  from the training set  $E_{train}^s$  such that for any PPI  $e_{ij} \in E_{test}^s$ ,  $\hat{y}_{ij}$  is the set of predicted labels for  $e_{ij}$ . To investigate the efficiency and generalization under complex scenarios beyond the supervised learning setting, we introduce the settings of semi-supervised learning (SSL) and unsupervised domain adaptation (UDA). In the SSL setting, the training datasets consist of limited labeled data  $E_{train}^l$  and unlabeled data  $E_{train}^u$  due to label scarcity. In the UDA setting, the model trained on  $E_{train}^s$  is tested on the unseen data  $E_{test}^t$  with different distribution.

### 3.2 Overview

Fig. 1 depicts the overview of our proposed SemiGNN-PPI framework. We first construct the multi-graph encoding (MGE) module to effectively leverage available labeled data,

which includes a protein graph encoding (PGE) network for exploring protein relations and a label graph encoding (LGE) network for learning label dependencies. To exploit knowledge from unlabeled data, we build a teacher network with the same architecture as the student network. During teacher-student training, multiple graph consistency constraints at both node and edge levels are utilized to enhance knowledge distillation for self-ensemble multi-graph learning.

### 3.3 Multi-Graph Encoding

**Protein-Graph Encoding.** Early works [Yang *et al.*, 2020; Lv *et al.*, 2021] have demonstrated the effectiveness of graph neural networks (GNNs) on PPI prediction. Considering the correlation of PPIs, we use proteins as nodes and PPIs as edges to build the PPI graph  $G = (P, E)$ . Then, the PPI prediction can be formulated from  $f(e_{ij}|G, \theta) \rightarrow \hat{y}_{ij}$  to  $f(e_{ij}|G, \theta) \rightarrow \hat{y}_{ij}$ . GNNs take the graph structure and sequence-based protein attributes as inputs to model high-level compact representation of the nodes (proteins), denoted by  $H \in \mathbb{R}^{|P| \times d}$  where  $h_p = H[p, :]$  is the latent representation of node  $p$ , and  $d$  is the dimensionality of protein features. In general, GNNs follow a recursive neighborhood aggregation scheme to iteratively update the representation of each node by aggregating and transforming the representations of its neighboring nodes. After  $l$  iterations, the transformed feature of node  $p$  can be denoted as:

$$h_p^{(l)} = \phi^{(l)}(h_p^{(l-1)}, f^{(l)}(\{h_u^{(l-1)} : u \in \mathcal{N}_k(p)\})), \quad (1)$$

where  $\mathcal{N}_k(p)$  denotes the set of  $k$ -hop neighbors of the node  $p$ ;  $f^{(l)}$  and  $\phi^{(l)}$  are an aggregation function and a combination function, respectively. Following Graph Isomorphism Network (GIN) [Xu *et al.*, 2019], we adopt the summation function to aggregate the representations of neighboring nodes and use the multi-layer perceptrons (MLPs) to update the aggregated features. Then, the update rule of the hidden node

features with a learnable parameter  $\epsilon$  in PGE is defined as:

$$h_p^{(l)} = g^l((1 + \epsilon^l) \cdot h_p^{(l-1)} + \sum_{u \in \mathcal{N}_k(p)} h_u^{(l-1)}). \quad (2)$$

**Label-Graph Encoding.** In multi-label PPI prediction, correlations exist among different types of interactions, *i.e.*, some PPI types may appear together frequently while others rarely appear together. Following [Chen *et al.*, 2019b], we model the interdependencies between different PPI types (labels) using a graph and learn inter-dependent classifiers with Graph Convolutional Network (GCN), which can be directly applied to protein features for multi-type PPI prediction. GCN aims to learn a function  $f(\cdot, \cdot)$  on the graph with  $t$  nodes. Each GCN layer can be formulated as follows:

$$h_c^{(l+1)} = f(h_c^{(l)}, A), A \in \mathbb{R}^{t \times t}, \quad (3)$$

where  $h_c^{(l+1)} \in \mathbb{R}^{t \times d'_l}$  and  $h_c^{(l)} \in \mathbb{R}^{t \times d_l}$  are the learned  $d'_l$ -dimensional node features from current layer and the  $d_l$ -dimensional node features from previous layer, respectively.  $A$  is the corresponding correlation matrix. With the convolutional operation,  $f(\cdot, \cdot)$  can be further expressed as:

$$h_c^{(l+1)} = \delta(\widehat{A}h_c^{(l)}W^l), \quad (4)$$

where  $\delta(\cdot)$  is a non-linear function set as LeakyReLU following [Chen *et al.*, 2019b],  $\widehat{A}$  is the normalized version of  $A$  and  $W^l \in \mathbb{R}^{d_l \times d'_l}$  is a transformation matrix. We leverage stacked GCNs to learn inter-dependent classifiers  $W$ . The first GCN layer takes word embeddings  $E_l \in \mathbb{R}^{t \times d_l}$  of labels and the correlation matrix  $A \in \mathbb{R}^{t \times t}$  as inputs. Considering that PPI type names are semantic, we apply the BioWordVec model [Zhang *et al.*, 2019] pretrained on the biomedical corpus for generating word embeddings  $E_l$  of each PPI type to better capture their semantics. To construct the label correlation matrix  $A$ , we compute the conditional probability of different labels within the training dataset. To avoid noises and over-smoothing, we binarize  $A$  with a threshold  $\tau$  and then re-weight it with a weight  $p$  to obtain  $\widehat{A}$ .

**Multi-Graph Based Classifier Learning.** By applying the learned classifiers  $W = \{w_i\}_{i=1}^t$  from label graph encoding (LGE) to the learned representations from protein graph encoding (PGE) for the PPI  $e_{ij}$ , we can obtain the predicted scores  $\hat{y}_{ij}$ , expressed as:

$$\hat{y}_{ij} = W(h_{p_i} \cdot h_{p_j}). \quad (5)$$

We use the traditional multi-label classification loss function to update the whole network in an end-to-end manner. The loss function can be written as:

$$\mathcal{L}_{sup} = \sum_{c=1}^t (y^c \log(\sigma(\hat{y}^c)) + (1 - y^c) \log(1 - \sigma(\hat{y}^c))),$$

where  $\sigma(\cdot)$  is the sigmoid function. Our model learns the aggregated features by combining protein neighbors and models the label correlations by learning inter-dependent classifiers simultaneously to improve the model generalization. In multi-graph learning, the learned classifiers are expected to be neighborhood aware at both feature and label levels.

### 3.4 Self-ensemble Graph Learning

To leverage unlabeled data, we adopt the mean teaching architecture for unsupervised learning, as shown in Fig. 1. We construct a teacher network  $f_t$  with the same architecture as the student network  $f_s$  based on self-ensembling [Tarvainen and Valpola, 2017]. Specifically, in each training iteration  $k$ , we update the teacher model weights  $\theta'$  with the exponential moving average (EMA) weights of the student model  $\theta$  by leveraging the momentum updating mechanism:

$$\theta'_k = m\theta'_{k-1} + (1 - m)\theta_k, \quad (6)$$

where  $m$  is momentum. During training, the student model is encouraged to be consistent with the teacher predictions for the inputs with different augmentations. Because of the non-euclidean graph structure, image augmentations such as crop and rotation cannot be directly applied to graphs. To facilitate self-ensemble graph learning, we construct two graph data augmentation methods at both the edge and node levels, *i.e.*, **Edge Manipulation** and **Node Manipulation** to augment graph topological and attribute information [You *et al.*, 2020]. **Edge Manipulation (EM):** To improve the robustness against connectivity variations, we randomly replace a certain percentage of edges in the input to the student and teacher models, since some edges (PPIs) between different nodes (proteins) may be unidentified or wrong in experimental procedures. Specifically, we follow an i.i.d. uniform distribution to randomly replace  $em_s\%$  and  $em_t\%$  of edges in the input to the student and the teacher, respectively. Different from [You *et al.*, 2020], we replace the dropped edge by linking the node with one of its neighbor's neighboring nodes for maintaining global structural information, *i.e.*, node  $p_s$  with a dropped edge connecting to node  $p_t$  could be linked to  $p_u \in \{p_u | e_{ut} = 1\}$ . **Node Manipulation (NM):** To improve the robustness against attribute missing, we randomly remove  $nm_s\%$  and  $nm_t\%$  of node features, mask them with zeros and feed them into the student and teacher models respectively, to expect the model to effectively learn the features even in the presence of missing attribute information. We construct two graph views with augmentations above to feed the student and teacher networks separately, and encourage them to generate consistent predictions using  $\ell_2$  loss:

$$\mathcal{L}_{con} = \|f_t(E_u|G, \theta'_k, \xi') - f_s(E_u|G, \theta_k, \xi)\|_2, \quad (7)$$

where  $E_u$  is unlabeled PPIs in a batch.  $\xi'$  and  $\xi$  are different augmentation operations. We randomly comprise the different augmentations in our experiments to avoid overfitting and improve model generalization.

### 3.5 Graph Consistency Constraint

The consistency regularization enforces instance-wise invariance in the prediction space towards different augmentations on the same input, describing the PPI interactions between samples. For the graph-based PPI prediction task, we also need to optimize the model in the feature space, as protein nodes in the testing set differ from the training set and PPI is performed as the relationships between proteins by feature representations extracted from neighboring proteins. Therefore, we model the fine-grained structural protein-protein relations in the feature embedding space [Ma *et al.*, 2022]. We

Method	SHS27k			SHS148k			STRING			
	Random	DFS	BFS	Random	DFS	BFS	Random	DFS	BFS	
ML	RF	78.45 <sub>0.88</sub>	35.55 <sub>2.22</sub>	37.67 <sub>1.57</sub>	82.10 <sub>0.20</sub>	43.26 <sub>3.43</sub>	38.96 <sub>1.94</sub>	88.91 <sub>0.08</sub>	70.80 <sub>0.45</sub>	55.31 <sub>1.02</sub>
	LR	71.55 <sub>0.93</sub>	48.51 <sub>1.87</sub>	43.06 <sub>5.05</sub>	67.00 <sub>0.07</sub>	51.09 <sub>2.09</sub>	47.45 <sub>1.42</sub>	67.74 <sub>0.16</sub>	61.28 <sub>0.53</sub>	50.54 <sub>2.00</sub>
DL	DPPI	73.99 <sub>5.04</sub>	46.12 <sub>3.02</sub>	41.43 <sub>0.56</sub>	77.48 <sub>1.39</sub>	52.03 <sub>1.18</sub>	52.12 <sub>8.70</sub>	94.85 <sub>0.13</sub>	66.82 <sub>0.29</sub>	56.68 <sub>1.04</sub>
	DNN-PPI	77.89 <sub>4.97</sub>	54.34 <sub>1.30</sub>	48.90 <sub>7.24</sub>	88.49 <sub>0.48</sub>	58.42 <sub>2.05</sub>	57.40 <sub>9.10</sub>	83.08 <sub>0.11</sub>	64.94 <sub>0.93</sub>	53.05 <sub>0.82</sub>
	PIPR	83.31 <sub>0.75</sub>	57.80 <sub>3.24</sub>	44.48 <sub>4.44</sub>	90.05 <sub>2.59</sub>	63.98 <sub>0.76</sub>	61.83 <sub>10.23</sub>	94.43 <sub>0.10</sub>	67.45 <sub>0.34</sub>	55.65 <sub>1.60</sub>
Graph	GNN-PPI	87.91 <sub>0.39</sub>	74.72 <sub>5.26</sub>	63.81 <sub>1.79</sub>	92.26 <sub>0.10</sub>	82.67 <sub>0.85</sub>	71.37 <sub>5.33</sub>	95.43 <sub>0.10</sub>	91.07 <sub>0.58</sub>	78.37 <sub>5.40</sub>
	GNN-PPI*	88.87 <sub>0.23</sub>	75.68 <sub>3.95</sub>	68.84 <sub>3.16</sub>	92.13 <sub>0.10</sub>	83.77 <sub>1.34</sub>	69.02 <sub>3.07</sub>	94.94 <sub>0.17</sub>	90.62 <sub>0.23</sub>	79.76 <sub>2.43</sub>
M-Graph	SemiGNN-PPI	<b>89.51</b> <sub>0.46</sub>	<b>78.32</b> <sub>3.15</sub>	<b>72.15</b> <sub>2.87</sub>	<b>92.40</b> <sub>0.22</sub>	<b>85.45</b> <sub>1.17</sub>	<b>71.78</b> <sub>3.56</sub>	<b>95.57</b> <sub>0.08</sub>	<b>91.23</b> <sub>0.26</sub>	<b>80.84</b> <sub>2.05</sub>

Table 1: Performance of SemiGNN-PPI and baseline methods over different datasets and data partition schemes. GNN-PPI: reported results in the original paper. GNN-PPI\*: reproduced GNN-PPI results. The scores are presented in the format of  $\text{mean}_{\text{std}}$ .

denote the features extracted from protein-graph encoding as  $z_s$  and  $z_t$  for the student and teacher networks, respectively. **Edge matching:** We construct the student embedding graph  $G_e^s$  and the teacher embedding graph  $G_e^t$  by calculating all pairwise Pearson’s correlation coefficient (PCC) between nodes in the same batch. Then, we enforce the student network to encode consistent instance-wise correlations with the teacher network in the embedding feature space by applying the edge matching loss:

$$\mathcal{L}_{\text{edge}} = \|\text{Adj}(G_e^s) - \text{Adj}(G_e^t)\|_2, \quad (8)$$

where Adj refers to the adjacency matrix. **Node matching:** We further formulate the edge embedding graph  $G_e^{st}$  by calculating all pairwise PCC between student encoding  $z_s$  and teacher encoding  $z_t$  in the same batch. To explicitly align encoding of the same protein from the teacher and the student network, we design a node matching loss:

$$\mathcal{L}_{\text{node}} = \|\text{diag}(\text{Adj}(G_e^{st})) - \text{diag}(I)\|_2, \quad (9)$$

where diag is an operator to create a block-diagonal matrix with the off-diagonal elements of 0, and  $I$  refers to the identity matrix. In this regard, we jointly leverage labeled and unlabeled data with graph learning in both protein and label spaces and consistency regularization in both prediction and feature spaces for PPI prediction. The overall objective function is defined as:

$$\mathcal{L} = \mathcal{L}_{\text{sup}} + \lambda_{\text{con}}\mathcal{L}_{\text{con}} + \lambda_{\text{edge}}\mathcal{L}_{\text{edge}} + \lambda_{\text{node}}\mathcal{L}_{\text{node}}, \quad (10)$$

where  $\lambda_{\text{con}}$ ,  $\lambda_{\text{edge}}$  and  $\lambda_{\text{node}}$  are scaling factors for  $\mathcal{L}_{\text{con}}$ ,  $\mathcal{L}_{\text{edge}}$  and  $\mathcal{L}_{\text{node}}$ , respectively.

## 4 Experiment

### 4.1 Dataset

We perform extensive experiments on three datasets, *i.e.*, STRING, SHS148k, and SHS27k. First, we use the multi-label PPI data of Homo sapiens from the STRING database [Szklarczyk *et al.*, 2019] for training and evaluation, including 15,355 proteins and 593,397 PPIs. The PPIs are annotated with 7 types, *i.e.*, Activation, Binding, Catalysis, Expression, Inhibition, Post-translational modification (Ptm), and Reaction. Each PPI is labeled with at least one of them. Moreover, we use two subsets of Homo sapiens PPIs from STRING, *i.e.*, SHS27k, and SHS148k [Chen *et al.*, 2019a], to further validate the proposed approach. SHS27k contains 1,690 proteins and 7,624 PPIs, while SHS148k contains 5,189 proteins and 44,488 PPIs.

### 4.2 Experimental Details

**Experimental Settings.** We follow partition algorithms in GNN-PPI [Lv *et al.*, 2021], including random, breath-first search (BFS), and depth-first search (DFS) to split the trainsets and testsets. For in-depth analysis, PPIs in the testset can be divided into **BS subset** (both proteins of the PPI are present in the labeled trainset), **ES subset** (either one protein of the PPI is present in the labeled trainset), and **NS subset** (neither of the proteins is present in the labeled trainset). The BFS and DFS partition schemes create more challenging paradigms than the random partitioning by including more ES and NS proteins in the testsets for the inter-novel protein interactions [Lv *et al.*, 2021]. In fully supervised experiments, we select 20% of the whole dataset for testing using the partition schemes mentioned above and use the rest for training. To simulate the label scarcity scenario, we randomly select 5%, 10%, and 20% samples from the trainset as the labeled data while keeping the rest as the unlabeled data. To assess the generalization capacity of our method, we evaluate our method trained with one dataset on another dataset, *i.e.*, a trainset-heterologous testset.

**Evaluation Metrics.** We use the F1 score to evaluate the model performance for multi-label PPI prediction. The score is micro-averaged over all 7 classes. The means and variances of F1 scores over three repeated experiments are reported as results, formatted as  $\text{mean}_{\text{std}}$ .

**Model Training.** 1) *Base train:* We follow GNN-PPI [Lv *et al.*, 2021] for protein-independent encoding to extract protein features from protein sequences as inputs to our framework. We initialize the multi-graph encoding network using the labeled data for 300 epochs with an initial learning rate of 0.001 and the Adam optimizer. 2) *Joint train:* Then, we train the self-ensemble graph learning framework on both labeled and unlabeled trainsets for 300 epochs. For label graph construction, we select the binarization threshold  $\tau = 0.05$  and the re-weighting factor  $p = 0.25$ . We randomly comprise the different manipulations in our experiments to avoid overfitting and improve model generalization during joint training. For manipulation ratios, we use higher ratios for the student inputs so that the student can better distill knowledge from the teacher during self-ensemble learning. More specifically, the edge manipulation ratios  $em_s\%$  and  $em_t\%$  are fixed at 10% and 5%, respectively. The node manipulation rates  $nm_s\%$  and  $nm_t\%$  are set to 10% and 5%, respectively. To scale the components of the loss function, we set the value of  $\lambda_{\text{con}}$ ,  $\lambda_{\text{edge}}$  and  $\lambda_{\text{node}}$  as 0.02, 0.01 and 0.003, respectively. More details are shown in Supplementary Material.

Method	STRING				SHS148k				SHS27k			
	5%	10%	20%	100%	5%	10%	20%	100%	5%	10%	20%	100%
Partition Scheme = Random												
GNN-PPI	89.94 <sub>0.29</sub>	92.38 <sub>0.51</sub>	93.30 <sub>0.56</sub>	94.94 <sub>0.17</sub>	79.19 <sub>0.67</sub>	82.86 <sub>0.49</sub>	86.67 <sub>0.22</sub>	92.13 <sub>0.10</sub>	52.04 <sub>3.32</sub>	60.28 <sub>12.26</sub>	79.44 <sub>1.19</sub>	88.87 <sub>0.23</sub>
Ours	<b>90.55</b> <sub>0.10</sub>	<b>92.66</b> <sub>0.59</sub>	<b>93.90</b> <sub>0.41</sub>	<b>95.57</b> <sub>0.08</sub>	<b>79.50</b> <sub>0.31</sub>	<b>83.48</b> <sub>0.30</sub>	<b>87.38</b> <sub>0.24</sub>	<b>92.40</b> <sub>0.22</sub>	<b>57.97</b> <sub>1.13</sub>	<b>62.67</b> <sub>11.26</sub>	<b>81.01</b> <sub>0.47</sub>	<b>89.51</b> <sub>0.46</sub>
Partition Scheme = DFS												
GNN-PPI	86.60 <sub>0.37</sub>	87.91 <sub>0.30</sub>	89.42 <sub>0.46</sub>	90.62 <sub>0.23</sub>	68.77 <sub>11.20</sub>	78.36 <sub>2.23</sub>	80.96 <sub>1.61</sub>	83.77 <sub>1.34</sub>	53.41 <sub>1.64</sub>	58.43 <sub>2.27</sub>	65.73 <sub>4.18</sub>	75.68 <sub>3.95</sub>
Ours	<b>87.54</b> <sub>0.06</sub>	<b>88.98</b> <sub>0.26</sub>	<b>90.23</b> <sub>0.12</sub>	<b>91.23</b> <sub>0.26</sub>	<b>69.94</b> <sub>9.57</sub>	<b>81.12</b> <sub>0.98</sub>	<b>83.63</b> <sub>0.86</sub>	<b>85.45</b> <sub>1.17</sub>	<b>58.48</b> <sub>1.11</sub>	<b>61.18</b> <sub>1.98</sub>	<b>70.31</b> <sub>2.38</sub>	<b>78.32</b> <sub>3.15</sub>
Partition Scheme = BFS												
GNN-PPI	71.35 <sub>4.67</sub>	74.94 <sub>2.35</sub>	79.99 <sub>2.75</sub>	79.76 <sub>2.43</sub>	61.42 <sub>3.29</sub>	62.51 <sub>3.07</sub>	67.10 <sub>3.48</sub>	69.02 <sub>3.07</sub>	57.93 <sub>4.11</sub>	56.84 <sub>12.19</sub>	61.18 <sub>6.58</sub>	68.84 <sub>3.16</sub>
Ours	<b>73.35</b> <sub>4.90</sub>	<b>76.94</b> <sub>2.53</sub>	<b>81.39</b> <sub>2.44</sub>	<b>80.84</b> <sub>2.05</sub>	<b>64.86</b> <sub>2.97</sub>	<b>68.76</b> <sub>1.62</sub>	<b>71.06</b> <sub>3.35</sub>	<b>71.78</b> <sub>3.56</sub>	<b>60.15</b> <sub>2.09</sub>	<b>66.13</b> <sub>2.01</sub>	<b>67.69</b> <sub>8.47</sub>	<b>72.15</b> <sub>2.87</sub>

 Table 2: Performance comparison of different methods under different label ratios. The scores are presented in the format of mean<sub>std</sub>.

Method	% Labels	Random Partition			DFS Partition		BFS Partition	
		BS (92.66%)	ES (6.95%)	NS (0.39%)	ES (75.95%)	NS (24.05%)	ES (85.70%)	NS (14.30%)
GNN-PPI	100	89.17	72.44	50.00	77.81	63.44	71.03	44.80
		<b>89.68</b>	<b>72.93</b>	50.00	<b>81.75</b>	<b>66.32</b>	<b>75.14</b>	<b>57.00</b>
SemiGNN-PPI	100	89.17	72.44	50.00	77.81	63.44	71.03	44.80
		<b>89.68</b>	<b>72.93</b>	50.00	<b>81.75</b>	<b>66.32</b>	<b>75.14</b>	<b>57.00</b>
GNN-PPI	20	BS (73.18%)	ES (24.98%)	NS (1.84%)	ES (72.87%)	NS (27.13%)	ES (47.71%)	NS (52.29%)
		83.46	70.10	43.68	64.40	54.21	<b>59.04</b>	66.33
SemiGNN-PPI	20	BS (73.18%)	ES (24.98%)	NS (1.84%)	ES (72.87%)	NS (27.13%)	ES (47.71%)	NS (52.29%)
		<b>84.09</b>	<b>71.95</b>	<b>45.78</b>	<b>73.30</b>	<b>55.46</b>	58.10	<b>73.82</b>
GNN-PPI	10	BS (55.80%)	ES (38.03%)	NS (6.16%)	ES (63.36%)	NS (36.64%)	ES (41.14%)	NS (58.86%)
		79.64	69.64	38.41	56.13	53.85	36.02	47.89
SemiGNN-PPI	10	BS (55.80%)	ES (38.03%)	NS (6.16%)	ES (63.36%)	NS (36.64%)	ES (41.14%)	NS (58.86%)
		<b>80.22</b>	<b>70.33</b>	<b>41.67</b>	<b>61.07</b>	<b>57.90</b>	<b>57.39</b>	<b>72.73</b>
GNN-PPI	5	BS (38.16%)	ES (47.61%)	NS (14.23%)	ES (46.63%)	NS (53.37%)	ES (43.18%)	NS (56.82%)
		53.43	44.33	40.64	53.85	49.62	56.10	51.95
SemiGNN-PPI	5	BS (38.16%)	ES (47.61%)	NS (14.23%)	ES (46.63%)	NS (53.37%)	ES (43.18%)	NS (56.82%)
		<b>59.76</b>	<b>57.82</b>	<b>42.71</b>	<b>58.25</b>	<b>56.25</b>	<b>58.18</b>	<b>58.60</b>

Table 3: Analysis on performance between GNN-PPI and SemiGNN-PPI over BS, ES, and NS subsets in the SHS27k dataset. The ratios of the subsets are annotated in brackets. The BS subsets are empty under DFS and BFS partitions and are omitted for brevity.

**Baseline Methods.** We compare SemiGNN-PPI with several representative methods in PPI prediction, including: **Machine Learning (ML)** methods include RF [Wong *et al.*, 2015] and LR [Silberberg *et al.*, 2014], which take commonly handcrafted protein features including AC [Guo *et al.*, 2008] and CTD [Du *et al.*, 2017] as inputs. **Deep Learning (DL)** approaches include DNN-PPI [Li *et al.*, 2018], PIPR [Chen *et al.*, 2019a], and GNN-PPI [Lv *et al.*, 2021], which take amino acid sequence-based features as inputs (More details are illustrated in the Supplementary Material). It is noted that GNN-PPI adopts graph learning to leverage protein correlations, achieving state-of-the-art performance on multi-type PPI prediction. In this regard, we extensively compare our method with GNN-PPI in different scenarios and settings.

### 4.3 Results and Analysis

**Benchmark Analysis.** In Table 1, we compare our methods with other baseline methods under different partition schemes and various datasets. It is observed that graph-based methods, *i.e.*, GNN-PPI and SemiGNN-PPI outperform other ML and DL methods, even under more challenging BFS and DFS partitions with more unseen proteins. It can be attributed to graph learning, which can better capture correlations between proteins despite the existence of more unknown proteins. Furthermore, our method incorporates multiple graphs (M-Graph) for feature learning, achieving state-of-the-art performance in multi-type PPI prediction. Especially, under challenging evaluations with small datasets, *e.g.*, SHS27k-DFS, our method achieves much higher F1 scores than GNN-PPI, since self-ensemble graph learning can effectively improve the model robustness against complex scenarios. Moreover, the number of parameters is 1.09M (GNN-PPI) and 1.13M

(ours), and the inference time on SHS27k is 0.050s (GNN-PPI) and 0.058s (ours). GNN-PPI and our method have comparable performance in the two metrics, showing the scalability of the proposed method.

**Label Efficiency.** To demonstrate the feasibility of our method under the label scarcity scenario, we present experimental results under different label ratios in Table 2. We can see that GNN-PPI receives severe performance degradation with fewer labels. In comparison, our method achieves better performance under all scenarios with different datasets, label ratios, and partition schemes. Remarkably, our method under some scenarios, *e.g.*, SHS148k-BFS-20% can achieve comparable performance with GNN-PPI using 100% labeled data, indicating the annotation efficiency of our method. To further analyze the model performance on inter-novel-protein interaction prediction, we make an in-depth performance comparison between GNN-PPI and SemiGNN-PPI in different subsets (BS/ES/NS) of the testset. As shown in Table 3, the BS subset comprises most of the whole testset under the random partition, which cannot reflect the prediction performance on the inter-novel-protein interactions. In contrast, the proportions of ES and NS subsets increase under label scarcity and other partition schemes; in these settings, SemiGNN-PPI consistently outperforms GNN-PPI in both ES and NS subsets by a large margin, which demonstrates the effectiveness of SemiGNN-PPI for inter-novel-protein interaction prediction.

**Performance on Different PPI Types.** To study the per-class prediction performance, we present the model performance on different PPI types with corresponding type ratios in Table 4. It is observed that the PPI types are unbalanced with some under-represented types, such as Ptmold, Inhibition, and Expression. Nevertheless, SemiGNN-PPI outper-

PPI Type	Type Ratio	Random Partition		DFS Partition		BFS Partition	
		GNN-PPI	SemiGNN-PPI	GNN-PPI	SemiGNN-PPI	GNN-PPI	SemiGNN-PPI
Reaction	40.61%	89.58 <sub>0.15</sub>	<b>90.16</b> <sub>0.43</sub>	81.90 <sub>1.65</sub>	<b>85.86</b> <sub>0.71</sub>	61.62 <sub>1.29</sub>	<b>64.92</b> <sub>5.73</sub>
Binding	52.71%	88.28 <sub>0.48</sub>	<b>89.46</b> <sub>0.57</sub>	83.52 <sub>1.41</sub>	<b>86.39</b> <sub>0.67</sub>	70.00 <sub>4.10</sub>	<b>72.43</b> <sub>6.33</sub>
Ptmod	20.99%	87.04 <sub>0.29</sub>	<b>87.42</b> <sub>0.33</sub>	77.94 <sub>1.67</sub>	<b>82.99</b> <sub>1.44</sub>	65.92 <sub>5.52</sub>	<b>71.32</b> <sub>5.04</sub>
Activation	42.51%	85.15 <sub>0.38</sub>	<b>85.26</b> <sub>0.46</sub>	73.48 <sub>2.74</sub>	<b>77.95</b> <sub>1.19</sub>	67.44 <sub>8.43</sub>	<b>68.04</b> <sub>8.06</sub>
Inhibition	20.20%	87.21 <sub>0.18</sub>	<b>88.09</b> <sub>0.31</sub>	72.46 <sub>1.11</sub>	<b>78.12</b> <sub>2.62</sub>	60.20 <sub>4.62</sub>	<b>67.71</b> <sub>7.21</sub>
Catalysis	44.67%	89.36 <sub>0.44</sub>	<b>90.35</b> <sub>0.31</sub>	82.30 <sub>0.80</sub>	<b>85.77</b> <sub>1.29</sub>	65.70 <sub>4.42</sub>	<b>73.39</b> <sub>6.33</sub>
Expression	7.69%	<b>47.85</b> <sub>0.79</sub>	46.99 <sub>0.22</sub>	<b>34.96</b> <sub>3.74</sub>	32.45 <sub>5.96</sub>	<b>31.81</b> <sub>6.87</sub>	28.99 <sub>4.90</sub>
Macro-Average	-	82.07 <sub>0.39</sub>	<b>82.53</b> <sub>0.38</sub>	72.37 <sub>1.87</sub>	<b>74.16</b> <sub>2.09</sub>	60.38 <sub>5.03</sub>	<b>63.29</b> <sub>5.29</sub>
Micro-Average	-	86.67 <sub>0.22</sub>	<b>87.38</b> <sub>0.24</sub>	80.96 <sub>1.61</sub>	<b>83.63</b> <sub>0.86</sub>	67.10 <sub>3.48</sub>	<b>71.06</b> <sub>3.35</sub>

Table 4: Per-class results in the SHS148k dataset with 20% training labels. The type ratios are calculated over the whole dataset.

forms GNN-PPI on most PPI types, especially for relatively imbalanced types (82.99 vs. 77.94 in Ptmod-DFS and 67.71 vs. 60.20 in Inhibition-BFS). It is noted that lower performance is achieved with our method in the type Expression, which could be due to inaccurate label correlations captured with extremely low co-occurrence with other labels, which is still a direction to explore in our future work.

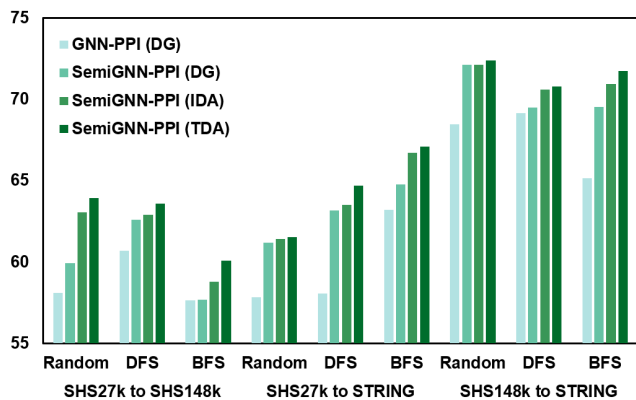


Figure 2: Performance comparison on trainset-heterologous testsets. DG: domain generalization. IDA: inductive domain adaptation. TDA: transductive domain adaptation.

**Model Generalization.** To access the generalization capability of the proposed method, we test the model trained using small datasets, *e.g.*, SHS27k on big datasets, *e.g.*, STRING in three evaluation settings: 1) *Domain Generalization (DG)*: The model is directly tested on the unseen dataset. 2) *Inductive Domain Adaptation (IDA)*: The model has access to unlabeled training data in the trainset-heterologous dataset during training. 3) *Transductive domain adaptation*: The model has access to the whole unlabeled trainset-heterologous dataset during training. In Fig. 2, we can observe that our method outperforms GNN-PPI in all partition schemes when tested on unseen datasets. Moreover, our model can effectively leverage unlabeled data, achieving better adaptation performance in both inductive and transductive setups.

**Ablation Study.** We investigate the effectiveness of different components in SemiGNN-PPI in Fig. 3. We can see that all components, *i.e.*, label graph encoding (LGE), self-ensemble (SE), and graph consistency constraint (GCC) positively contribute to the performance improvements. It is noted that too few labels *e.g.*, 10% may influence model initialization, limiting self-ensemble graph learning, while the per-

formance gains are more evenly distributed among various components with 10% or more labeled data. Particularly, the proposed GCC can further enhance the results from the self-ensemble by providing stronger regularization in the feature space. Moreover, we have performed one experiment for each augmentation strategy (F1-score) on SHS27k (20% labeled) under random partition, *i.e.*, random edge dropout (78.92), random node dropout (79.04), centrality-based [Tang *et al.*, 2015] node and edge manipulation (81.08), and ours (81.11), which show that our strategy is comparable with centrality-based manipulation and outperforms others.

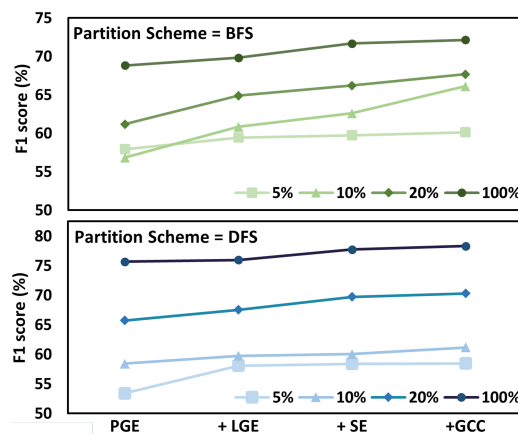


Figure 3: Results of ablation studies on different components of SemiGNN-PPI using the SHS27k dataset.

## 5 Conclusion

In this paper, we propose a novel self-ensembling multi-graph neural network (SemiGNN-PPI) for efficient and generalizable multi-type PPI prediction, which models both protein correlations and label dependencies by constructing and processing graphs at protein and label levels. To leverage unlabeled PPI data, we integrate GNN into Mean Teacher for self-ensemble graph learning, in which multiple graph consistency constraints are designed to align the teacher and student graphs in the feature embedding space for optimized consistency regularization. Extensive experiments have demonstrated the superiority in model performance, label efficiency and generalization ability of SemiGNN-PPI over state-of-the-art methods by large margins.

## Acknowledgements

The first two authors contributed equally to this work. This research was funded by Competitive Research Programme “NRF-CRP22-2019-0003”, National Research Foundation Singapore, and partially supported by A\*STAR core funding.

## References

- [Acuner Ozbabacan *et al.*, 2011] Saliha Acuner Ozbabacan, Hatice Billur Engin, Attila Gursoy, and Ozlem Keskin. Transient protein–protein interactions. *Protein engineering, design and selection*, 24(9):635–648, 2011.
- [Bekker and Davis, 2020] Jessa Bekker and Jesse Davis. Learning from positive and unlabeled data: A survey. *Machine Learning*, 2020.
- [Browne *et al.*, 2007] Fiona Browne, Haiying Wang, Hui Zheng, and Francisco Azuaje. Supervised statistical and machine learning approaches to inferring pairwise and module-based protein interaction networks. In *IEEE International Symposium on Bioinformatics and BioEngineering*, 2007.
- [Chen and Liu, 2005] Xue-Wen Chen and Mei Liu. Prediction of protein-protein interactions using random decision forest framework. *Bioinformatics*, 2005.
- [Chen *et al.*, 2019a] Muhao Chen, Chelsea J-T Ju, Guangyu Zhou, Xuelu Chen, Tianran Zhang, Kai-Wei Chang, Carlo Zaniolo, and Wei Wang. Multifaceted protein–protein interaction prediction based on siamese residual rcnn. *Bioinformatics*, 35(14):i305–i314, 2019.
- [Chen *et al.*, 2019b] Zhao-Min Chen, Xiu-Shen Wei, Peng Wang, and Yanwen Guo. Multi-label image recognition with graph convolutional networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5177–5186, 2019.
- [Choi *et al.*, 2019] Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-ensembling with gan-based data augmentation for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6830–6840, 2019.
- [Du *et al.*, 2017] Xiuquan Du, Shiwei Sun, Changlin Hu, Yu Yao, Yuanting Yan, and Yanping Zhang. Deepppi: boosting prediction of protein–protein interactions with deep neural networks. *Journal of chemical information and modeling*, 57(6):1499–1510, 2017.
- [Fields and Song, 1989] Stanley Fields and Ok-kyu Song. A novel genetic system to detect protein–protein interactions. *Nature*, 340(6230):245–246, 1989.
- [Guo *et al.*, 2008] Yanzhi Guo, Lezheng Yu, Zhining Wen, and Menglong Li. Using support vector machine combined with auto covariance to predict protein–protein interactions from protein sequences. *Nucleic acids research*, 2008.
- [Guo *et al.*, 2019] Chuan Guo, Ali Mousavi, Xiang Wu, Daniel N Holtmann-Rice, Satyen Kale, Sashank Reddi, and Sanjiv Kumar. Breaking the glass ceiling for embedding-based classifiers for large output spaces. *Advances in Neural Information Processing Systems*, 2019.
- [Hashemifar *et al.*, 2018] Somaye Hashemifar, Behnam Neyshabur, Aly A Khan, and Jinbo Xu. Predicting protein–protein interactions through sequence-based deep learning. *Bioinformatics*, 34(17):i802–i810, 2018.
- [Ho *et al.*, 2002] Yuen Ho, Albrecht Gruhler, Adrian Heilbut, Gary D Bader, Lynda Moore, Sally-Lin Adams, Anna Millar, Paul Taylor, Keiryn Bennett, Kelly Boutilier, et al. Systematic identification of protein complexes in *saccharomyces cerevisiae* by mass spectrometry. *Nature*, 415(6868):180–183, 2002.
- [Kipf and Welling, 2016] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2016.
- [Laine and Aila, 2017] Samuli Laine and Timo Aila. Temporal ensembling for semi-supervised learning. In *International Conference on Learning Representations*, 2017.
- [Li *et al.*, 2018] Hang Li, Xiu-Jun Gong, Hua Yu, and Chang Zhou. Deep neural network based predictions of protein interactions using primary sequences. *Molecules*, 23(8):1923, 2018.
- [Lin and Chen, 2013] Xiaotong Lin and Xue-wen Chen. Heterogeneous data integration by tree-augmented naïve bayes for protein-protein interactions prediction. *Proteomics*, 13(2):261–268, 2013.
- [Liu *et al.*, 2021] Weiwei Liu, Haobo Wang, Xiaobo Shen, and Ivor W Tsang. The emerging trends of multi-label learning. *IEEE transactions on pattern analysis and machine intelligence*, 44(11):7955–7974, 2021.
- [Lu *et al.*, 2022] Yingzhou Lu, Chiung-Ting Wu, Sarah J Parker, Zuolin Cheng, Georgia Saylor, Jennifer E Van Eyk, Guoqiang Yu, Robert Clarke, David M Herrington, and Yue Wang. Cot: an efficient and accurate method for detecting marker genes among many subtypes. *Bioinformatics Advances*, 2022.
- [Luo *et al.*, 2015] Xin Luo, Zhuhong You, Mengchu Zhou, Shuai Li, Hareton Leung, Yunni Xia, and Qingsheng Zhu. A highly efficient approach to protein interactome mapping based on collaborative filtering framework. *Scientific reports*, 5(1):1–10, 2015.
- [Lv *et al.*, 2021] Guofeng Lv, Zhiqiang Hu, Yanguang Bi, and Shaoting Zhang. Learning unknown from correlations: Graph neural network for inter-novel-protein interaction prediction. In *IJCAI International joint conference on artificial intelligence*, 2021.
- [Ma *et al.*, 2022] Yuchen Ma, Yanbei Chen, and Zeynep Akata. Distilling knowledge from self-supervised teacher by embedding graph alignment. In *33rd British Machine Vision Conference*. BMVA Press, 2022.
- [Margolin *et al.*, 2006] Adam A Margolin, Ilya Nemenman, Katia Basso, Chris Wiggins, Gustavo Stolovitzky, Riccardo Dalla Favera, and Andrea Califano. Aracne: an al-



- gorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. In *BMC bioinformatics*, volume 7, pages 1–15. Springer, 2006.
- [Min *et al.*, 2017] Seonwoo Min, Byunghan Lee, and Sungroh Yoon. Deep learning in bioinformatics. *Briefings in bioinformatics*, 18(5):851–869, 2017.
- [Petta *et al.*, 2016] Ioanna Petta, Sam Lievens, Claude Libert, Jan Tavernier, and Karolien De Bosscher. Modulation of protein–protein interactions for the development of novel therapeutics. *Molecular Therapy*, 2016.
- [Pio *et al.*, 2020] Gianvito Pio, Michelangelo Ceci, Francesca Prisciandaro, and Donato Malerba. Exploiting causality in gene network reconstruction based on graph embedding. *Machine Learning*, 2020.
- [Pio *et al.*, 2022] Gianvito Pio, Paolo Mignone, Giuseppe Magazzù, Guido Zampieri, Michelangelo Ceci, and Claudio Angione. Integrating genome-scale metabolic modelling and transfer learning for human gene regulatory network reconstruction. *Bioinformatics*, 2022.
- [Qu and Hickey, 2022] Xiaodong Qu and Timothy J. Hickey. Eeg4home: A human-in-the-loop machine learning model for eeg-based bci. In *Augmented Cognition*, 2022.
- [Silberberg *et al.*, 2014] Yael Silberberg, Martin Kupiec, and Roded Sharan. A method for predicting protein-protein interaction types. *PLoS One*, 9(3):e90904, 2014.
- [Skrabanek *et al.*, 2008] Lucy Skrabanek, Harpreet K Saini, Gary D Bader, and Anton J Enright. Computational prediction of protein–protein interactions. *Molecular biotechnology*, 38(1):1–17, 2008.
- [Soleymani *et al.*, 2022] Farzan Soleymani, Eric Paquet, Herna Viktor, Wojtek Michalowski, and Davide Spinello. Protein–protein interaction prediction with deep learning: A comprehensive review. *Computational and Structural Biotechnology Journal*, 2022.
- [Sun *et al.*, 2017] Tanlin Sun, Bo Zhou, Luhua Lai, and Jianfeng Pei. Sequence-based prediction of protein protein interaction using a deep-learning algorithm. *BMC bioinformatics*, 18(1):1–8, 2017.
- [Szklarczyk *et al.*, 2019] Damian Szklarczyk, Annika L Gable, David Lyon, Alexander Junge, Stefan Wyder, Jaime Huerta-Cepas, Milan Simonovic, Nadezhda T Doncheva, John H Morris, Peer Bork, et al. String v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic acids research*, 2019.
- [Tajbakhsh *et al.*, 2020] Nima Tajbakhsh, Laura Jeyaseelan, Qian Li, Jeffrey N Chiang, Zhihao Wu, and Xiaowei Ding. Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. *Medical Image Analysis*, 63:101693, 2020.
- [Tang *et al.*, 2015] Yu Tang, Min Li, Jianxin Wang, Yi Pan, and Fang-Xiang Wu. Cytonca: a cytoscape plugin for centrality analysis and evaluation of protein interaction networks. *Biosystems*, 127:67–72, 2015.
- [Tarvainen and Valpola, 2017] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 2017.
- [Wang *et al.*, 2016] Jiang Wang, Yi Yang, Junhua Mao, Zhiheng Huang, Chang Huang, and Wei Xu. Cnn-rnn: A unified framework for multi-label image classification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2285–2294, 2016.
- [Wang *et al.*, 2020] Ya Wang, Dongliang He, Fu Li, Xiang Long, Zhichao Zhou, Jinwen Ma, and Shilei Wen. Multi-label classification with label graph superimposing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 12265–12272, 2020.
- [Wong *et al.*, 2015] Leon Wong, Zhu-Hong You, Shuai Li, Yu-An Huang, and Gang Liu. Detection of protein-protein interactions from amino acid sequences using a rotation forest model with a novel pr-lpq descriptor. In *International Conference on Intelligent Computing*, 2015.
- [Xu *et al.*, 2019] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [Xu *et al.*, 2022] Kaixin Xu, Liyang Liu, Ziyuan Zhao, Zeng Zeng, and Veeravalli Bharadwaj. Object-aware self-supervised multi-label learning. In *IEEE International Conference on Image Processing*, 2022.
- [Yang *et al.*, 2020] Fang Yang, Kunjie Fan, Dandan Song, and Huakang Lin. Graph-based prediction of protein-protein interactions with attributed signed graph embedding. *BMC bioinformatics*, 21(1):1–16, 2020.
- [You *et al.*, 2020] Yuning You, Tianlong Chen, Yongduo Sui, Ting Chen, Zhangyang Wang, and Yang Shen. Graph contrastive learning with augmentations. *Advances in Neural Information Processing Systems*, 33:5812–5823, 2020.
- [Zhang *et al.*, 2019] Yijia Zhang, Qingyu Chen, Zhihao Yang, Hongfei Lin, and Zhiyong Lu. Biowordvec, improving biomedical word embeddings with subword information and mesh. *Scientific data*, 6(1):1–9, 2019.
- [Zhang *et al.*, 2021] Yabin Zhang, Haojian Zhang, Bin Deng, Shuai Li, Kui Jia, and Lei Zhang. Semi-supervised models are strong unsupervised domain adaptation learners. *arXiv preprint arXiv:2106.00417*, 2021.
- [Zhao *et al.*, 2021] Ziyuan Zhao, Kaixin Xu, Shumeng Li, Zeng Zeng, and Cuntai Guan. Mt-uda: Towards unsupervised cross-modality medical image segmentation with limited source labels. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 293–303. Springer, 2021.
- [Zhao *et al.*, 2022] Ziyuan Zhao, Fangcheng Zhou, Kaixin Xu, Zeng Zeng, Cuntai Guan, and S Kevin Zhou. Le-uda: Label-efficient unsupervised domain adaptation for medical image segmentation. *IEEE Transactions on Medical Imaging*, 2022.