

# Graph Neural Networks: Taxonomy, Advances and Trends

Yu Zhou<sup>1,2\*</sup>, Haixia Zheng<sup>1,2</sup> & Xin Huang<sup>1,2</sup>

<sup>1</sup>*College of Data Science, Taiyuan University of Technology, Taiyuan 030024, China;*

<sup>2</sup>*Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, Taiyuan 030024, China*

---

**Abstract** Graph neural networks provide a powerful toolkit for embedding real-world graphs into low-dimensional spaces according to specific tasks. Up to now, there have been several surveys on this topic. However, they usually lay emphasis on different angles so that the readers can not see a panorama of the graph neural networks. This survey aims to overcome this limitation, and provide a comprehensive review on the graph neural networks. First of all, we provide a novel taxonomy for the graph neural networks, and then refer to up to 400 relevant literatures to show the panorama of the graph neural networks. All of them are classified into the corresponding categories. In order to drive the graph neural networks into a new stage, we summarize four future research directions so as to overcome the facing challenges. It is expected that more and more scholars can understand and exploit the graph neural networks, and use them in their research community.

**Keywords** Graph Convolutional Neural Network, Graph Recurrent Neural Network, Graph Representation Learning, Graph Generative Model, Graph Neural Network

---

**Citation** Yu Zhou. Graph Neural Networks: Taxonomy, Advances and Trends. *Sci China Inf Sci*, for review

---

## 1 Introduction

Graph, as a complex data structure, consists of nodes (or vertices) and edges (or links). It can be used to model lots of complex systems in real world, e.g. social networks, protein-protein interaction networks, brain networks, road networks, physical interaction networks and knowledge graph etc. Thus, Analyzing the complex networks becomes an intriguing research frontier. With the rapid development of deep learning techniques, many scholars employ the deep learning architectures to tackle the graphs. Graph Neural Networks (GNNs) emerge under these circumstances. Up to now, the GNNs have evolved into a prevalent and powerful computational framework for tackling irregular data such as graphs and manifolds.

The GNNs can learn task-specific node/edge/graph representations via hierarchical iterative operators so that the traditional machine learning methods can be employed to perform graph-related learning tasks, e.g. node classification, graph classification, link prediction and clustering etc. Although the GNNs has attained substantial success over the graph-related learning tasks, they still face great challenges. Firstly, the structural complexity of graphs incurs expensive computational cost on large graphs. Secondly, perturbing the graph structure and/or initial features incurs sharp performance decay. Thirdly, the Wesfeiler-Leman (WL) graph isomorphism test impedes the performance improvement of the GNNs. At last, the blackbox work mechanism of the GNNs hinders safely deploying them to real-world applications.

---

\* Corresponding author (email: zhouyu@tyut.edu.cn)

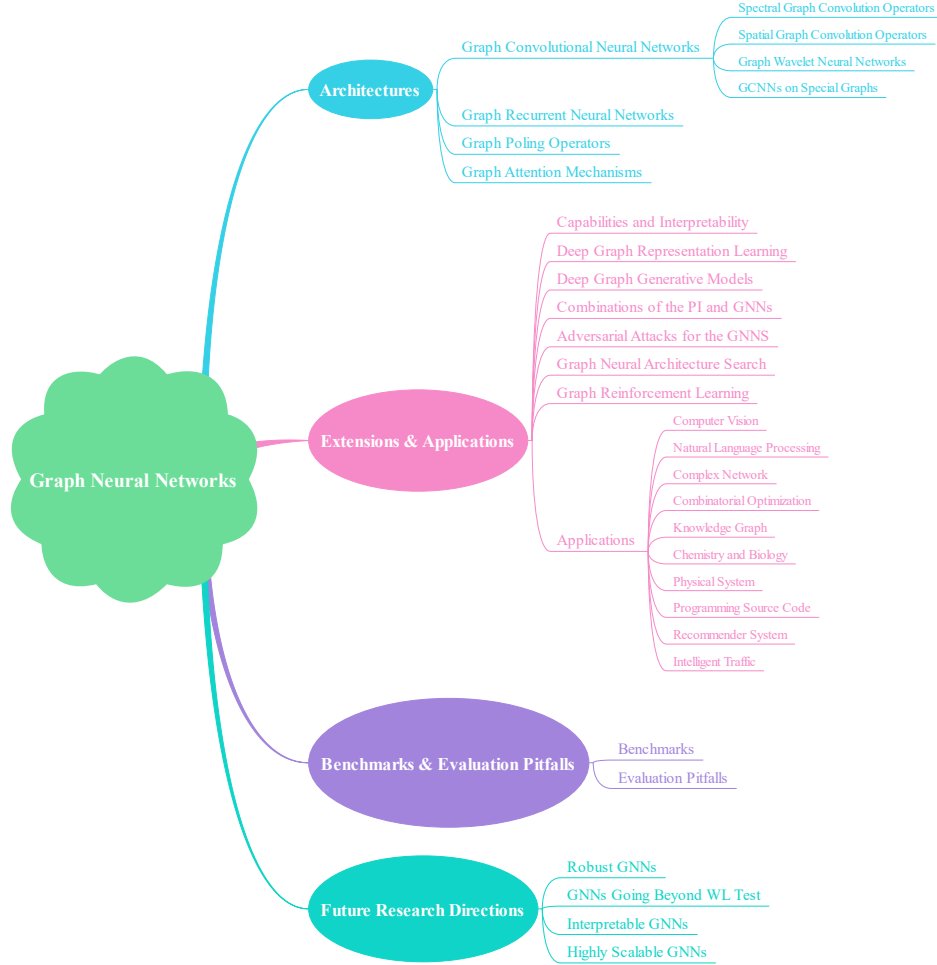
In this paper, we generalize the conventional deep architectures to the non-Euclidean domains, and summarize the architectures, extensions and applications, benchmarks and evaluation pitfalls and future research directions of the graph neural networks. Up to now, there have been several surveys on the GNNs. However, they usually discuss the GNN models from different angles and with different emphasises. To the best of our knowledge, the first survey on the GNNs was conducted by Michael M. Bronstein et al [1]. Peng Cui et al [2] reviewed different kinds of deep learning models applied to graphs from three aspects: semi-supervised learning methods including graph convolutional neural networks, unsupervised learning methods including graph auto-encoders, and recent advancements including graph recurrent neural networks and graph reinforcement learning. This survey laid emphasis on semi-supervised learning models, i.e. the spatial and spectral graph convolutional neural networks, yet comparatively less emphasis on the other two aspects. Due to the space limit, this survey only listed a few of key applications of the GNNs, but ignored the diversity of the applications. Maosong Sun et al [3] provided a detailed review of the spectral and spatial graph convolutional neural networks from three aspects: graph types, propagation step and training method, and divided its applications into three scenarios: structural scenarios, non-structural scenarios and other scenarios. However, this article did not involve the other GNN architectures such as graph auto-encoders, graph recurrent neural networks and graph generative networks. Philip S. Yu et al [4] conducted a comprehensive survey on the graph neural networks, and investigated available datasets, open-source implementations and practical applications. However, they only listed a few of core literatures on each research topic. Davide Bacciu et al [367] gives a gentle introduction to the field of deep learning for graph data. The goal of this article is to introduce the main concepts and building blocks to construct neural networks for graph data, and therefore it falls short of an exposition of recent works on graph neural networks.

It is noted that all of the aforementioned surveys do not concern capability and interpretability of GNNs, combinations of the probabilistic inference and GNNs, and adversarial attacks on graphs. In this article, we provide a panorama of GNNs for readers from 4 perspectives: architectures, extensions and applications, benchmarks and evaluations pitfalls, future research directions, as shown in Fig. 1. For the architectures of GNNs, we investigate the studies on graph convolutional neural networks (GCNNs), graph pooling operators, graph attention mechanisms and graph recurrent neural networks (GRNNs). The extensions and applications demonstrate some notable research topics on the GNNs through integrating the above architectures. Specifically, this perspective includes the capabilities and interpretability, deep graph representation learning, deep graph generative models, combinations of the Probabilistic Inference (PI) and the GNNs, adversarial attacks for GNNs, Graph Neural Architecture Search and graph reinforcement learning and applications. In summary, our article provides a complete taxonomy for GNNs, and comprehensively review the current advances and trends of the GNNs. These are our main differences from the aforementioned surveys.

**Contributions.** Our main contributions boils down to the following three-fold aspects.

1. We propose a novel taxonomy for the GNNs, which has three levels. The first includes architectures, benchmarks and evaluation pitfalls, and applications. The architectures are classified into 9 categories, the benchmarks and evaluation pitfalls into 2 categories, and the applications into 10 categories. Furthermore, the graph convolutional neural networks, as a classic GNN architecture, are again classified into 6 categories.
2. We provide a comprehensive review of the GNNs. All of the literatures fall into the corresponding categories. It is expected that the readers not only understand the panorama of the GNNs, but also comprehend the basic principles and various computation modules of the GNNs through reading this survey.
3. We summarize four future research directions for the GNNs according to the current facing challenges, most of which are not mentioned the other surveys. It is expected that the research on the GNNs can progress into a new stage by overcoming these challenges.

**Roadmap.** The remainder of this paper is organized as follows. First of all, we provide some basic



**Figure 1** The overall architecture of this paper.

notations and definitions that will be often used in the following sections. Then, we start reviewing the GNNs from 4 aspects: architectures in section 3, extensions and applications in section 4, benchmarks and evaluation pitfalls in section 5 and future research directions in section 6. Finally, we conclude our paper.

## 2 Preliminaries

In this section, we introduce relevant notations so as to conveniently describe the graph neural network models. A simple graph can be denoted by  $G = (V, E)$  where  $V$  and  $E$  respectively denote the set of  $N$  nodes (or vertices) and  $M$  edges. Without loss of generality, let  $V = \{v_1, \dots, v_N\}$  and  $E = \{e_1, \dots, e_M\}$ . Each edge  $e_j \in E$  can be denoted by  $e_j = (v_{s_j}, v_{r_j})$  where  $v_{s_j}, v_{r_j} \in V$ . Let  $A_G$  denote the adjacency matrix of  $G$  where  $A_G(s, r) = 1$  iff there is an edge between  $v_s$  and  $v_r$ . If  $G$  is edge-weighted,  $A_G(s, r)$  equals the weight value of the edge  $(v_s, v_r)$ . If  $G$  is directed,  $(v_{s_j}, v_{r_j}) \neq (v_{r_j}, v_{s_j})$  and therefore  $A_G$  is asymmetric. A directed edge  $e_j = (v_{s_j}, v_{r_j})$  is also called an arch, i.e.  $e_j = \langle v_{s_j}, v_{s_j} \rangle$ . Otherwise  $(v_{s_j}, v_{r_j}) = (v_{r_j}, v_{s_j})$  and  $A_G$  is symmetric. For a node  $v_s \in V$ , let  $N_G(v_s)$  denote the set of neighbors of  $v_s$ , and  $d_G(v_s)$  denote the degree of  $v_s$ . If  $G$  is directed, let  $N_G^+(v_s)$  and  $N_G^-(v_s)$  respectively denote the incoming and outgoing neighbors of  $v_s$ , and  $d_G^+(v_s)$  and  $d_G^-(v_s)$  respectively denote the incoming and outgoing degree of  $v_s$ . Given a vector  $a = (a_1, \dots, a_N) \in \mathbb{R}^N$ ,  $\text{diag}(a)$  (or  $\text{diag}(a_1, \dots, a_N)$ ) denotes a diagonal matrix consisting of the elements  $a_n, n = 1, \dots, N$ .

A vector  $x \in \mathbb{R}^N$  is called a 1-dimensional graph signal on  $G$ . Similarly,  $X \in \mathbb{R}^{N \times d}$  is called a  $d$ -dimensional graph signal on  $G$ . In fact,  $X$  is also called a feature matrix of nodes on  $G$ . Without loss of generality, let  $X[j, k]$  denote the  $(j, k)$ -th entry of the matrix  $X \in \mathbb{R}^{N \times d}$ ,  $X[j, :] \in \mathbb{R}^d$  denote the feature vector of the node  $v_j$  and  $X[:, j]$  denote the 1-dimensional graph signal on  $G$ . Let  $\mathbb{I}_N$  denote a  $N \times N$  identity matrix. For undirected graphs,  $L_G = D_G - A_G$  is called the Laplacian matrix of  $G$ , where  $D_G[r, r] = \sum_{c=1}^N A_G[r, c]$ . For a 1-dimensional graph signal  $x$ , its smoothness  $s(x)$  is defined as

$$s(x) = x^T L_G x = \frac{1}{2} \sum_{r,c=1}^N A_G(r, c) (x[r] - x[c])^2. \quad (1)$$

The normalization of  $L_G$  is defined by  $\bar{L}_G = \mathbb{I}_N - D_G^{-\frac{1}{2}} A_G D_G^{-\frac{1}{2}}$ .  $\bar{L}_G$  is a real symmetric semi-positive definite matrix. So, it has  $N$  ordered real non-negative eigenvalues  $\{\lambda_n : n = 1, \dots, N\}$  and corresponding orthonormal eigenvectors  $\{u_n : n = 1, \dots, N\}$ , namely  $\bar{L}_G = U \Lambda U^T$  where  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$  and  $U = (u_1, \dots, u_N)$  denotes a orthonormal matrix. Without loss of generality,  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N = \lambda_{\max}$ . The eigenvectors  $u_n, n = 1, \dots, N$  are also called the graph Fourier bases of  $G$ . Obviously, the graph Fourier basis are also the 1-dimensional graph signal on  $G$ . The graph Fourier transform [5] for a given graph signal  $x$  can be denoted by

$$\hat{x} \triangleq \mathcal{F}(x) = U^T x. \quad (2)$$

The inverse graph Fourier transform can be correspondingly denoted by

$$x \triangleq \mathcal{F}^{-1}(\hat{x}) = U \hat{x}. \quad (3)$$

Note that the eigenvalue  $\lambda_n$  actually measures the smoothness of the graph Fourier mode  $u_n$ . Throughout this paper,  $\rho(\cdot)$  is used to denote an activation function, and  $\bowtie$  denote the concatenation of at least two vectors. We somewhere use the function  $\text{Concat}(\cdot)$  to denote the concatenation of two vectors as well.

### 3 Architectures

In this section, we conduct circumstantial investigations for the GNN architectures so as to provide users with available ingredients to construct complex GNN models.

#### 3.1 Graph Convolutional Neural Networks (GCNNs)

The GCNNs play pivotal roles on tackling the irregular data (e.g. graph and manifold). The GCNNs are motivated by the Convolutional Neural Networks (CNNs) to learn hierarchical representations of irregular data. There have been some efforts to generalize the CNN to graphs [27, 31, 48]. However, they are usually computationally expensive and cannot capture spectral or spatial features. Below, we introduce the GCNNs from the next 6 aspects: spectral GCNNs, spatial GCNNs, Graph wavelet neural networks and GCNNs on special graphs.

##### 3.1.1 Spectral Graph Convolution Operators

The spectral graph convolution operator is defined via the graph Fourier transform. For two graph signals  $x$  and  $y$  on  $G$ , their spectral graph convolution  $x *_G y$  is defined by

$$\begin{aligned} x *_G y &= \mathcal{F}^{-1}(\mathcal{F}(x) \otimes \mathcal{F}(y)) \\ &= U (U^T x \otimes U^T y) \\ &= U \text{diag}(U^T y) U^T x, \end{aligned} \quad (4)$$

where  $\otimes$  denotes the element-wise Hadamard product [7, 18, 27]. The spectral graph convolution can be rewritten as

$$x *_G f_\theta = U f_\theta U^T x,$$

where  $f_\theta$  is a diagonal matrix consisting of the learnable parameters. That is, the signal  $x$  is filtered by the spectral graph filter (or graph convolution kernel)  $f_\theta$ . For a  $d^{(l)}$ -dimensional graph signal  $X^{(l)}$  on  $G$ , the output  $X^{(l+1)}$  yielded by a graph convolution layer, namely  $d^{(l+1)}$ -dimensional graph signal on  $G$ , can be written as

$$X^{(l+1)}[:, k] = \rho \left( \sum_{j=1}^{d^{(l)}} U f_{\theta,j,k}^{(l)} U^T X^{(l)}[:, j] \right), \quad (5)$$

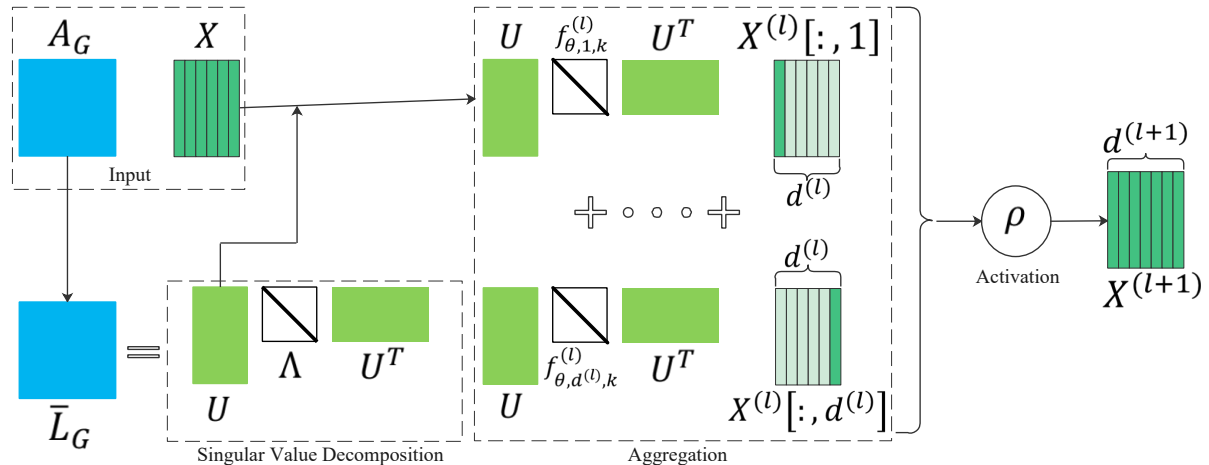
where  $f_{\theta,j,k}^{(l)}$  is a spectral graph filter, i.e. a  $N \times N$  diagonal matrix consisting of learnable parameters corresponding to the  $j$ -th graph signal at  $l$ -th layer and the  $k$ -th graph signal at  $(l+1)$ -th layer. The computational framework of the spectral GCNN in Eq. (5) is demonstrated in Fig. 2. It is worth noting that the calculation of the above graph convolution layer takes  $O(N^3)$  time and  $O(N^2)$  space to perform the eigendecomposition of  $\bar{L}_G$  especially for real large graphs. More generally, a generic graph filter  $F_{\theta,j,k}^{(l)} \in \mathbb{R}^{N \times N}$  can be employed to reformulate the graph convolution layer, i.e.

$$X^{(l+1)}[:, k] = \rho \left( \sum_{j=1}^{d^{(l)}} F_{\theta,j,k}^{(l)} X^{(l)}[:, j] \right),$$

where  $k = 1, \dots, d^{(l+1)}$ . The literature [127] gives a comprehensive analysis of when work matters by transforming different classical network structures to GCNNs and provides guidelines for choosing appropriate graph network frameworks. GraphMix [179], as a regularization technique, is inspired by interpolation based data augmentation techniques [387, 388], and is adjusted for the graph-structured data. It augments the vanilla GCNN with a Fully Connected Network (FCN) via a parameter sharing strategy. The FCN loss is computed using the manifold mixup, which is defined as follows

$$\mathcal{L} = \mathbb{E}_{(x,y) \sim \mathcal{D}} \mathbb{E}_{(x',y') \sim \mathcal{D}} \mathbb{E}_{\lambda \sim \text{Beta}(\alpha,\alpha)} l(f(\text{Mix}_\lambda(g(x), g(x'))), \text{Mix}_\lambda(y, y')),$$

where  $g(\cdot)$  is a function mapping input vectors to hidden states,  $f(\cdot)$  is a function mapping hidden states to predicted output, and  $\text{Mix}_\lambda(a, b) = \lambda \otimes a + (\mathbf{1} - \lambda) \otimes b$ . The GCNN loss is computed in the standard manner. These two losses are optimized in an alternating fashion. The manifold mixup is used in the FCN so as to facilitate learning better features.



**Figure 2** Computational framework of the spectral GCNN.

**Spectral Graph Filter.** Many studies focus on designing different spectral graph filters. Mikael Henaff et al [27] argues that smooth spectral graph filter coefficients result in spatially-localized filters and therefore use parametric filters of the form

$$f_\theta = \text{diag} \left\{ \sum_{j=1}^r \alpha_j \beta_j(\lambda_1), \dots, \sum_{j=1}^r \alpha_j \beta_j(\lambda_N) \right\},$$

where  $\beta_j(\lambda_k)$  and  $\alpha_j$  respectively denote the interpolation kernels and coefficients. In order to circumvent the eigendecomposition, the spectral graph filter  $f_\theta$  can be formulated as a  $K$ -localized polynomial of the eigenvalues of the normalized graph Laplacian  $\bar{L}_G$  [6, 8, 78], i.e.

$$f_\theta = f_\theta(\Lambda) \triangleq \sum_{k=0}^{K-1} \theta_k \Lambda^k. \quad (6)$$

In practice, the  $K$ -localized Chebyshev polynomial [78] is a favorable choice of formulating the spectral graph filter, i.e.

$$f_\theta(\Lambda) = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\Lambda}),$$

where the Chebyshev polynomial is defined as

$$\begin{aligned} T_0(x) &= 1, & T_1(x) &= x \\ T_k(x) &= 2xT_{k-1}(x) - T_{k-2}(x) \end{aligned} \quad (7)$$

and  $\tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - \mathbb{I}_N$ . The reason why  $\tilde{\Lambda} = \frac{2}{\lambda_{\max}} \Lambda - \mathbb{I}_N$  is because it can map eigenvalues  $\lambda \in [0, \lambda_{\max}]$  into  $[-1, 1]$ . This filter is  $K$ -localized in the sense that it leverages information from nodes which are at most  $K$ -hops away. In order to further decrease the computational cost, the 1st-order Chebyshev polynomial is used to define the spectral graph filter. Specifically, it lets  $\lambda_{\max} \approx 2$  (because the largest eigenvalue of  $\bar{L}_G$  is less than or equal to 2 [374]) and  $\theta = \theta_0 = -\theta_1$ . Moreover, the renormalization trick is used here to mitigate the limitations of the vanishing/exploding gradient, namely substituting  $\tilde{D}_G^{-\frac{1}{2}} \tilde{A}_G \tilde{D}_G^{-\frac{1}{2}}$  for  $\mathbb{I}_N + D_G^{-\frac{1}{2}} A_G D_G^{-\frac{1}{2}}$  where  $\tilde{A}_G = A_G + \mathbb{I}_N$  and  $\tilde{D}_G = \text{diag} \left( \sum_{k=1}^N \tilde{A}[1, k], \dots, \sum_{k=1}^N \tilde{A}[N, k] \right)$ . As a result, the Graph Convolutional Network (GCN) [6, 375] can be defined as

$$X^{(l+1)} = \rho \left( \tilde{D}_G^{-\frac{1}{2}} \tilde{A}_G \tilde{D}_G^{-\frac{1}{2}} X^{(l)} \Theta^{(l)} \right). \quad (8)$$

The literature [19] proposes a dual GCNN which incorporates local consistency and global consistency of the semi-supervised learning into an objective function. The local consistency embodied by the adjacency matrix  $A$  is defined to be  $\text{Conv}_A^{(l)}(X) = X^{(l)} = \rho \left( \tilde{D}_G^{-\frac{1}{2}} \tilde{A}_G \tilde{D}_G^{-\frac{1}{2}} X^{(l-1)} W^{(l-1)} \right)$ , whereas the global consistency embodied by the Positive Pointwise Mutual Information (PPMI) matrix  $P$  [377, 378] is defined to be  $\text{Conv}_P^{(l)}(X) = X^{(l)} = \rho \left( D_G^{-\frac{1}{2}} P D_G^{-\frac{1}{2}} X^{(l-1)} W^{(l-1)} \right)$ .

The Chebyshev spectral graph filter suffers from a drawback that the spectrum of  $\bar{L}_G$  is linearly mapped into  $[-1, 1]$ . This drawback makes it hard to specialize in the low frequency bands. In order to mitigate this problem, Michael M. Bronstein et al [12] proposes the Cayley spectral graph filter. Given that  $i$  denotes the imaginary unit, i.e.  $i^2 = -1$ , the Cayley spectral graph filter is defined to be

$$f_{c,h}(\bar{L}_G) = c_0 \mathbb{I}_N + 2\text{Re} \left( \sum_{j=0}^r c_j \mathcal{C}(h \bar{L}_G)^j \right),$$

where the order- $r$  Cayley polynomial  $f_{c,h}(\lambda) = c_0 + 2\text{Re} \left( \sum_{j=0}^r c_j \mathcal{C}(\lambda)^j \right)$  and the Cayley transform  $\mathcal{C}(\lambda) = \frac{\lambda-i}{\lambda+i}$ . The spectral zoom parameter  $h$  in  $g_{c,h}(\cdot)$  can dilate the spectrum of  $\bar{L}_G$ . In addition to

the aforementioned spectral graph filters, there are also other available filters, e.g. the adaptive graph filter [8] via metric learning (Mahalanobis distance and Gaussian kernel), the Lanczos-based multi-scale graph filter [14], the simple graph convolution filter [22], the accelerated graph filter using Lanczos [30], the cross graph convolution [37], the topology adaptive graph filter [38,39], the fractional generalized sigma-based filter [44]. MixHop [90] is a higher-order graph convolution architecture via sparsified neighborhood mixing. It replaces the graph convolution layer in Eq. (8) with  $X^{(l+1)} = \bigotimes_{p \in P} \rho \left( \tilde{A}_G^j X^{(l)} W_j^{(l)} \right)$ , where the hyper-parameter  $P$  is a set of integer adjacency powers. The literature [90] proves that the MixHop GCN, rather than the vanilla GCN, can represent two-hop delta operators, and is capable of representing general layer-wise neighborhood mixing. Some studies are motivated by the capsule network [400] to construct capsule-inspired graph neural networks [52–54]. CapsGNN [52] consists of three key blocks: basic node capsules extraction block, high level graph capsules extraction block and graph classification block. The first takes the layer-wise outputs of a GCN, and outputs primary capsules. The second apply an attention module to scale the primary capsules, and the last is used to perform the graph classification task via the dynamic routing mechanism.

**Overcoming Time and Memory Challenges.** A chief challenge for GCNNs is that their training cost is strikingly expensive, especially on huge and sparse graphs. The reason is that the GCNNs require full expansion of neighborhoods for the feed-forward computation of each node, and large memory space for storing intermediate results and outputs. There are two approaches, namely sampling and decomposition, to overcome the time and memory challenges for spectral GCNNs. Applying sampling techniques to every node is a straightforward approach to mitigating above issues [9–11, 23, 25]. The literature [9] accelerates the training of GCNs through performing an adaptive layer-wise sampling method (neighbor sampling). The literature [10] interprets the graph convolution in Eq. (8) as integral transforms of embedding functions with respect to a probability measure, and then uses importance sampling to consistently estimate the integrals. The sampling techniques can also be integrated into the mini-batch Stochastic Gradient Descent (mini-batch SGD) method [11, 23]. The literature [11] develops control variate based algorithms allowing sampling an arbitrary small neighbor size when training the GCNs via the mini-batch SGD. The literature [23] employs a graph clustering method to partition a graph into disjoint subgraphs, and then samples a cluster as a mini-batch when performing SGD. The literature [25] addresses the time and memory challenges by transforming the GCNs into regular CNNs via  $k$ -largest node selection and 1-D convolution kernels and performing mini-batch SGD via subgraph selection. The second approach to mitigating time and space costs is to decompose graphs into several small dense subgraphs [24]. Specifically, this paper employs the chordal decomposition technique to accelerate the training of GCNs on large-scale sparse networks. It contains two stages: (1) preprocessing, namely decomposing a large-scale sparse networks into several cliques and constructing a clique tree; (2) GCN training, which is performed clique by clique by traversing the clique tree.

**Depth Trap of Spectral GCNNs.** A bottleneck of GCNNs is that their performance maybe decrease with ever-increasing number of layer. This decay is often attributed to three factors: (1) overfitting resulting from the ever-increasing number of parameters; (2) gradient vanishing/explosion during training; (3) oversmoothing making vertices from different clusters more and more indistinguishable. The reason for oversmoothing is that performing the Laplacian smoothing many times forces the features of vertices within the same connected component to stuck in stationary points [16]. The Laplacian smoothing is defined as

$$\begin{aligned} X^{(l+1)} &= (1 - \gamma)X^{(l)} + \gamma \tilde{D}^{-1} \tilde{A} X^{(l)} \\ &= X^{(l)} - \gamma (\mathbb{I}_N - \tilde{D}^{-1} \tilde{A}) X^{(l)}. \end{aligned} \tag{9}$$

A feasible approach to combating the oversmoothing is to capitalize on information beyond the  $K$ -hop neighbors. The literature [33] proposes a novel architecture which incorporates the student-teacher semi-supervised framework [376] into the classic GCN. Its overall loss function is defined to be

$$\mathcal{L}(\Theta_s, \Theta_t; A, A', \mathcal{D}_L, \mathcal{D}_U) = \sum_{(x,y) \in \mathcal{D}_L} l_s + \sum_{x \in \mathcal{D}_L \cup \mathcal{D}_U} l_{st},$$

where  $\mathcal{D}_L$  and  $\mathcal{D}_U$  denote the labeled and unlabeled data respectively,  $l_s = -\sum_{c=1}^C y_c \log GCN_s(A, x; \Theta_s)_c$  denotes the cross entropy of the student network and  $l_{st} = D_{KL}(GCN_t(A, x; \Theta_t \| GCN'_s(A', x; \Theta_s))$  denotes the KL divergence between the teacher network and the perturbed student network. It is noted that this model can discover much more information from unlabeled vertices and learn from the global graph topology. The literature [373] proposes a novel, fast and easily-implementing normalization layer, named PAIRNORM, in order to mitigating the oversmoothing issues. Specifically, suppose  $X^{(L)} = GCNN^{(L)}(X, A, \Theta)$  is the output of a GCNN with  $L$  layers. The PAIRNORM can be written as

$$\begin{aligned} \widehat{X}^{(L)}[j, :] &= X^{(L)}[j, :] - \frac{1}{N} \sum_{k=1}^N X^{(L)}[k, :] && \text{(Center),} \\ \dot{X}^{(L)}[j, :] &= s \cdot \frac{\widehat{X}^{(L)}[j, :]}{\sqrt{\frac{1}{N} \sum_{k=1}^N \|\widehat{X}^{(L)}[k, :]\|_2^2}} && \text{(Scale),} \end{aligned}$$

where  $s$  is a hyper-parameter determining the total pairwise squared distance. The literature [379] proposes a novel and flexible technique, named DROPEdge, to mitigate both the oversmoothing and overfitting issues. Specifically, epoch-wise DROPEdge randomly drops out a certain rate of edges of the input graph at each training epoch, whereas layer-wise one performs the dropping-edge operations at each of of the GCN. For the gradient vanishing issues, a feasible solution to mitigate it is to add skip connections between two inconsecutive layers [398]. In practice, the highway networks [399] can be employed to control the tradeoff of how much multi-hop neighborhood information should be passed to the current node.

### 3.1.2 Spatial Graph Convolution Operators

To the best of our knowledge, the original spatial GCNN model constitutes a transition function and an update function [47, 76, 77, 380]. In order to ensure the uniqueness of states, the transition function must be a contraction map. It is noted that the universal approximation theorem for graph neural networks holds as well. That is, under mild generic conditions, most of practically useful functions on graphs can be approximated in probability by GNNs up to any prescribed degree of accuracy [77]. In the following, we firstly introduce a generic framework of the spatial GCNN, and then discuss its variants sequentially.

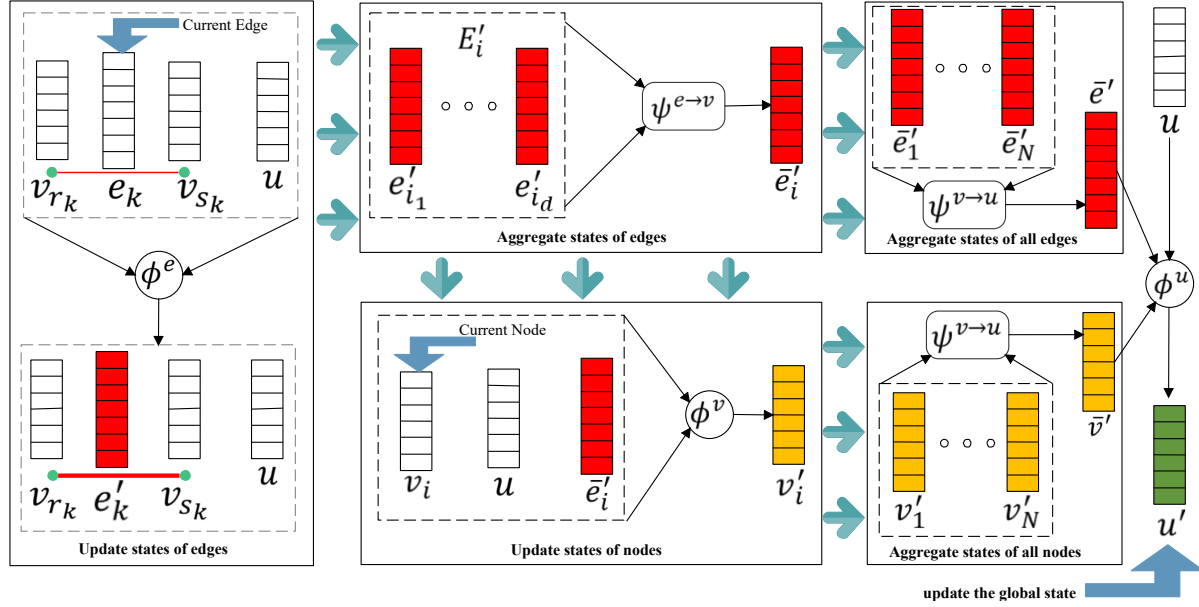
Graph networks (GNs) as generic architectures with relational inductive bias [85] provide an elegant interface for learning entities, relations and structured knowledge. Specifically, GNs are composed of GN blocks in a sequential, encode-process-decode or recurrent manner. GN blocks contain three kinds of update functions, namely  $\phi^e(\cdot), \phi^v(\cdot), \phi^u(\cdot)$ , and three kinds of aggregation functions, namely  $\psi^{e \rightarrow v}(\cdot), \psi^{e \rightarrow u}(\cdot), \psi^{v \rightarrow u}(\cdot)$ . The iterations are described as follows.

$$\begin{aligned} e'_k &= \phi^e(e_k, v_{r_k}, v_{s_k}, u) & \bar{e}'_i &= \psi^{e \rightarrow v}(E'_i), \\ v'_i &= \phi^v(v_i, \bar{e}'_i, u) & \bar{e}' &= \psi^{e \rightarrow u}(E'), \\ u' &= \phi^u(u, \bar{e}', \bar{v}') & \bar{v}' &= \psi^{v \rightarrow u}(V'), \end{aligned} \tag{10}$$

where  $e_k$  is an arch from  $v_{s_k}$  to  $v_{r_k}$ ,  $E'_i = \{(e'_k, s_k, r_k) : r_k = i, k = 1, \dots, M\}$ ,  $V' = \{v'_i : i = 1, \dots, N\}$  and  $E' = \{(e'_k, s_k, r_k) : k = 1, \dots, M\}$ , see Fig. 3. It is noted that the aggregation functions should be invariant to any permutations of nodes or edges. In practice, the GN framework can be used to implement a wide variety of architectures in accordance with three key design principles, namely flexible representations, configurable within-block structure and flexible multi-block architectures. The first principle means that GNs support highly flexible graph representations in terms of the representations of the attributes and the structure of the graph itself. The second means that the GN block can be configured in different manners. This also reflects the general adaption and universality of the GN framework. The last means that the GN blocks can be organized in different architectures, e.g. sequential architectures, encode-process-decode architectures and recurrent architectures. GNs can be employed to learn object-



and relation-centric representations of complex systems so as to help them implement relational inductive bias [100, 102]. Simon S. Du et al [49] present a new class of graph kernels, namely Graph Neural Tangent Kernels (GNTKs), based on GNs. GNTKs not only have the expressive power of GNs, but also inherit advantages of graph kernels. Below, we respectively introduce three prevalent variants of the GNs, namely Message Passing Neural Networks (MPNNs) [93], Non-local Neural Networks (NLNNs) [102] and GraphSAGE [104].



**Figure 3** Computational framework of the spatial GCNN.

**Variants of GNs—MPNNs.** MPNNs [93] have two phases, a message passing phase and a readout phase. The message passing phase is defined by a message function  $M_l$  (playing the role of the composition of the update function  $\psi^{e \rightarrow v}(\cdot)$  and the update function  $\phi^e(\cdot)$ ) and a vertex update function  $U_l$  (playing the role of the update function  $\phi^v(\cdot)$ ). Specifically,

$$\begin{aligned} m_v^{(l+1)} &= \sum_{u \in N_G(v)} M_l(x_v^{(l)}, x_u^{(l)}, e_{v,u}) \\ x_v^{(l+1)} &= U_l(x_v^{(l)}, m_v^{(l+1)}) \end{aligned},$$

where  $e_{v,u}$  denotes the feature vector of the edge with two endpoints  $v$  and  $u$ . The readout phase computes a universal feature vector for the whole graph using a readout function  $R(\cdot)$ , i.e.

$$u = R\left(\left\{x_v^{(L)} : v \in V\right\}\right).$$

The readout function  $R(\cdot)$  should be invariant to permutations of nodes. Note that a lot of GCNNs can be regarded as special forms of the MPNN such as convolutional networks for learning molecular fingerprints [103], gated graph sequence neural networks [122], interaction networks [336] and molecular graph convolutions [195].

**Variants of GNs—NLNNs.** NLNNs [102] give a general definition of non-local operations [381] which is a flexible building block and can be easily integrated into convolutional/recurrent layers. Specifically, the generic non-local operation is defined as

$$y_s = \frac{1}{\mathcal{C}(x_s)} \sum_t f(x_s, x_t) g(x_t), \quad (11)$$

where  $s$  is the index of the output position,  $t$  is the index enumerating all possible positions,  $x$  and  $y$  are input and output signal respectively,  $f(\cdot, \cdot)$  denotes the affinity between  $x_s$  and  $x_t$ , and  $\mathcal{C}(x_s) =$

$\sum_t f(x_s, x_t)$  is a normalization factor. The affinity function  $f(\cdot, \cdot)$  can be instantiated by the following forms

1. Gaussian:  $f(x_s, x_t) = e^{x_s^T x_t}$ ;
2. Embedded Gaussian:  $f(x_s, x_t) = e^{\theta(x_s)^T \eta(x_t)}$ , where  $\theta(x_s) = W_\theta x_s$  and  $\eta(x_t) = W_\eta x_t$ ;
3. Dot Product:  $f(x_s, x_t) = \theta(x_s)^T \eta(x_t)$ ;
4. Concatenation:  $f(x_s, x_t) = \text{ReLU}(w_f^T [\theta(x_s), \eta(x_t)])$ .

Finally, the non-local building block is defined as  $z_s = W_z y_s + x_s$ , where " $+x_s$ " denotes a residual connection. It is noted that  $f(\cdot, \cdot)$  and  $g(\cdot)$  play the role of  $\phi^e(e_k, v_{r_k}, v_{s_k}, u)$ , and the summation in formula (11) plays the role of  $\psi^{e \rightarrow v}(E'_i)$ .

**Variants of GNs—GraphSAGE.** GraphSAGE (SAMPLE and AGGREGATE) [104] is a general inductive framework capitalizing on node feature information to efficiently generate node embedding vectors for previously unseen nodes. Specifically, GraphSAGE is composed of an aggregation function  $\text{AGGREGATE}^{(l)}(\cdot)$  and an update function  $\text{UPDATE}^{(l)}$ , i.e.

$$\begin{aligned} x_{N_G(v)}^{(l)} &= \text{AGGREGATE}^{(l)} \left( \left\{ x_u^{(l-1)} : u \in N_G(v) \right\} \right) \\ x_v^{(l)} &= \text{UPDATE}^{(l)} \left( \left\{ x_v^{(l-1)}, x_{N_G(v)}^{(l)} \right\} \right) \end{aligned}$$

where  $N_G(v)$  denotes a fixed-size set of neighbors of  $v$  uniformly sampling from its whole neighbors so as to reduce the time and memory cost. The aggregation function can be instantiated by the following forms

1. Mean Aggregator:  $x_v^{(l)} = \sigma \left( W \cdot \text{MEAN} \left( \left\{ x_v^{(l-1)} \right\} \cup \left\{ x_u^{(l-1)} : u \in N_G(v) \right\} \right) \right)$ ;
2. LSTM Aggregator: applying the LSTM model [160] to a random permutation of neighbors of  $v$ ;
3. Pooling Aggregator:  $\text{AGGREGATE}^{(l)} = \max \left( \left\{ \sigma \left( W x_u^{(l)} + b \right) : u \in N_G(v) \right\} \right)$ .

Note that the aggregation function and update function play the role of  $\psi^{e \rightarrow v}(E'_i)$  and  $\phi^v(v_i, \tilde{e}'_i, u)$  in formula (10) respectively.

**Variants of GNs—Hyperbolic GCNNs.** The Euclidean GCNNs aim to embed nodes in a graph into a Euclidean space. This will incur a large distortion especially when embedding real-world graphs with scale-free and hierarchical structure. Hyperbolic GCNNs pave an alternative way of embedding with little distortion. The  $n$ -dimensional hyperbolic space [69, 193], denoted as  $\mathbb{H}_K^n$ , is a unique, complete, simply connected  $d$ -dimensional Riemannian manifold with constant negative sectional curvature  $-\frac{1}{K}$ ,

$$\mathbb{H}_K^n = \{x \in \mathbb{R}^{n+1} : \langle x, x \rangle_{\mathcal{M}} = -K, x_0 > 0\},$$

where the Minkowski inner produce  $\langle x, y \rangle_{\mathcal{M}} = -x_0 y_0 + \sum_{j=1}^d x_j y_j, \forall x, y \in \mathbb{R}^{n+1}$ . Its tangent space centered at point  $x$  is denoted to be  $\mathcal{T}_x \mathbb{H}_K^n = \{v \in \mathbb{R}^{n+1} : \langle x, v \rangle_{\mathcal{M}} = 0\}$ . Given  $x \in \mathbb{H}_K^n$ , let  $u \in \mathcal{T}_x \mathbb{H}_K^n$  be unit-speed. The unique unit-speed geodesic  $\gamma_{x \rightarrow u}(\cdot)$  such that  $\gamma_{x \rightarrow u}(0) = x$  and  $\dot{\gamma}_{x \rightarrow u}(0) = u$  is denoted as  $\gamma_{x \rightarrow u}(t) = \cosh \left( \frac{t}{\sqrt{K}} \right) x + \sqrt{K} \sinh \left( \frac{t}{\sqrt{K}} \right) u, t > 0$ . The intrinsic distance between two points  $x, y \in \mathbb{H}_K^n$  is then equal to

$$d_{\mathcal{M}}^K(x, y) = \sqrt{K} \text{arcosh} \left( -\frac{\langle x, y \rangle_{\mathcal{M}}}{K} \right).$$

Therefore, the above  $n$ -dimensional hyperbolic space with constant negative sectional curvature  $-\frac{1}{K}$  is usually denoted as  $(\mathbb{H}_K^n, d_{\mathcal{M}}^K(\cdot, \cdot))$ . In particular,  $\mathbb{H}_1^n$ , i.e.  $K = 1$ , is called the hyperboloid model of the hyperbolic space. Hyperbolic Graph Convolutional Networks (HGCM) [15] benefit from the expressiveness of both GCNNs and hyperbolic embedding. It employs the exponential and logarithmic maps of the

hyperboloid model, respectively denoted as  $\exp_x^K(\cdot)$  and  $\log_x^K(\cdot)$ , to realize the mutual transformation between Euclidean features and hyperbolic ones. Let  $\|v\|_{\mathcal{M}} = \langle v, v \rangle_{\mathcal{M}}, v \in \mathcal{T}_x \mathbb{H}_K^n$ . The  $\exp_x^K(\cdot)$  and  $\log_x^K(\cdot)$  are respectively defined to be

$$\begin{aligned} \exp_x^K(v) &= \cosh\left(\frac{\|v\|_{\mathcal{M}}}{\sqrt{K}}\right) x + \sqrt{K} \sinh\left(\frac{\|v\|_{\mathcal{M}}}{\sqrt{K}}\right) \frac{v}{\|v\|_{\mathcal{M}}} \\ \log_x^K(y) &= d_{\mathcal{M}}^K(x, y) \frac{y + \frac{1}{K} \langle x, y \rangle_{\mathcal{M}} x}{\|y + \frac{1}{K} \langle x, y \rangle_{\mathcal{M}} x\|_{\mathcal{M}}}, \end{aligned}$$

where  $x \in \mathbb{H}_K^n$ ,  $v \in \mathcal{T}_x \mathbb{H}_K^n$  and  $y \in \mathbb{H}_K^n$  such that  $y \neq 0$  and  $y \neq x$ . The HGNN architecture is composed of three components, namely a hyperbolic feature transform, an attention-based aggregation and a non-linear activation with different curvatures. They are respectively defined to be

$$\begin{aligned} h_j^{H,l} &= \left(W^{(l)} \otimes^{K_{l-1}} x_j^{H,l-1}\right) \oplus^{K_{l-1}} b^{(l)} && \text{(hyperbolic feature transform),} \\ y_j^{H,l} &= \text{AGGREGATE}^{K_{l-1}}(h^{H,l})_j && \text{(attention-based aggregation),} \\ x_j^{H,l} &= \exp_o^{K_l} \left(\rho \left(\log_o^{K_{l-1}} \left(y_j^{H,l}\right)\right)\right) && \text{(non-linear activation with different curvatures),} \end{aligned}$$

where  $o = (\sqrt{K}, 0, \dots, 0) \in \mathbb{H}_K^n$ , the subscript  $j$  denotes the indices of nodes, the superscript  $l$  denotes the layer of the HGNN. The linear transform in hyperboloid manifold is defined to be  $W \otimes^K x^H = \exp_o^K \left(W \log_o^K(x^H)\right)$  and  $x^H \oplus^K b = \exp_{x^H}^K \left(P_{o \rightarrow x^H}^K(b)\right)$ , where  $P_{o \rightarrow x^H}^K(b)$  is the parallel transport from  $\mathcal{T}_o \mathbb{H}_K^n$  to  $\mathcal{T}_{x^H} \mathbb{H}_K^n$ . The attention-based aggregation is defined to be

$$\text{AGGREGATE}^K(x^H)_j = \exp_{x_j^H}^K \left( \sum_{k \in N_G(j)} \omega_{j,k} \log_{x_j^H}^K(x_k^H) \right),$$

where the attention weight  $\omega_{j,k} = \text{Softmax}_{k \in N_G(j)} \left( \text{MLP} \left( \log_o^K(x_j^H) \bowtie \log_o^K(x_k^H) \right) \right)$ .

**Higher-Order Spatial GCNNs.** the aforementioned GCNN architectures are constructed from the microscopic perspective. They only consider nodes and edges, yet overlook the higher-order substructures and their connections, i.e subgraphs consisting of at least 2 nodes. Here, we introduce the studies on the  $k$ -dimensional GCNNs [88]. Specifically, they take higher-order graph structures at multiple scales into consideration by leveraging the  $k$ -Weisfeiler-Leman ( $k$ -WL) graph isomorphism test so that the message passing is performed directly between subgraph structures rather than individual nodes. Let  $\{\{\dots\}\}$  denote a multiset,  $\text{HASH}(\cdot)$  a hashing function and  $C_{l,k}^{(l)}(s)$  the node coloring (label) of  $s = (s_1, \dots, s_k) \in V^k$  at the  $l$ -th time. Moreover, let  $N_G^j(s) = \{(s_1, \dots, s_{j-1}, r, s_{j+1}, \dots, s_k) : r \in V\}$ . The  $k$ -WL is computed by

$$C_{l,k}^{(l+1)}(s) = \text{HASH} \left( C_{l,k}^{(l)}(s), \left( c_1^{(l+1)}(s), \dots, c_k^{(l+1)}(s) \right) \right),$$

where

$$c_j^{(l+1)} = \text{HASH} \left( \left\{ \left\{ C_{l,k}^{(l)}(s') : s' \in N_G^j(s) \right\} \right\} \right).$$

Let  $N_G(s) = \bigcup_{j=1}^k N_G^j(s)$ . The local neighborhood  $N_{G,L}(s)$  is composed of all  $t \in N_G(s)$  such that  $(v, w) \in E$  for the unique  $v \in s - t$  and the unique  $w \in t - s$ . The  $k$ -GCNN computes new features of  $s \in V^k$  by multiple computational layers. Each layer is computed by

$$X_k^{(l+1)}[s, :] = \rho \left( X_k^{(l)}[s, :] W_1^{(l)} + \sum_{t \in N_G(s)} X_k^{(l)}[t, :] W_2^{(l)} \right).$$

In practice, the local  $k$ -GCNNs is usually used to scale  $k$ -GCNNs to larger datasets and to prevent overfitting, in practice, i.e.

$$X_k^{(l+1)}[s, :] = \rho \left( X_k^{(l)}[s, :] W_1^{(l)} + \sum_{t \in N_{G,L}(s)} X_k^{(l)}[t, :] W_2^{(l)} \right).$$

Furthermore, in order to learn hierarchical representations, the initial input of  $k$ -GCNN is just equal to the output features learned by a  $(k - 1)$ -GCNN, i.e.

$$X_k^{(0)}[s, :] = \rho \left( X^{iso}[s, :] \bowtie \sum_{u \in s} X_{k-1}^{(L_{k-1})}[u, :] \right).$$

Note that we abuse the notation  $X_k^{(l)}[s, :]$  to let it denote the feature vector corresponding to  $s \in V^k$ .

**Variants of GNNs—Miscellaneous.** In addition to the aforementioned GN framework and its variants, there are also other spatial GCNNs which is defined from other viewpoints. Diffusion-Convolutional Neural Network (DCNN) [20], which is permutation-invariant under graph isomorphism, employs  $H$ -hop transition matrix to capture distant connectivity information. Position-aware Graph Neural Networks (P-GNNs) [46] can compute position-aware node embedding vectors via learning a non-linear distance-weighted aggregation scheme over sampled anchor sets. Memory-based graph networks, namely Memory-based Graph Neural Network (MemGNN) and Graph Memory Network (GMN) [84], are an efficient memory layer for GNNs, which can jointly learn node hierarchical representations. Graph Partition Neural Networks (GPNNs) [35] are able to tackle extremely large graphs by virtue of a heuristic graph partition algorithm. Specifically, they alternate between nodes in partitions and globally propagating information between the partitions. Edge-Conditioned Convolution (ECC) [43] proposes to condition filter weights on the respective edge labels via a filter-generating network, borrowed from dynamic filter networks [385], mapping from edge labels to edge-specific weight matrix. The literature [21] argues that spatial graph convolution should have three properties, namely seed-oriented, degree-aware and order-free, and proposes a generic degree-specific spatial GCNN named DEMO-Net via hash kernel [386]. Column networks [29] is a novel deep learning model for collective classification in multi-relational graphs. Each entity aggregates information from relations of different types associated with it, and then update its state at the next layer via the aggregated information and the state at the current layer. The literature [32] defines spatial convolution filters as a polynomial of functions of the graph adjacency matrix.

**Invariance and Equivariance.** Let  $P \star A_G = P^T A_G P$  denote the tensor resulting from renumbering the nodes in  $V_G$  according to the permutation matrix  $P$ . Permutation-invariance refers to that a function  $f : \mathbb{R}^{n^k} \rightarrow \mathbb{R}$  (e.g. the aggregation function) is independent of any permutations of node/edge indices [50, 389], i.e.  $f(P \star A_G) = f(A_G)$  where  $P$  is a permutation matrix and  $A_G \in \mathbb{R}^{n^k}$  is a  $k$ -order tensor of edges or multi-edges in the (hyper-)graph  $G$ . Permutation-equivariance refers to that a function  $f : \mathbb{R}^{n^k} \rightarrow \mathbb{R}^{n^l}$  coincides with permutations of node/edge indices [50, 389], i.e.  $f(P \star A_G) = P \star f(A_G)$  where  $P$  and  $A_G$  are defined as similarly as the permutation-invariance. The literature [50] considers a specific class of invariant and equivariant GCNNs with a single hidden layer, for which we prove new universality theorems. The one-layer GCNNs are obtained by summing channels formed by applying an equivariant linear operator, a pointwise non-linearity and either an invariant or equivariant linear operator. In particular, for permutation-invariant aggregation functions, A straightforward choice is to take sum/ max /average/concatenation as heuristic aggregation schemes. Nevertheless, these aggregation functions treat all the neighbors of a vertex equivalently so that they cannot precisely distinguish the structural effects of different neighbors to the target vertex. That is, the aggregation functions should extract and filter graph signals aggregated from neighbors of different hops away and different importance. GeniePath [34] proposes a scalable approach for learning adaptive receptive fields of GCNNs. It, as an adaptive path layer, is composed of two complementary functions, namely adaptive breadth function and adaptive depth function. The former learns the importance of different sized neighborhoods, whereas the latter extracts and filters graph signals aggregated from neighbors of different hops away. More specifically, the adaptive breadth function is defined as follows.

$$h_{v_j}^{\text{temp}} = \tanh \left( (W^{(t)})^T \sum_{v_k \in N_G(v_j) \cup \{v_j\}} \alpha(h_{v_j}^{(t)}, h_{v_k}^{(t)}) \cdot h_{v_k}^{(t)} \right),$$

where

$$\alpha(x, y) = \text{Softmax}_y (\alpha^T \tanh (W_x^T x + W_y^T y)).$$

The adaptive depth function is defined as a LSTM [160], i.e.

$$\begin{aligned}
 i_{v_j} &= \sigma \left( \left( W_i^{(t)} \right)^T h_{v_j}^{\text{temp}} \right) & f_{v_j} &= \sigma \left( \left( W_f^{(t)} \right)^T h_{v_j}^{\text{temp}} \right) \\
 o_{v_j} &= \sigma \left( \left( W_o^{(t)} \right)^T h_{v_j}^{\text{temp}} \right) & \tilde{C}_{v_j} &= \tanh \left( \left( W_c^{(t)} \right)^T h_{v_j}^{\text{temp}} \right) \\
 C_{v_j}^{(t+1)} &= f_{v_j} \otimes C_{v_j}^{(t)} + i_{v_j} \otimes \tilde{C}_{v_j} & h_{v_j}^{(t+1)} &= o_{v_j} \otimes \tanh \left( C_{v_j}^{(t+1)} \right).
 \end{aligned} \tag{12}$$

GEOM-GCN [17] proposes a novel permutation-invariant geometric aggregation scheme consisting of three modules, namely node embedding, structural neighborhood, and bi-level aggregation (i.e. low-level aggregation, high-level aggregation and non-linear aggregation). The geometric aggregation scheme can overcome two fundamental weaknesses, namely losing structural information of nodes and failing to capture long-range dependencies in disassortative graphs. PiNet [41] is an end-to-end spatial GCNN architecture that utilizes the permutation equivariance of graph convolutions in order to learn permutation-invariant graph representations. It is composed of a pair of double-stacked message passing layers, namely attention-oriented message passing layers and feature-oriented message passing layers, combined by a matrix product.

**Depth Trap of Spatial GCNNs.** Like the spectral GCNNs, spatial GCNNs also go through a depth trap. That is, the performance gets decayed as the depth of spatial GCNNs increases. As stated previously, the depth trap results from oversmoothing, overfitting and gradient vanishing/explosion. In order to escape from the depth trap, DeepGCN [13] borrows concepts from CNNs, specifically residual/dense connections [382, 383] and dilated aggregation [384], and adapts them to the spatial GCNN architecture. That is, DeepGCN has three instantiations, namely ResGCN, DenseGCN and dilated graph convolution. ResGCN is inspired by the ResNet [382], which is defined to be

$$\begin{aligned}
 G^{(l+1)} &\triangleq \mathcal{H}(G^{(l)}, \Theta^{(l)}) \\
 &= \mathcal{F}(G^{(l)}, \Theta^{(l)}) + G^{(l)},
 \end{aligned}$$

where  $\mathcal{F}(\cdot, \cdot)$  is a map taking a graph as input and outputting an updated graph representation. In practice, it can be computed by spectral or spatial GCNNs. DenseGCN collectively exploit information from different GCNN layers like the DenseNet [383], which is defined to be

$$\begin{aligned}
 G^{(l+1)} &\triangleq \mathcal{H}(G^{(l)}, \Theta^{(l)}) \\
 &= \text{Concat}(\mathcal{F}(G^{(l)}, \Theta^{(l)}), G^{(l)}) \\
 &= \text{Concat}(\mathcal{F}(G^{(l)}, \Theta^{(l)}), \dots, \mathcal{F}(G^{(0)}, \Theta^{(0)}), G^{(0)}),
 \end{aligned}$$

where  $\text{Concat}(\cdot)$  is a vertex-wise concatenation function that densely fuse input graph  $G^{(0)}$  with all the intermediate GCNN layer outputs. The dilated aggregation [384] can magnify the receptive field of spatial GCNNs by a dilation rate  $d$ . More specifically, let  $N_G^{(k,d)}(v)$  denote the set of  $k$   $d$ -dilated neighbors of vertex  $v$  in  $G$ . If  $(u_1, u_2, \dots, u_{k \times d})$  are the first sorted  $k \times d$  nearest neighbors, then

$$N_G^{(k,d)}(v) = \{u_1, u_{1+d}, \dots, u_{1+(k-1)d}\}.$$

Thereby, we can construct a new graph  $G^{(k,d)} = (V^{(k,d)}, E^{(k,d)})$  where  $V^{(k,d)} = V$  and  $E^{(k,d)} = \{\langle v, u \rangle : v \in V, u \in N_G^{(k,d)}(v)\}$ . The aggregation and update functions of the spital GCNN are applied to  $G^{(k,d)}$  so as to obtain a dilated graph convolution layer. The jumping knowledge networks [89] adopts neighborhood aggregation with skip connections to integrate information from different layers. More formally, each node is updated as

$$\begin{aligned}
 x_{N_G(v)}^{(l)} &= \rho \left( W^{(l)} \text{AGGREGATE} \left( \left\{ x_u^{(l-1)} : u \in N_G(v) \right\} \right) \right) \\
 x_v^{(l)} &= \text{COMBINE} \left( x_v^{(l-1)}, x_{N_G(v)}^{(l)} \right),
 \end{aligned}$$

where the COMBINE function embodies the skip connection between different layers.

### 3.1.3 Graph Wavelet Neural Networks

As stated previously, the spectral and spatial GCNNs are respectively inspired by the graph Fourier transform and message-passing mechanism. Here, we introduce a new GCNN architecture from the perspective of the Spectral Graph Wavelet Transform (SGWT) [189]. First of all, the SGWT is determined by a graph wavelet generating kernel  $g: \mathbb{R}^+ \rightarrow \mathbb{R}^+$  with the property  $g(0) = 0, g(+\infty) = \lim_{x \rightarrow \infty} g(x) = 0$ . An feasible instance of  $g(\cdot)$  is parameterized by two integers  $\alpha$  and  $\beta$ , and two positive real numbers  $x_1$  and  $x_2$  determining the transition regions, i.e.

$$g(x; \alpha, \beta, x_1, x_2) = \begin{cases} x_1^{-\alpha} x^\alpha & x < x_1 \\ s(x) & x_1 \leq x \leq x_2 \\ x^{-\beta} x_2^\beta & x > x_2, \end{cases}$$

where  $s(x)$  is a cubic polynomial whose coefficients can be determined by the continuity constraints  $s(x_1) = s(x_2) = 1$ ,  $s'(x_1) = \frac{\alpha}{x_1}$  and  $s'(x_2) = -\frac{\beta}{x_2}$ . Given the graph wavelet generating kernel  $g(\cdot)$  and a scaling parameter  $s \in \mathbb{R}^+$ , the spectral graph wavelet operator  $\Psi_g^s$  is defined to be

$$\Psi_g^s = U g(s\Lambda) U^T, \quad (13)$$

where  $g(s\Lambda) = g(\text{diag}(s\lambda_1, \dots, s\lambda_N))$ . A graph signal  $x \in \mathbb{R}^N$  on  $G$  can thereby be filtered by the spectral graph wavelet operator, i.e.

$$\mathcal{W}_{g,s}^x = \Psi_g^s x \in \mathbb{R}^N. \quad (14)$$

The  $j$ -th entry of  $\mathcal{W}_{g,s}^x$  in formula (14) is equal to  $\mathcal{W}_{g,s}^x[j] = (\Psi_g^s x)[j]$ ,  $j = 1, \dots, N$ . Let  $\psi_g^{s,j} = \Psi_g^s \delta_j$  where  $\delta_j(k) = 1$  if  $k = j$  and 0 otherwise. It is therefore easy to prove that

$$\mathcal{W}_{g,s}^x[j] = \langle \psi_g^{s,j}, x \rangle,$$

where  $\langle \cdot, \cdot \rangle$  denotes the inner product of two vectors. Note that directly computing the spectral graph wavelet transformation via Eq. (14) is computationally expensive especially for large graphs since the eigendecomposition of the Laplacian  $\bar{L}_G$  commonly requires  $O(N^3)$  computational complexity. Here, we can also employ the  $K$ -localized Chebyshev polynomial [78] to fast approximate the spectral graph wavelet transformation, i.e.

$$\begin{aligned} \Psi_g^s &= \frac{1}{2} c_0^s + \sum_{k=1}^K c_k^s T_k(\tilde{L}) \\ c_k^s &= 2e^{-s} J_k(-s), \end{aligned} \quad (15)$$

where  $\tilde{L} = \frac{2}{\lambda_{\max}} \bar{L} - \mathbb{I}_N$ ,  $T_j(x)$  is the Chebyshev polynomial (see formula 7) and  $J_k(-s)$  is the Bessel function of the first kind.

**Graph Wavelet Neural Networks.** The above spectral graph wavelet operator  $\Psi_g^s$  can be employed to construct a Graph Wavelet Neural Network (GWNN) [191]. Let  $\Psi_g^{-s} \triangleq (\Psi_g^s)^{-1}$ . The graph wavelet based convolution is defined to be

$$x *_G y = \Psi_g^{-s} (\Psi_g^s x \otimes \Psi_g^s y)$$

by replacing the graph Fourier transform in formula (4) with the graph wavelet transform. The GWNN is composed of multiple layers of the graph wavelet based convolution. The structure of the  $l$ -th layer is defined as

$$X^{(l+1)}[:, j] = \rho \left( \sum_{k=1}^{d^{(l)}} \Psi_g^{-s} \Theta_{j,k}^{(l)} \Psi_g^s X^{(l)}[:, k] \right), \quad (16)$$

where  $\Theta_{j,k}^{(l)}$  is a diagonal filter matrix learned in spectral domain. In order to reduce the number of learnable parameters, Eq. (16) can be rewritten as a matrix form, i.e.

$$X_{\text{temp}}^{(l)} = X^{(l)} W^{(l)} \quad (\text{Feature Transformation})$$

$$X^{(l+1)} = \rho \left( \Psi_g^{-s} \Theta \Psi_g^s X_{\text{temp}}^{(l)} \right) \quad (\text{Graph Wavelet based Convolution}).$$

In practice, the learnable filter matrix  $\Theta$  can be replaced with the  $K$ -localized Chebyshev Polynomial in formula (15) so as to eschew the eigendecomposition of the Laplacian  $\bar{L}_G$ .

**Graph Scattering Networks.** The literature [190] utilizes a special instance of the graph wavelet operator  $\Psi_g^s$  in Eq. (13) to construct a graph scattering network, and proves its covariance and approximate invariance to permutations and stability to graph manipulations. Suppose  $X = (x_1, \dots, x_d) \in \mathbb{R}^{N \times d}$  is a  $d$ -dimensional graph signal on  $G$ . The graph scattering transform is defined to be

$$S[\mathcal{P}_J]X = (S[\mathcal{P}_J]x_j)_{j=1}^d, \quad W[\mathcal{P}_J]X = (W[\mathcal{P}_J]x_j)_{j=1}^d$$

where  $S[\mathcal{P}_J]x_j = (S[p]x_j)_{p \in \mathcal{P}_J}$ ,  $W[\mathcal{P}_J]x_j = (W[p]x_j)_{p \in \mathcal{P}_J}$  and  $\mathcal{P}_j = \bigcup_{m=0}^{\infty} \{p : \text{length}(p) = m\}$  denotes the collection of all paths of finite length. The windowed scattering transform of a graph signal  $x$  on  $G$  with respect to a path  $p$  of length  $m$  is defined by

$$S[p]x = U \text{diag}(\phi(2^{-J}\lambda_1), \dots, \phi(2^{-J}\lambda_N)) U^T W[p]x,$$

where  $\phi(\cdot)$  is a scaling function, and  $W[p]x = W[j_m]W[j_{m-1}] \cdots W[j_1]x$ . The one-step propagator  $W[\psi_j] : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is defined by  $W[\psi_j]x = f * \psi_j = W \text{diag}(\psi(2^{-j}\lambda_1), \dots, \psi(2^{-j}\lambda_N)) W^T x$ , where  $\psi_j(n) = \psi(2^{-j}n)$  is a wavelet function.

### 3.1.4 GCNNs on Special Graphs

The aforementioned GCNNs (possibly with graph pooling operators or attention mechanisms) aims at learning the node/edge/graph representations of input graphs (directed or undirected, weighted or unweighted). The graphs extracted from the real-world may have more additional characteristics, e.g. spatial-temporal graphs, heterogeneous graphs, hyper-graphs, signed graphs and so all. The literature [28] proposes a GCNN for signed graphs by leveraging balance theory to correctly aggregate and propagate information across positive and negative links. Below, we introduce the GCNN architectures on these special graphs.

**Heterogeneous Graphs.** Heterogeneous Graphs are composed of nodes and edges of different types, and each type of edges is called a relation between two endpoints. For example, a bibliographic information network contains at least 4 types of nodes, namely Author, Paper, Venue and Term, and at least 3 types of edges, namely Author-Paper, Term-Paper and Venue-Paper [395]. The heterogeneity and rich semantic information brings great challenges for designing heterogeneous graph convolutional neural networks. In general, heterogeneous graphs can be denoted as  $H = (V, E, \nu, \zeta)$ , where  $\nu(v)$  denotes the type of node  $v \in V$  and  $\zeta(e)$  denotes the type of edge  $e \in E$ . Let  $\mathcal{T}^v$  and  $\mathcal{T}^e$  respectively denote the set of node types and edge types. Below, we summarize the heterogeneous GCNNs from two perspectives, namely vanilla heterogeneous GCNNs and heterogeneous graph attention mechanism.

**(1) Vanilla Heterogeneous GCNNs.** The Heterogeneous Graph Neural Networks (HetGNNs) [61] aims to resolve the issue of jointly considering heterogeneous structural information as well as heterogeneous contents information of nodes. It firstly samples a fixed size of strongly correlated heterogeneous neighbors for each node via a Random Walk with Restart (RWR) and groups them into different node types. Then, it aggregates feature information of those sampled neighboring nodes via a bi-directional Long Short Term Memory (LSTM) and attention mechanism. Running RWR with a restart probability  $p$  from node  $v$  will yield a collection of a fixed number of nodes, denoted as  $\text{RWR}(v)$ . For each node type  $t$ , the  $t$ -type neighbors  $N_G^t(v)$  of node  $v$  denotes the set of top- $k_t$  nodes from  $\text{RWR}(v)$  with regard to frequency. Let  $\mathcal{C}_v$  denote the heterogeneous contents of node  $v$ , which can be encoded as a fixed size embedding via a function  $f_1(v)$ , i.e.

$$f_1(v) = \frac{\sum_{j \in \mathcal{C}_v} \left[ \overrightarrow{\text{LSTM}}(\mathcal{F}\mathcal{C}_{\theta_x}(x_j)) \bowtie \overleftarrow{\text{LSTM}}(\mathcal{F}\mathcal{C}_{\theta_x}(x_j)) \right]}{|\mathcal{C}_v|},$$

where  $\mathcal{F}\mathcal{C}_{\theta_x}(\cdot)$  denotes feature transformer, e.g. identity or fully connected neural networks with parameter  $\theta_x$ , and the definitions of  $\overrightarrow{\text{LSTM}}$  and  $\overleftarrow{\text{LSTM}}$  refer to the literature [160], see Eq. (12). Now, the

content embedding of the  $t$ -type neighbors of node  $v$  can be aggregated as follows,

$$\begin{aligned} f_2^t(v) &= \text{AGGREGATE}^T(\{f_1(v') : v' \in N_G^t(v)\}) \\ &= \frac{\sum_{v' \in N_G^t(v)} [\overrightarrow{\text{LSTM}}(f_1(v')) \bowtie \overleftarrow{\text{LSTM}}(f_1(v'))]}{|N_G^t(v)|}. \end{aligned}$$

Let  $\mathcal{F}(v) = \{f_1(v)\} \cup \{f_2^t(v) : t \in \mathcal{T}^v\}$ . As a result, the output embedding of node  $v$  can be obtained via the attention mechanism, i.e.

$$\mathcal{E}_v = \sum_{f_j(v) \in \mathcal{F}(v)} \omega_{v,j} f_j(v),$$

where the attention weights is computed by

$$\omega_{v,j} = \frac{\exp(\rho(u^t [f_j(v) \bowtie f_1(v)]))}{\sum_{f_j(v) \in \mathcal{F}(v)} \exp(\rho(u^t [f_j(v) \bowtie f_1(v)]))}.$$

The objective of the HetGNN is the Maximum Likelihood Estimation (MLE) of the heterogeneous softmax function, and is optimized by employing the Negative Sampling (NE) techniques. The literature [62] aims to learn the deep relational features on heterogeneous graphs via the proposed GraphInception. This model firstly converts the input heterogeneous graph into a multi-channel graph (each meta path corresponds to a channel), and then employs graph inception module inspired by the Inception module [397] to automatically generate a hierarchy of relational features from simple to complex ones.

**(2) Heterogeneous Graph Attention Mechanism.** The literature [75] firstly proposes a hierarchical attention based heterogeneous GCNNs consisting of node-level and semantic-level attentions. The node-level attention aims to learn the attention weights of a node and its meta-path-based neighbors, and the semantic-level attention is able to learn the importance of different meta-paths. More specifically, given a meta path  $\Phi$ , the node-level attention weight of a node  $v_i$  and its meta-path-based neighbors  $v_j \in N_G^\Phi(v_i)$  is defined to be

$$\alpha_{j,k}^\Phi = \frac{\exp(\rho(a_\Phi^T(M_{\nu(v_j)}x_j \bowtie M_{\nu(v_k)}x_k)))}{\sum_{v_k \in N_G^\Phi(v_j)} \exp(\rho(a_\Phi^T(M_{\nu(v_j)}x_j \bowtie M_{\nu(v_k)}x_k)))},$$

where  $M_{\nu(v_j)}$  transforms the feature vectors of nodes of type  $\nu(v_j)$  in different vector spaces into a unified vector space. The embedding of node  $v_i$  under the meta path  $\Phi$  can be computed by

$$x_j^{\Phi,l+1} = \bowtie_{k=1}^K \rho \left( \sum_{k \in N_G^\Phi(v_j)} \omega_{j,k}^\Phi M_{\nu(v_k)} x_k^{\Phi,l} \right).$$

Given a meta-path set  $\{\Phi_0, \dots, \Phi_P\}$ , performing the node-level attention layers under these meta paths will yield a set of semantic-specific node representations, namely  $\{X^{\Phi_0}, \dots, X^{\Phi_P}\}$ . The semantic-level attention weight of the meta path  $\Phi_j$  is defined as

$$\beta^{\Phi_j} = \frac{\exp(\omega^{\Phi_j})}{\sum_{p=1}^P \exp(\omega^{\Phi_p})},$$

where  $\omega^{\Phi_p} = \frac{1}{|V|} \sum_{v_k \in V} q^T \tanh(Wx_k^{\Phi_p} + b)$ . The final embedding matrix  $X = \sum_{p=1}^P \beta^{\Phi_p} X^{\Phi_p}$ . Moreover, the literature [36, 68, 361] studies the GCNNs for multi-relational graphs. They usually establish attention-based neural architectures with different learnable parameters under different relations, and then update hidden representations of nodes by aggregating all the outputs from different architectures. some other studies [60, 114] borrows the concept of transformer from the literature [73] to learn the node representations of dynamic heterogeneous graphs via self-attention mechanisms and positional encoding heuristics.



**Spatio-Temporal Graphs.** Spatio-temporal graphs can be used to model traffic networks [59, 71] and skeleton networks [55, 57]. In general, a spatio-temporal graph is denoted as  $G_{ST} = (V_{ST}, E_{ST})$  where  $V_{ST} = \{v_{t,j} : t = 1, \dots, T, j = 1, \dots, N_{ST}\}$ . The edge set  $E_{ST}$  is composed of two types of edges, namely spatial edges and temporal edges. All spatial edges  $(v_{t,j}, v_{t,k}) \in E_{ST}$  are collected in the intra-frame edge set  $E_S$ , and all temporal edges  $(v_{t,j}, v_{t+1,j}) \in E_{ST}$  are collected in the inter-frame edge set  $E_T$ . The literature [59] proposes a novel deep learning framework, Spatio-Temporal Graph Convolutional Networks (STGCN), to tackle the traffic forecasting problem. Specifically, STGCN consists of several layers of spatio-temporal convolutional blocks, each of which has a "sandwich" structure with two temporal gated convolution layers (abbreviated as Temporal Gated-Conv) and a spatial graph convolution layer in between (abbreviated as Spatial Graph-Conv). The Spatial Graph-Conv exploits the conventional GCNNs to extract the spatial features, whereas the Temporal Gated-Conv the temporal gated convolution operator to extract temporal features. Suppose that the input of the temporal gated convolution for each node is a length- $M$  sequence with  $C_{in}$  channels, i.e.  $X \in \mathbb{R}^{M \times C_{in}}$ . The temporal gated convolution kernel  $\Gamma \in \mathbb{R}^{K \times C_{in} \times 2C_{out}}$  is used to filter the input  $Y$  to yield an output  $P \bowtie Q \in \mathbb{R}^{(M-K+1) \times (2C_{out})}$ , i.e.

$$\Gamma *_T Y = P \circledast \sigma(Q) \in \mathbb{R}^{(M-K+1) \times C_{out}},$$

where the sigmoid gate  $\sigma(Q)$  controls which inputs in  $P$  are relevant for discovering compositional structure and dynamic variances in time series. The Spatial Graph-Conv takes a tensor  $\underline{X}^{(l)} \in \mathbb{R}^{M \times N_{ST} \times C^{(l)}}$  as input, and outputs a tensor  $\underline{X}^{(l+1)} \in \mathbb{R}^{(M-2(K-1)) \times N_{ST} \times C^{(l+1)}}$ , i.e.

$$\underline{X}^{(l+1)} = \Gamma_1^{(l)} *_T \rho \left( \Theta^{(l)} *_G \left( \Gamma_0^{(l)} *_T X^{(l)} \right) \right),$$

where  $\Gamma_0^{(l)}, \Gamma_1^{(l)}$  are the upper and lower temporal kernel and  $\Theta^{(l)}$  is the spectral kernel of the graph convolution. Structural-RNN [56] is derived from the factor graph representation of the spatio-temporal graph. The factor graph representation has a factor  $\Psi_v(x_v, y_v)$  for each node  $v \in V_{ST}$  and a pairwise factor  $\Psi_e(x_e, y_{e[0]}, y_{e[1]})$  for each edge  $e \in E_{ST}$ . Note that the Structural-RNN assigns semantically similar nodes to the same factor so as to reduce the number of learnable parameters, and so do the edges. Each type of node factor is essentially a nodeRNN capturing the temporal dependencies of the nodes, whereas each type of pairwise factor is essentially an edgeRNN capturing the temporal dependencies of the edges. In order to capture the spatial dependencies of nodes and edges, the outputs of each type of edgeRNN is feed into the nodeRNN when the edge is incident to the node. The literatures [108, 109] propose a gated Graph Convolutional Recurrent Neural Network (GCRNN) combining the vanilla GCNNs and RNNs to learn the spatio-temporal representations of nodes on spatio-temporal graphs. That is, it employs the RNN architecture to extract the long short-term dependencies in the outputs  $(X^{(1)}, X^{(2)}, \dots, X^{(L)})$  of layers of the vanilla GCNN.

**Hypergraphs.** The aforementioned GCNN architectures are concerned with the conventional graphs consisting of pairwise connectivity between two nodes. However, there could be even more complicated connections between nodes beyond the pairwise connectivity, e.g. co-authorship networks. Under such circumstances, a hypergraph, as a generalization to the conventional graph, provides a flexible and elegant modeling tools to represent these complicated connections between nodes. A hypergraph is usually denoted as  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \omega)$ , where  $\mathcal{V} = \{v_1, \dots, v_N\}$  like the conventional graph,  $\mathcal{E} = \{e_1, \dots, e_M\}$  is a set of  $M$  hyperedges.  $\omega(e_k)$  denote weights of hyperedges  $e_k \in \mathcal{E}$ . A non-trivial hyperedge is a subset of  $\mathcal{V}$  with at least 2 nodes. In particular, a trivial hyperedge, called a self-loop, is composed of a single node. The hypergraph  $\mathcal{G}$  can also be denoted by an incidence matrix  $\mathcal{H}_{\mathcal{G}} \in \mathbb{R}^{N \times M}$ , i.e.

$$\mathcal{H}_{\mathcal{G}}[j, k] = \begin{cases} 0, & v_j \notin e_k \\ 1, & v_j \in e_k. \end{cases}$$

For a node  $v_j \in \mathcal{V}$ , its degree  $\deg_{\mathcal{V}}(v_j) = \sum_{e_k \in \mathcal{E}} \omega(e_k) \mathcal{H}_{\mathcal{G}}[j, k]$ . For a hyperedge  $e_k \in \mathcal{E}$ , its degree

$\deg_{\mathcal{E}}(e_k) = \sum_{v_j \in e_k} \mathcal{H}_{\mathcal{G}}[j, k]$ . Let

$$\begin{aligned} \mathcal{D}_{\mathcal{V}} &= \text{diag}(\deg_{\mathcal{V}}(v_1), \dots, \deg_{\mathcal{V}}(v_N)), \\ \mathcal{D}_{\mathcal{E}} &= \text{diag}(\deg_{\mathcal{E}}(e_1), \dots, \deg_{\mathcal{E}}(e_M)), \\ \mathcal{W}_{\mathcal{G}} &= \text{diag}(\omega(e_1), \dots, \omega(e_M)). \end{aligned}$$

According to the literature [65], the hypergraph Laplacian  $\mathcal{L}_{\mathcal{G}}$  of  $\mathcal{G}$  is defined to be

$$\mathcal{L}_{\mathcal{G}} = \mathcal{I}_N - \mathcal{D}_{\mathcal{V}}^{-\frac{1}{2}} \mathcal{H}_{\mathcal{G}} \mathcal{W}_{\mathcal{G}} \mathcal{D}_{\mathcal{E}}^{-1} \mathcal{H}_{\mathcal{G}}^T \mathcal{D}_{\mathcal{V}}^{-\frac{1}{2}}.$$

The hypergraph Laplacian can also be factorized by the eigendecomposition, i.e.

$$\mathcal{L}_{\mathcal{G}} = \mathcal{U} \Lambda \mathcal{U}^T,$$

where  $\Lambda = \text{diag}\{\lambda_1, \dots, \lambda_N\}$  and  $\mathcal{U}$  is a orthonormal matrix. The spectral hypergraph convolution operator, the Chebyshev hypergraph convolutional neural network and the hypergraph convolutional network can be defined in analogy to Eqs (4,6,8). The HyperGraph Neural Network (HGNN) architecture proposed in the literature [65] is composed of multiple layers of the hyperedge convolution, which is defined to be

$$X^{(l+1)} = \rho \left( \mathcal{D}_{\mathcal{V}}^{-\frac{1}{2}} \mathcal{H}_{\mathcal{G}} \mathcal{W}_{\mathcal{G}} \mathcal{D}_{\mathcal{E}}^{-1} \mathcal{H}_{\mathcal{G}}^T \mathcal{D}_{\mathcal{V}}^{-\frac{1}{2}} X^{(l)} \Theta^{(l)} \right).$$

Note that the HGNN essentially views each hyperedge as a complete graph so that the hypergraph is converted into a conventional graph. The HyperGCN [64] adopts a more efficient approach than the HGNN to performing the above conversion. For a  $d$ -dimensional graph signal  $\mathcal{X}^{(l)} \in \mathbb{R}^d$  on  $\mathcal{G}$  at the  $l$ th layer, the HyperGCN firstly computes the  $(j_e, k_e) = \arg \max_{j, k \in e} \|\mathcal{X}[j, :] - \mathcal{X}[k, :]\|^2$ . Let  $\mathcal{K}_e = \{m \in e : m \neq j_e, m \neq k_e\}$  denote the mediator set. Then, it constructs a new graph whose node set is  $\mathcal{V}$  and edge set includes the edges in  $\{(j_e, k_e)\} \cup \{(j_e, m) : m \in \mathcal{K}_e\} \cup \{(k_e, m) : m \in \mathcal{K}_e\}$ . Finally, the HyperGCN performs the graph convolution operators on the new graph to learn the node representations. The Directed Hypergraph Networks (DHN) [63] treats each hyperedge  $e \in \mathcal{E}$  as a node of a conventional graph  $G_{\mathcal{E}} = (\mathcal{E}, E_{\mathcal{E}})$ . Then, we employ the conventional GCNNs to update the hyperedge representations, i.e.  $\mathcal{X}_{\mathcal{E}}^{(l+1)} = f_{\text{GCNN}}(G_{\mathcal{E}}, \mathcal{X}_{\mathcal{E}}^{(l)})$ . At last, the node representations are interacted with the hyperedge representations via the following update formula,

$$\begin{aligned} \mathcal{X}_{\mathcal{V}}^{(l+1)} &= \rho \left( \mathcal{X}_{\mathcal{V}}^{(l)} \bowtie \left( \mathbb{I}_N \mathcal{X}_{\mathcal{E}}^{(l)} \Theta^{(l)} \right) \right), \\ \mathcal{X}_{\mathcal{E}}^{(l+1)} &= \rho \left( \mathbb{I}_N^T \mathcal{X}_{\mathcal{V}}^{(l+1)} \right). \end{aligned}$$

The HNHN [174] interleaves updating the node representations with the hyperedge representations by the following formulas,

$$\begin{aligned} \mathcal{X}_{\mathcal{V}}^{(l+1)} &= \rho \left( \mathcal{D}_{\mathcal{V}}^{-1} \mathcal{H}_{\mathcal{G}} \mathcal{X}_{\mathcal{E}}^{(l)} \Theta_{\mathcal{V}}^{(l)} \right), \\ \mathcal{X}_{\mathcal{E}}^{(l+1)} &= \rho \left( \mathcal{D}_{\mathcal{E}}^{-1} \mathcal{H}_{\mathcal{G}}^T \mathcal{X}_{\mathcal{V}}^{(l)} \Theta_{\mathcal{E}}^{(l)} \right). \end{aligned}$$

### 3.1.5 Summary

The aforementioned GCNN architectures have attracted considerable attention of many scholars all over the world, and been applied to many applications regarding the graphs. In practice, we can construct our own GCNNs by assembling a variety of modules introduced above. Additionally, some scholars also studies the GCNNs from some novel perspectives, e.g. the parallel computing framework of the GCNNs [40], the hierarchical covariant compositional networks [213], the transfer active learning for GCNNs [45] and quantum walk based subgraph convolutional neural network [51]. Note that these studies are closely related to the GCNNs yet fairly different from the ones introduced above.

### 3.2 Graph Pooling Operators

Graph pooling operators are very important and useful modules of the GCNNs especially for graph-level tasks such as graph classification. There are two kinds of graph pooling operators, namely global graph pooling operators and hierarchical graph pooling operators. The former aims to obtain the universal representations of input graphs, and the latter aims to capture adequate structural information for node representations.

**Global Graph Pooling Operators.** Global graph pooling operators pool all of representations of nodes into a universal graph representation. Many literatures [8, 21, 41] apply some simple global graph pooling operators, e.g. max/average/concatenate graph pooling, to performing graph-level classification tasks. Here, we introduce some more sophisticated global graph pooling operators in contrast to the simple ones. Relational pooling (RP) [98] provides a novel framework for graph representation with maximal representation power. Under this framework, all node embeddings can be aggregated via a learnable function to form a global embedding of  $G$ . Let  $X^{(v)} \in \mathbb{R}^{N \times d_v}$  and  $\underline{X}^{(e)} \in \mathbb{N}^{N \times N \times d_e}$  respectively denote node feature matrix and edge feature tensor. Tensor  $\underline{A}_G \in \mathbb{R}^{N \times N \times (1+d_e)}$  combines the adjacency matrix  $A_G$  of  $G$  with its edge feature tensor  $\underline{X}^{(e)}$ , i.e.  $\underline{A}_G[u, v, :] = \mathbb{I}_{(u,v) \in E_G} \bowtie \underline{X}^{(e)}[u, v, :]$ . After performing a permutation on  $V_G$ , the edge feature tensor  $\underline{A}_G^{(\pi, \pi)}[\pi(r), \pi(c), d] = \underline{A}_G[r, c, d]$  and the node feature matrix  $X_\pi^{(v)}[\pi(r), c] = X^{(v)}[r, c]$ . The joint RP permutation-invariant function for directed or non-directed graphs is defined as

$$\bar{f}(G) = \frac{1}{N!} \sum_{\pi \in \Pi_{|V|}} \vec{f}(\underline{A}_G^{\pi, \pi}, X_\pi^{(v)}),$$

where  $\Pi_{|V|}$  is the set of all distinct permutations of  $V_G$  and  $\vec{f}(\cdot, \cdot)$  is an arbitrary (possibly permutation-sensitive) vector-valued function of the graph  $G$ . Note that  $\vec{f}(\cdot, \cdot)$  can be denoted as Multi-Layer Perceptrons (MLPs), Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs) or Graph Neural Networks (GNNs). The literature [98] proves that  $\bar{f}(G)$  is the most expressive representation of  $G$  under some mild conditions, and provides approximation approaches to making RP computationally tractable. The literature [92] proposes a novel SortPooling layer which takes node features of a graph yielded by a spatial GCNN as input and outputs a sorted graph representation with a fixed size via traditional CNNs. Then, we concatenate all Weisfeiler-Lehman (WL) colors (i.e. feature vectors) of each node from all iteration and sort nodes in descending order according to their final WL colors. At last, CNNs are applied to the sorted representations for prediction. The literature [97] proposes a function space pooling method which maps a set of node representations to a function space. Specifically, this method firstly maps  $d$ -dimensional node representations to  $[0, 1]^d$  element-wise via a sigmoid function  $\sigma(x) = \frac{1}{1+e^{-x}}$ . Then, all the images are mapped to a function space via a  $d$ -dimensional Gaussian distribution with these images as centers, i.e.

$$g_u(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-u)^T(x-u)}{2\sigma^2}},$$

where  $u = \sigma(x)$  is an image of a node representation  $x$ . Finally, the function representation of  $G$  is defined as  $\beta(x) = \sum_{u \in S(D)} g_u(x)$ , where  $S(D)$  denotes the set of images of node representations, and can be approximated as a finite dimensional vector space by discretizing the function domain using a regular grid of elements.

**Hierarchical Graph Pooling Operators.** Hierarchical graph pooling operators group a set of proximal nodes into a super-node via graph clustering methods. Consequently, the original graph is coarsened to a new graph with coarse granularity. In practice, the hierarchical graph pooling operators are interleaved with the GCNN layers. In general, there are three kinds of approaches to performing the graph coarsening operations, namely invoking the existing graph clustering algorithms (e.g. spectral clustering [7] and Graclus [99]) and learning a soft cluster assignment and selecting the first  $k$  top-rank nodes.

(1) **Invoking existing graph clustering algorithms.** The graph clustering aims to assign proximal nodes to the same cluster. The proximity of two nodes can be measured by their structural similarity or

structural equivalence. After invoking a graph clustering algorithm, the output clusters can be viewed as super-nodes, and the connections between two clusters are viewed as super-edges. As a result, the coarsened graph is composed of all the super-nodes and super-edges. The hierarchical graph pooling operators aggregate the representations of nodes in super-nodes via aggregation functions such as max pooling and average pooling [91] to compute the representations of the super-nodes. The relationship between the original graph and the coarsened graph is referred to the literature [82] proposing the EigenPooling method. As stated previously,  $G$  contains  $N$  nodes and  $M$  edges, its adjacency matrix and a  $d$ -dimensional graph signal on  $G$  are respectively denoted as  $A_G$  and  $X$ . In order to construct a coarsened graph of  $G$ , a graph clustering method is employed to partition  $G$  into  $K$  disjoint clusters, namely  $\{G_k : k = 1, \dots, K\}$ . Suppose each cluster  $G_k$  has  $N_k$  nodes, namely  $\{v_{k,1}, \dots, v_{k,N_k}\}$ , and its adjacency matrix is denoted as  $A_{G_k}$ . The coarsened graph  $G_{\text{coar}}$  of  $G$  can be constructed by regarding the clusters  $G_k, k = 1, \dots, K$  as super-nodes and connections between two super-nodes as edges. For  $G_k$ , its sampling matrix  $C_k$  of size  $(N \times N_k)$  is defined by

$$C_k(s, t) = \begin{cases} 1 & \text{if node } v_{k,s} \text{ in } G_k \text{ is identical to vertex } v_t \text{ in } G \\ 0 & \text{otherwise} \end{cases}$$

On one hand,  $C_k$  can be used to down-sample a 1-dimensional graph signal  $x$  on  $G$  to obtain an contracted graph signal  $x_{G_k}$  on  $G_k$ , i.e.  $x_{G_k} = C_k^T x$ . On the other hand,  $C_k$  can also be used to up-sample a graph signal  $x_{G_k}$  on  $G_k$  to obtain a dilated graph  $G$ , i.e.  $x = C_k x_{G_k}$ . Furthermore, the adjacency matrix  $A_{G_k}$  of  $G_k$  can be computed by

$$A_{G_k} = C_k^T A_G C_k.$$

The intra-subgraph adjacency matrix of  $G$  is computed by

$$A_{\text{intra}} = \sum_{k=1}^K C_k A_{G_k} C_k^T.$$

Thereby, the inter-subgraph adjacency matrix of  $G$  can be computed by  $A_{\text{inter}} = A_G - A_{\text{intra}}$ . Let  $M_{\text{coar}} \in \mathbb{R}^{N \times K}$  denote the assignment matrix from  $G$  to  $G_{\text{coar}}$ . Its  $(j, k)$ -th entry is defined as

$$M_{\text{coar}}[j, k] = \begin{cases} 1 & \text{if } v_j \text{ in } G \text{ is grouped into } G_k \text{ in } G_{\text{coar}} \\ 0 & \text{otherwise} \end{cases}$$

As a result, the adjacency matrix  $A_{\text{coar}}$  of the coarsened graph  $G_{\text{coar}}$  is computed by

$$A_{\text{coar}} = M_{\text{coar}}^T A_{\text{inter}} M_{\text{coar}}.$$

In fact,  $A_{\text{coar}}$  can be written as  $A_{\text{coar}} = f(M_{\text{coar}}^T A_G M_{\text{coar}})$  as well, where  $f(\tilde{a}_{i,j}) = 1$  if  $\tilde{a}_{i,j} > 0$  and  $f(\tilde{a}_{i,j}) = 0$  otherwise. As stated previously,  $X$  is a  $d$ -dimensional graph signal on  $G$ . Then, a  $d$ -dimensional graph signal  $X_{\text{coar}}$  on  $G_{\text{coar}}$  can be computed by  $X_{\text{coar}} = M_{\text{coar}}^T X$ . EigenPooling [82] employs spectral clustering to obtain the coarsened graph, and then up-sample the Fourier basis of subgraphs  $G_k, k = 1, \dots, K$ . These Fourier basis are then organized into pooling operators with regard to ascending eigenvalues. Consequently, the pooled node feature matrix is obtained via concatenating the pooled results. The literature [42] proposes a novel deep Hierarchical Graph Convolutional Network (H-GCN) consisting of graph coarsening layers and graph refining layers. They employ both the structural equivalence grouping and structural similarity grouping to generate coarsened graphs. Suppose H-GCN has  $L$  layers. These two kinds of layers are symmetric in the sense that  $A^{(l)}$  is identical to  $A^{(L-l+1)}$ . The graph coarsening layers employ structural equivalence grouping and structural similarity grouping to compute the coarsened graph, and the graph refining layers aims to restore the original topological structure of the corresponding graph. H-GCN employs multi-channel mechanism to explore features in different subspaces, i.e.  $X^{(l)} = \sum_{c=1}^C w_c X_c^{(l)}$ .

**(2) Learning a soft cluster assignment.** STRUCTPOOL [390], as a structured graph pooling technique, regards the graph pooling as a graph clustering problem requiring the learning of a cluster assignment matrix via the feature matrix  $X$  and adjacency matrix  $A_G$ . Learning the cluster assignment matrix can be formulated as a Conditional Random Field (CRF) [391] based probabilistic inference. Specifically, the input feature matrix  $X$  is treated as global observation, and  $Y = \{Y_1, \dots, Y_N\}$  is a random field where  $Y_i \in \{1, \dots, K\}$  is a random variable indicating which clusters the node  $v_i$  is assigned to. As a result,  $(Y, X)$  can be characterized by a CRF model, i.e.

$$\begin{aligned} \mathbb{P}(Y|X) &= \frac{1}{Z(X)} \exp(-\mathcal{E}(Y|X)) \\ &= \frac{1}{Z(X)} \exp\left(\sum_{C \in \mathcal{C}_G} \psi_C(Y_C|X)\right) \end{aligned}$$

where  $\mathcal{E}(Y|X) = \sum_{C \in \mathcal{C}_G} \psi_C(Y_C|X)$  is called an energy function,  $\mathcal{C}_G$  is a set of cliques,  $\psi_C(Y_C|X)$  is a potential function and  $Z(X)$  is a partition function. The energy function  $\mathcal{E}(Y|X)$  can usually be characterized by an unary energy  $\psi_u(\cdot)$  and a pairwise energy  $\psi_p(\cdot, \cdot)$ , i.e.

$$\mathcal{E}(Y|X) = \sum_{s=1}^N \psi_u(y_s|X) + \sum_{s \neq t} \psi_p(y_s, y_t|X) a_{s,t}^l,$$

where  $a_{s,t}^l$  denotes the  $(s, t)$ -th entry of the  $l$ -hop adjacency matrix  $A_G^l$ . The unary energy matrix  $\Psi_u = (\psi_u(y_s|X))_{N \times K}$  can be obtained by a GCNN taking the global observation  $X$  and the adjacency  $A_G$  as input. The pairwise energy matrix  $\Psi_p = (\psi_p(y_s, y_t|X))_{K \times K}$  can be obtained by

$$\psi_p(y_s, y_t|X) = \mu(y_s, y_t) \frac{x_s^T x_t}{\sum_{j \neq s} x_s^T s_j},$$

where  $\mu(y_s, y_t)$  is a learnable compatibility function. Minimizing the energy function  $\mathcal{E}(Y|X)$  via mean-field approximation results in the most probable cluster assignment matrix  $M$  for a give graph  $G$ . As a result, we obtain a new graph  $A_{\text{coar}} = f(M^T A_G M)$  and  $X_{\text{coar}} = M^T X$ . DIFFPOOL [94] is a differentiable graph pooling operator which can generate hierarchical representations of graphs and can be incorporated into various GCNNs in an end-to-end fashion. It maps an adjacency matrix  $A_{G^{(l)}}$  and embedding matrix  $Z^{(l)}$  at the  $l$ -th layer to a new adjacency matrix  $A_{G^{(l+1)}}$  and a coarsened feature matrix  $X^{(l+1)}$ , i.e.

$$(A_{G^{(l+1)}}, X^{(l+1)}) = \text{DIFFPOOL}(A_{G^{(l)}}, Z^{(l)}).$$

More specifically,

$$\begin{aligned} X^{(l+1)} &= (M^{(l)})^T Z^{(l)}, \\ A_{G^{(l+1)}} &= (M^{(l)})^T A_{G^{(l)}} M^{(l)}. \end{aligned}$$

Note that the assignment matrix  $M^{(l)}$  and embedding matrix  $Z^{(l)}$  are respectively computed by two separate GCNNs, namely embedding GCNN and pooling GCNN, i.e.

$$\begin{aligned} Z^{(l)} &= \text{GCNN}_{\text{embed}}(A_{G^{(l)}}, X^{(l)}), \\ M^{(l)} &= \text{Softmax}(\text{GCNN}_{\text{pool}}(A_{G^{(l)}}, X^{(l)})). \end{aligned}$$

**(3) Selecting the first  $k$  top-rank nodes.** The literature [86] proposes a novel Self-Attention Graph Pooling operator (abbreviated as SAGPool). Specifically, SAGPool firstly employs the GCN [6] to calculate the self-attention scores, and then invokes the top-rank function to select the top  $[kN]$  node indices, i.e.

$$\begin{aligned} Z &= \rho\left(\tilde{D}_G^{-\frac{1}{2}} \tilde{A}_G \tilde{D}_G^{-\frac{1}{2}} X \Theta\right), & \text{idx} &= \text{top-rank}(Z, [kN]), \\ Z_{\text{mask}} &= Z_{\text{idx}}, & X' &= X_{\text{idx}}, \\ X_{\text{out}} &= X' \otimes Z_{\text{mask}}, & A_{\text{out}} &= A_{\text{idx, idx}}. \end{aligned}$$

As a result, the selected top- $\lceil kN \rceil$  node indices are employed to extract the output adjacency matrix  $A_{\text{out}}$  and feature matrix  $X_{\text{out}}$ . In order to exploit the expressive power of an encoder-decoder architecture like U-Net [72], the literature [95] proposes a novel graph pooling (gPool) layer and a graph unpooling (gUnpool) layer. The gPool adaptively selects top- $k$  ranked node indices by the down-sampling technique to form a coarsened graph ( $A_{\text{coar}} \in \mathbb{R}^{N \times N}$  and  $X_{\text{coar}} \in \mathbb{R}^{N \times d}$ ) based on scalar projection values on a learnable projection vector, i.e.

$$\begin{aligned} y &= \frac{Xp}{\|p\|}, & \text{idx} &= \text{top-rank}(y, k) \\ \tilde{y} &= \tanh(y_{\text{idx}}), & \tilde{X}_{\text{coar}} &= X_{\text{idx},:} \\ A_{\text{coar}} &= A_{\text{idx},\text{idx}}, & X_{\text{coar}} &= \tilde{X}_{\text{coar}} \otimes (\tilde{y}\mathbf{1}_C^T), \end{aligned}$$

where  $y \in \mathbb{R}^d$ . The gUnpool performs the inverse operation of the gPool layer so as to restore the coarsened graph into its original structure. To this end, gUnpool records the locations of nodes selected in the corresponding gPool layer, and then restores the selected nodes to their original positions in the graph. Specifically, we have

$$X_{\text{refine}} = \text{Distribute}(0_{N \times d}, X_{\text{coar}}, \text{idx}),$$

where the function  $\text{Distribute}(\cdot, \cdot, \cdot)$  distributes row vectors in  $X_{\text{coar}}$  into  $0_{N \times d}$  feature matrix according to the indices  $\text{idx}$ . Note that row vectors of  $X_{\text{refine}}$  with indices in  $\text{idx}$  are updated by the ones in  $X_{\text{coar}}$ , whereas other row vectors remain zero. The literature [96] adopts the similar pooling strategy as gPool to learn the hierarchical representations of nodes. The global representation of the input graph can be computed by

$$\begin{aligned} u &= \sum_{l=1}^L u^{(l)} \\ &= \sum_{l=1}^L \left( \frac{1}{N^{(l)}} \sum_{j=1}^{N^{(l)}} x_j^{(l)} \bowtie \max_j x_j^{(l)} \right), \end{aligned}$$

where  $x_j^{(l)}$  denotes the feature vector of the node  $v_j$  in the graph  $G^{(l)}$  with  $N^{(l)}$  nodes.

### 3.3 Graph Attention Mechanisms

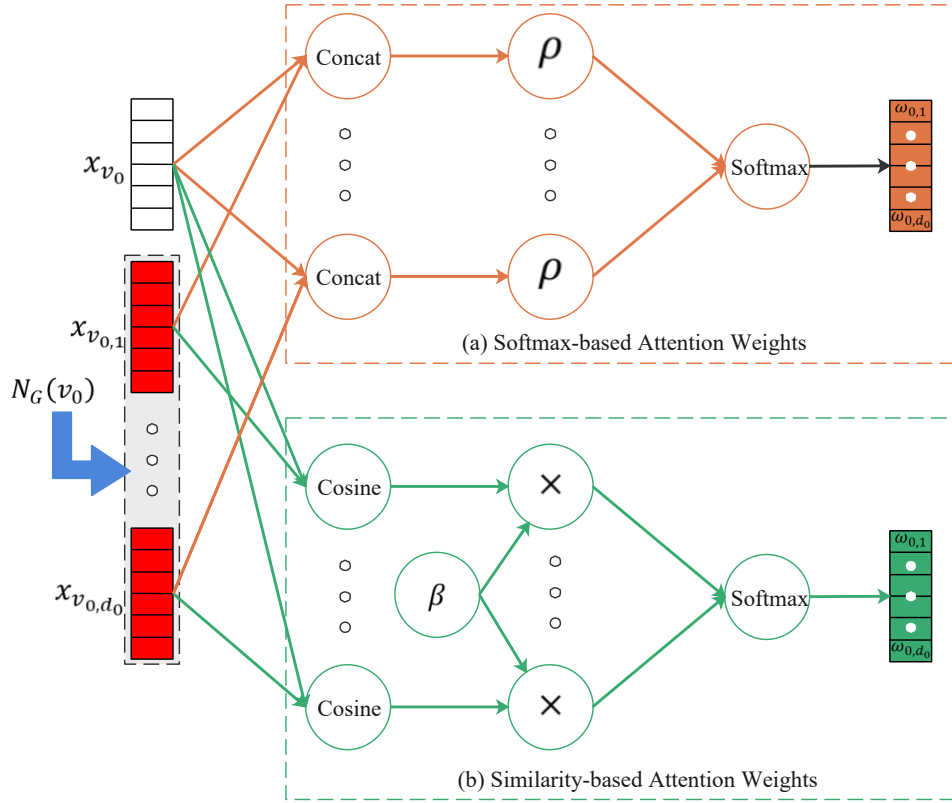
Attention mechanisms, firstly introduced in the deep learning community, guide deep learning models to focus on the task-relevant part of its inputs so as to make precise predictions or inferences [73, 392, 393]. Recently, applying the attention mechanisms to GCNNs has gained considerable attentions so that various attention techniques have been proposed. Below, we summarize the existing attention techniques on graphs and categorize them into 4 classes [394], namely softmax-based graph attention, similarity-based graph attention, spectral graph attention and attention-guided walk. Without loss of generality, the neighbors of a given node  $v_0$  in  $G$  are denoted as  $v_1, \dots, v_{d_0}$ , and their current feature vectors are respectively denoted as  $x_0, x_1, \dots, x_{d_0}$ , where  $d_0 = d_G(v_0)$ .

**Softmax-based Graph Attention.** The softmax-based graph attention is typically implemented by employing a softmax with learnable weights [66, 110] to measure the relevance of  $v_j, j = 1, \dots, d_G(v_0)$  to  $v_0$ . More specifically, the softmax-based attention weights between  $v_0$  and  $v_j$  can be defined as

$$\begin{aligned} \omega_{0,j} &= \text{Softmax}([e_{0,1}, \dots, e_{0,d_G(v_0)}]) \\ &= \frac{\exp(\rho(a^T(Wx_0 \bowtie Wx_j)))}{\sum_{k=1}^{d_G(v_0)} \exp(\rho(a^T(Wx_0 \bowtie Wx_k)))}, \end{aligned} \quad (17)$$

where  $e_{0,j} = \exp(\rho(a^T(Wx_0 \bowtie Wx_j)))$ ,  $a$  is a learnable attention vector and  $W$  is a learnable weight matrix, see Fig. 4(a). As a result, the new feature vector of  $v_0$  can be updated by

$$x'_0 = \rho \left( \sum_{j=1}^{d_G(v_0)} \omega_{0,j} Wx_j \right). \quad (18)$$



**Figure 4** Two kinds of graph attention mechanisms.

In practice, multi-head attention mechanisms are usually employed to stabilize the learning process of the single-head attention [66]. Suppose there are  $H$  independent attention strategies executing the update of Eq. (18). The updated feature vectors, respectively denoted by  $x_0^{(h)} = \rho\left(\sum_{j=1}^{d_G(v_0)} \omega_{0,j}^{(h)} W^{(h)} x_j\right)$ ,  $h = 1, \dots, H$ , can be combined by two kinds of functions, namely concatenation and average. The concatenation based multi-head attention is computed by

$$\begin{aligned} x'_0 &= \bowtie_{h=1}^H x_0^{(h)} \\ &= \bowtie_{h=1}^H \rho\left(\sum_{j=1}^{d_G(v_0)} \omega_{0,j}^{(h)} W^{(h)} x_j\right). \end{aligned}$$

The average based multi-head attention is computed by

$$x'_0 = \rho\left(\frac{1}{H} \sum_{h=1}^H \sum_{j=1}^{d_G(v_0)} \omega_{0,j}^{(h)} W^{(h)} x_j\right).$$

The conventional multi-head attention mechanism treats all the attention heads equally so that feeding the output of an attention that captures a useless representation maybe mislead the final prediction of the model. The literature [74] computes an additional soft gate between 0 and 1 to assign different weights to heads, and gets the formulation of the gated multi-head attention mechanism. The Graph Transformer (GTR) [110] can capture long-range dependencies of dynamic graphs with softmax-based attention mechanism by propagating features within the same graph structure via an intra-graph message passing. The Graph-BERT [367] is essentially a pre-training method only based on the graph attention mechanism without any graph convolution or aggregation operators. Its key component is called a graph transformer based encoder, i.e.

$$\begin{aligned} X^{(l+1)} &= \text{G-Transformer}(X^{(l)}) \\ &= \text{Softmax}\left(\frac{QK^T}{\sqrt{d_h}}\right)V, \end{aligned}$$

where  $Q = X^{(l)}W_Q^{(l+1)}$ ,  $K = X^{(l)}W_K^{(l+1)}$  and  $V = X^{(l)}W_V^{(l+1)}$ . The Graph2Seq [106] is a general end-to-end graph-to-sequence neural encoder-decoder model converting an input graph to a sequence of vectors with the attention based LSTM model. It is composed of a graph encoder, a sequence decoder and a node attention mechanism. The sequence decoder takes outputs (node and graph representations) of the graph encoder as input, and employs the softmax-based attention to compute the context vector sequence.

**Similarity-based Graph Attention.** The similarity-based graph attention depends on the cosine similarities of the given node  $v_0$  and its neighbors  $v_j, j = 1, \dots, d_G(v_0)$ . More specifically, the similarity-based attention weights are computed by

$$\omega_{0,j} = \frac{\exp(\beta \cdot \cos(Wx_0, Wx_j))}{\sum_{k=1}^{d_G(v_0)} \exp(\beta \cdot \cos(Wx_0, Wx_k))}, \quad (19)$$

where  $\beta$  is learnable bias and  $W$  is a learnable weight matrix, see Fig. 4(b). It is well known that  $\cos(x, y) = \frac{\langle x, y \rangle}{\|x\|_2 \|y\|_2}$ . Attention-based Graph Neural Network (AGNN) [67] adopts the similarity-based attention to construct the propagation matrix  $P^{(l)}$  capturing the relevance of  $v_j$  to  $v_i$ . As a result, the output hidden representation  $X^{(l+1)}$  at the  $(l+1)$ -th layer is computed by

$$X^{(l+1)} = \rho \left( P^{(l)} X^{(l)} W^{(l)} \right),$$

where  $P^{(l)}(i, j) \triangleq \omega_{i,j}$  is defined in formula (19).

**Spectral Graph Attention.** The Spectral Graph Attention Networks (SpGAT) aims to learn representations for different frequency components regarding weighted filters and graph wavelet bases [362]. The eigenvalues of the normalized graph Laplacian  $\bar{L}_G$  can be treated as frequencies on the graph  $G$ . As stated in the Preliminary section,  $0 = \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_N = \lambda_{\max}$ . The SpGAT firstly extracts the low-frequency component  $B_L = \{u_1, \dots, u_n\}$  and the high-frequency component  $B_H = \{u_{N-n+1}, \dots, u_N\}$  from the graph Fourier bases  $\{u_1, u_2, \dots, u_N\}$ . So, we have

$$\begin{aligned} X_L &= X^{(l)} \Theta_L, & X_H &= X^{(l)} \Theta_H \\ X^{(l+1)} &= \rho \left( \text{AGGREGATE} \left( B_L F_L B_L^T X_L, B_H F_H B_H^T X_H \right) \right), \end{aligned} \quad (20)$$

where  $F_L$  and  $F_H$  are respectively used to measure the importance of the low- and high-frequency. In practice, we exploit a re-parameterization trick to accelerate the training. More specifically, we replace  $F_L$  and  $F_H$  respectively with the learnable attention weights  $\Omega_L = \text{diag}(\omega_L, \dots, \omega_L)$  and  $\Omega_H = \text{diag}(\omega_H, \dots, \omega_H)$  so as to reduce the number of learnable parameters. To ensure that  $\omega_L$  and  $\omega_H$  are positive and comparable, we normalize them by the softmax function, i.e.

$$\omega_L = \frac{\exp(\omega_L)}{\exp(\omega_L) + \exp(\omega_H)}, \quad \omega_H = \frac{\exp(\omega_H)}{\exp(\omega_L) + \exp(\omega_H)}.$$

In addition to the attention weights, another important issue is how to choose the low- and high-frequency components  $B_L$  and  $B_H$ . A natural choice is to use the graph Fourier bases, yet the literatures [191, 194] conclude that utilizing the spectral graph wavelet operators can achieve better embedding results than the graph Fourier bases. Therefore, we substitute  $B_L$  and  $B_H$  in Eq. (14) for the spectral graph wavelet operator  $\Psi_{L,g}^s$  and  $\Psi_{H,g}^s$ , i.e.

$$X^{(l+1)} = \rho \left( \text{AGGREGATE} \left( \left( \Psi_{L,g}^s \right) F_L \left( \Psi_{L,g}^s \right)^{-1} X_L, \left( \Psi_{H,g}^s \right) F_H \left( \Psi_{H,g}^s \right)^{-1} X_H \right) \right).$$

**Attention-guided Walk.** The two aforementioned kinds of attention mechanisms focuses on incorporating task-relevant information from the neighbors of a given node into the updated representations of the pivot. Here, we introduce a new attention mechanism, namely attention-guided walk [81], which has different purpose from the softmax- and similarity-based attention mechanisms. Suppose a walker walks along the edges of the graph  $G$  and he currently locates at the node  $v_i$ . The hidden representation



$x^{(t)}$  of  $v_t$  is computed by a recurrent neural network  $f_x(s^{(t)}, x^{(t-1)}; \Theta_x)$  taking the step embedding  $s^{(t)}$  and internal representation of the historical information from the previous step  $x^{(t-1)}$  as input, i.e.

$$x^{(t)} = f_x(s^{(t)}, x^{(t-1)}; \Theta_x).$$

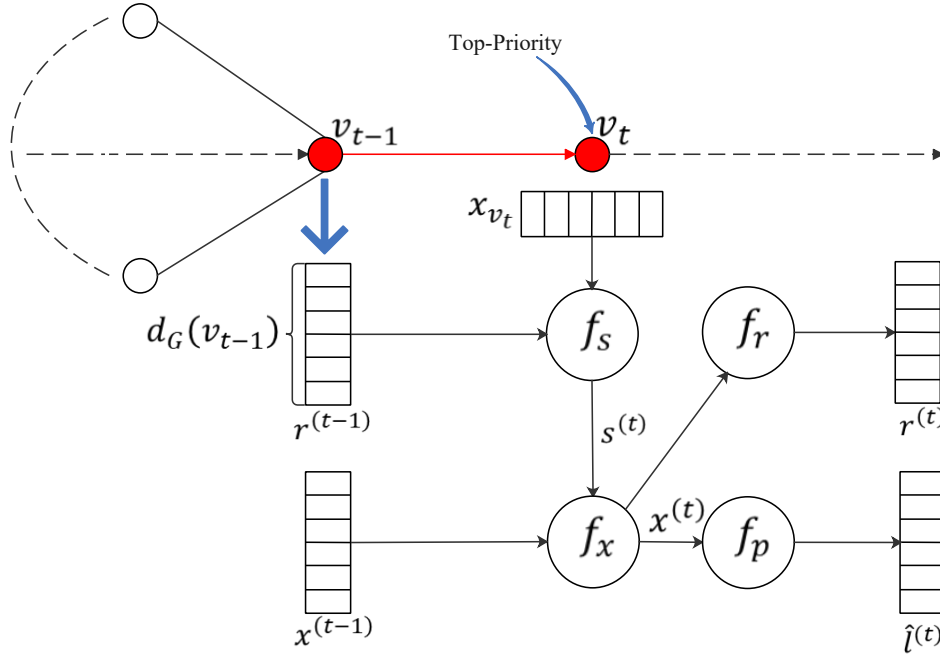
The step embedding  $s^{(t)}$  is computed by a step network  $f_s(r^{(t-1)}, x_{v_t}; \Theta_s)$  taking the ranking vector  $r^{(t-1)}$  and the input feature vector  $x_{v_t}$  of the top-priority node  $v_t$  as input, i.e.

$$s^{(t)} = f_s(r^{(t-1)}, x_{v_t}; \Theta_s).$$

The hidden representation  $x^{(t)}$  is then fed into a ranking network  $f_r(x^{(t)}; \Theta_r)$  and a predicting network  $f_p(x^{(t)}; \Theta_p)$ , i.e.

$$\begin{aligned} r^{(t)} &= f_r(x^{(t)}; \Theta_r), \\ \hat{l}^{(t)} &= f_p(x^{(t)}; \Theta_p). \end{aligned}$$

The ranking network  $f_r(x^{(t)}; \Theta_r)$  determines which neighbors of  $v_t$  should be prioritized in the next step, and the predicting network  $f_p(x^{(t)}; \Theta_p)$  makes a prediction on graph labels. Now,  $x^{(t)}$  and  $r^{(t)}$  are fed into the next node to compute its hidden representations. Fig. 5 shows the computational framework of the attention-guided walk.



**Figure 5** The computational framework of the attention-guided walk.

### 3.4 Graph Recurrent Neural Networks

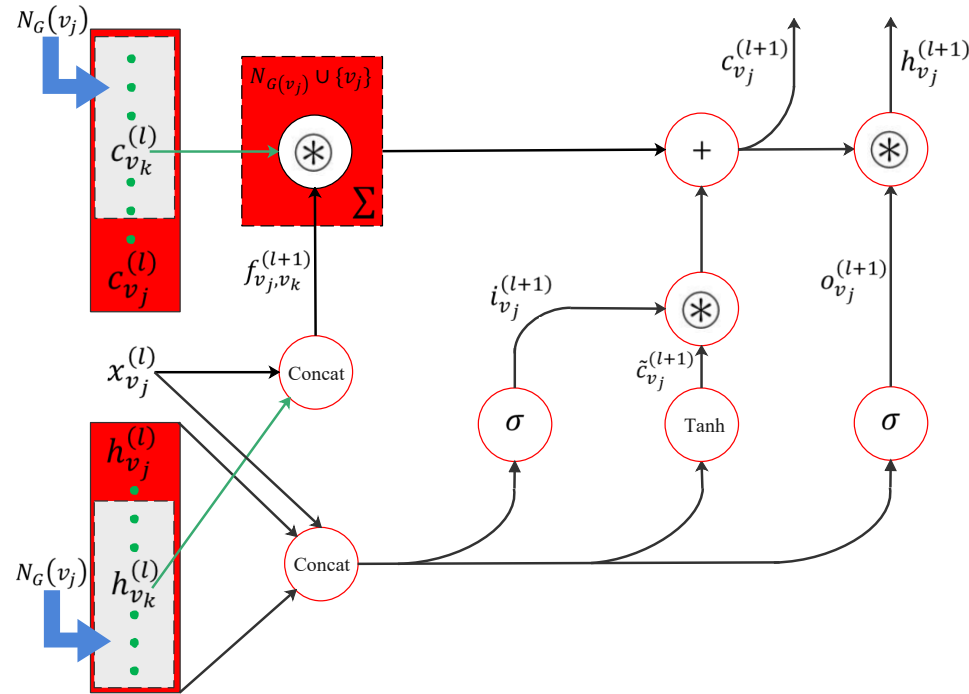
In addition to the GCNNs, there are another classic GNN architectures, namely Graph Recurrent Neural Networks (GRNNs), which generalize the Recurrent Neural Networks (RNNs) to the graphs. In general, the GRNN can be formulated by Eq.  $h_{v_j}^{(l+1)} = \text{GRNN}(x_{v_j}^{(l)}, \{h_{v_k}^{(l)} : v_k \in N_G(v_j) \cup \{v_j\}\})$ . Below, we introduce the GRNN architectures in details.

**Graph LSTM.** The Graph Long Short Term Memroy (Graph LSTM) [117–121, 123, 143] is a generalization to the vanilla LSTM from the sequential data or multi-dimensional data to general graph-structured data. Specifically, the graph LSTM updates the hidden states and cell states of nodes by the

following formula,

$$\begin{aligned}
 i_{v_j}^{(l+1)} &= \sigma \left( W_i x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_i h_{v_k}^{(l)} + b_i \right), \\
 o_{v_j}^{(l+1)} &= \sigma \left( W_o x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_o h_{v_k}^{(l)} + b_o \right) \\
 \tilde{c}_{v_j}^{(l+1)} &= \tanh \left( W_c x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_c h_{v_k}^{(l)} + b_c \right), \\
 f_{v_j, v_k}^{(l+1)} &= \sigma \left( W_f x_{v_j}^{(l)} + U_f h_{v_k}^{(l)} + b_f \right) \\
 c_{v_j}^{(l+1)} &= i_{v_j}^{(l+1)} \otimes \tilde{c}_{v_j}^{(l+1)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} f_{v_j, v_k}^{(l+1)} \otimes c_{v_k}^{(l)}, \\
 h_{v_j}^{(l+1)} &= o_{v_j}^{(l+1)} \otimes \tanh(c_{v_j}^{(l+1)}).
 \end{aligned}$$

see the Fig. 6. The literature [267] develops a general framework, named structure-evolving LSTM, for learning interpretable data representations via the graph LSTM. It progressively evolves the multi-level node representations by stochastically merging two adjacent nodes with high compatibilities estimated by the adaptive forget gate of the graph LSTM. As a result, the new graph is produced with a Metropolis-Hastings sampling method. The Gated Graph Sequence Neural Networks (GGs-NNs) [122] employs Gated Recurrent Unit (GRU) [115] to modify the vanilla GCNNs so that it can be extended to output sequence.



**Figure 6** Computational framework of the Graph LSTM.

**GRNNs for Dynamic Graphs.** A dynamic graph is the one whose structure (i.e. adding a node, removing a node, adding an edge, removing an edge) or node/edge representations change over time. The GRNN models provide a powerful tool to tackle dynamic graphs. The Dynamic Graph Neural Network (DGNN) model [107] can model the dynamic information as nodes or edges adding. More specifically, the DGNN is composed of two key components: an update component and a propagation component. Suppose an edge  $(v_s, v_g, t)$  is added to the input dynamic graph at time  $t$ . Let  $t-$  denotes a time before the time  $t$ . The update component consists of three sequential units: the interacting unit, the S- or G-update unit and the merging unit. The interacting unit takes the source and target representations before the time  $t$  as input, and outputs the joint representation of the interaction, i.e.  $x_e^{(t)} = \rho \left( W_s x_s^{(t-)} + W_g x_g^{(t-)} + b_e \right)$ . The S- and G-update units employ the LSTM [160] to respectively

update the cell states and hidden states of the source and target, i.e.

$$\begin{aligned} (C_{v_s}^{(t)}, h_{v_s}^{(t)}) &= \text{LSTM}_s (C_{v_s}^{(t-)}, h_{v_s}^{(t-)}, \Delta t_s), \\ (C_{v_g}^{(t)}, h_{v_g}^{(t)}) &= \text{LSTM}_s (C_{v_g}^{(t-)}, h_{v_g}^{(t-)}, \Delta t_g). \end{aligned}$$

The merging unit adopts the similar functions to the interacting unit to respectively merge  $h_{v_s}^{(t)}$  and  $h_{v_s}^{(t-)}$ , and  $h_{v_g}^{(t)}$  and  $h_{v_g}^{(t-)}$ . The propagation component can propagate information from two interacting nodes ( $v_s$  and  $v_g$ ) to influenced nodes (i.e. their neighbors). It also consists of three units: the interacting unit, the propagation unit and the merge unit, which are defined similarly to the update component except that they have different learnable parameters. The literature [105] addresses the vertex- and graph-focused prediction tasks on dynamic graphs whose vertex set does not change over time by combining GCNs, LSTMs and fully connected layers. The Variational Graph Recurrent Neural Networks (VGRNNs) [113] is essentially a variational graph auto-encoder whose encoder integrates the GCN and RNN into a graph RNN (GRNN) framework and decoder is a joint probability distribution of a multi-variate Gaussian distribution and Bernoulli Distribution. Note that the semi-implicit variational inference is employed to approximate the posterior in practice so as to generate the node embedding.

**GRNNs based on Vanilla RNNs.** The GRNNs based on vanilla RNNs firstly employ random walk techniques or traversal methods, e.g. Breadth-First Search (BFS) and Depth-First Search (DFS), to obtain a collection of node sequences, and then leverage RNN model, e.g. LSTM and GRU, to capture the long short-term dependencies. The literature [116] performs joint random walks on attributed networks, and utilizes them to boost the node representation learning. The proposed framework is called GraphRNA consisting of two key components, namely a collaborative walking mechanism AttriWalk and a tailored deep embedding architecture for joint random walks, named Graph Recurrent Network (GRN). Suppose  $\mathcal{A}_G$  denotes the node-attribute matrix of size  $N \times M$ . The AttriWalk admits the transition matrix of size  $\mathbb{R}_+^{(N+M) \times (N+M)}$  which is written as

$$\mathcal{T} = \begin{bmatrix} \alpha \mathcal{A}_G & (1 - \alpha) \mathcal{A}_G \\ (1 - \alpha) \mathcal{A}_G^T & 0 \end{bmatrix}.$$

After obtaining the sequences via the collaborative random walk, the bi-directional GRU [115] and pooling operator are employed to learn the global representations of sequences. The literature [111] leverages the BFS node ordering and truncation strategy to obtain a collection of node representation sequences, and then uses the GRU model and variational auto-regression regularization to perform the graph classification.

## 4 Extensions and Applications

The aforementioned architectures are essentially ingredients of constructing the GNNs. Below, we investigate the directions and tasks of the GNNs from the next 7 aspects: capability and interpretability, deep graph representation learning, deep graph generative models, combinations of the PI and GNNs, adversarial attacks for the GNNs, graph neural architecture search and graph reinforcement learning.

### 4.1 Capability and Interpretability

The GCNNs have achieved tremendous empirical successes over the supervised, semi-supervised and unsupervised learning on graphs. Recently, many studies start to put an eye on the capability and interpretability of the GCNNs.

**Capability.** The capability of the GCNNs refers to their expressive power. If two graphs are isomorphic, they will obviously output the same representations of nodes/edges/graph. Otherwise, they should output different representations. However, two non-isomorphic graphs maybe output the same representations in practice. This is the theoretical limitations of the GCNNs. As described in the

literatures [83, 88, 360], The 1-hop spatial GCNNs (1-GCNNs) have the same expressive power as the 1-dimensional Weisfeiler-Leman (1-WL) graph isomorphism test in terms of distinguishing non-isomorphic graphs. The 1-WL iteratively update the colors of nodes according to the following formula

$$C_l^{(l+1)}(v) = \text{HASH} \left( C_l^{(l)}(v), \left\{ \left\{ C_l^{(l)}(u) : u \in N_G(v) \right\} \right\} \right).$$

According to the literatures [83, 88], we have that the 1-GCNN architectures do not have more power in terms of distinguishing two non-isomorphic graphs than the 1-WL heuristic. Nevertheless, they have equivalent power if the aggregation and update functions are injective. In order to overcome the theoretical limitations of the GCNNs, the literature [83] proposes a Graph Isomorphism Network (GIN) architecture, i.e.

$$\begin{aligned} A_v^{(l+1)} &= \text{AGGREGATE}^{(l+1)} \left( \left\{ \left\{ X^{(l)}[u, :] : u \in N_G(v) \right\} \right\} \right) \\ &\triangleq \sum_{u \in N_G(v)} X^{(l)}[u, :] \\ X^{(l+1)}[v, :] &= \text{UPDATE}^{(l+1)} \left( X^{(l)}[v, :], A_v^{(l+1)} \right) \\ &\triangleq \text{MLP} \left( (1 + \epsilon^{(l+1)}) X^{(l)}[v, :] + A_v^{(l+1)} \right), \end{aligned}$$

where  $\epsilon^{(l)}$  is a scalar parameter. The literature [163] studies the expressive power of the spatial GCNNs, and presents two results: (1) The spatial GCNNs are shown to be a universal approximator under sufficient conditions on their depth, width, initial node features and layer expressiveness; (2) The power of the spatial GCNNs is limited when their depth and width is restricted. In addition, there are some other studies on the capability of the GCNNs from different perspectives, e.g. the first order logic [164],  $p$ -order graph moments  $M_p(A_G) = \prod_{j=1}^p (A_G \cdot W_j + B_j)$  [167] and algorithmic alignment with the dynamic programming [401].

**Interpretability.** Interpretability plays a vital role in constructing a reliable and intelligent learning systems. Although some studies have started to explore the interpretability of the conventional deep learning models, few of studies put an eye on the interpretability of the GNs [85]. The literature [165] bridges the gap between the empirical success of the GNs and lack of theoretical interpretations. More specifically, it considers two classes of techniques: (1) gradient based explanations, e.g. sensitivity analysis and guided back-propagation; (2) decomposition based explanations, e.g. layer-wise relevance propagation and Taylor decomposition. The GNNExplainer [166] is a general and model-agnostic approach for providing interpretable explanations for any spatial GCNN based model in terms of graph machine learning tasks. Given a trained spatial GCNN model  $\Phi$  and a set of predictions, the GNNExplainer will generate a single-instance explanation by identifying a subgraph of the computation graph and a subset of initial node features, which are the most vital for the prediction of the model  $\Phi$ . In general, the GNNExplainer can be formulated as the following optimization problem

$$\max_{G_S, X_S^F} I(Y, (G_S, X_S^F)) = H(Y) - H(Y|G = G_S, X = X_S^F), \quad (21)$$

where  $I(\cdot, \cdot)$  denotes the mutual information of two random variables,  $G_S$  is a small subgraph of the computation graph and  $X_S^F$  is a small subset of node features  $\{X^F[j, :] : v_j \in G_S\}$ . The entropy term  $H(Y)$  is constant because the spatial GCNN model  $\Phi$  is fixed. Therefore, Eq. (21) is equivalent to minimizing the conditional entropy  $H(Y|G = G_S, X = X_S^F)$ . In order to improve the tractability and computational efficiency of the GNNExplainer, the final optimization framework is reformulated as

$$\min_{M, F} - \sum_{c=1}^C \mathbb{I}[y = c] \log P_{\Phi}(Y = y|G = A_G \otimes \sigma(M), X = X_S^F).$$

In addition, the GNNExplainer also provides multi-instances explanations based on graph alignments and prototypes so as to answer questions like "How did a GCNN predict that a given set of nodes all have label  $c$ ?"

## 4.2 Deep Graph Representation Learning

Graph representation learning (or called network embedding) is a paradigm of unsupervised learning on graphs. It gains a large amount of popularity since the DeepWalk was published in 2014 [154]. Subsequently, many studies exploit deep learning techniques to learn low-dimensional representations of nodes [139]. In general, the network embedding via the vanilla deep learning techniques learn low-dimensional feature vectors of nodes by utilizing either stacked auto-encoders to reconstruct the adjacent or positive point-wise mutual information features [124, 132, 133, 136, 141, 142] or RNNs to capture long and short-term dependencies of node sequences yielded by random walks [134, 135]. In the following, we introduce the network embedding approaches based on GNNs.

**Network Embedding based on GNNs.** In essence, the GNNs provides an elegant and powerful framework for learning node/edge/graph representations. The majority of the GCNNs and GRNNs are concerned with semi-supervised learning (i.e. node-focused tasks) or supervised learning (i.e. graph-focused) tasks. Here, we review the GNN based unsupervised learning on graphs. In general, the network embedding based on GNNs firstly utilize the GCNNs and variational auto-encoder to generate gaussian-distributed hidden states of nodes, and then reconstruct the adjacency matrix and/or the feature matrix of the input graph [87, 127, 129, 131, 138, 148]. A representative approach among these ones is the Variational Graph Auto-Encoder (VGAE) [87] consisting of a GCN based encoder and an inner product decoder. The GCN based encoder is defined to be

$$\begin{aligned} q(Z|X, A_G) &= \prod_{j=1}^N q(Z[j, :]|X, A_G) \\ &= \prod_{j=1}^N \mathcal{N}(Z[j, :]|\mu_j, \text{diag}(\sigma_j^2)), \end{aligned}$$

where  $\mu_j = \text{GCN}_\mu(X, A_G)$  and  $\log \sigma_j = \text{GCN}_\sigma(X, A_G)$ . The inner product decoder is defined to be

$$\begin{aligned} p(A_G|Z) &= \prod_{j=1}^N \prod_{k=1}^N p(A_G[j, k]|Z[j, :], Z[k, :]) \\ &= \prod_{j=1}^N \prod_{k=1}^N \sigma(Z[j, :]Z[k, :]^T). \end{aligned}$$

They adopt the evidence lower bound [151] as their objective function. The adversarially regularized (variational) graph autoencoder [129] extends the VGAE by adding an adversarial regularization term to the evidence lower bound. The literature [128] proposes a symmetric graph convolutional autoencoder which produces a low-dimensional latent nodes representations. Its encoder employs the Laplacian smoothing in Eq. (9) to jointly encode the structural and attributed information, and its decoder is designed based on Laplacian sharpening as the counterpart of the Laplacian smoothing of the encoder. The Laplacian sharpening is defined to be

$$\begin{aligned} X^{(l+1)} &= (1 + \gamma)X^{(l)} - \gamma D^{-1}AX^{(l)} \\ &= X^{(l)} + \gamma(\mathbb{I}_N - D^{-1}A)X^{(l)}, \end{aligned}$$

which allows utilizing the graph structure in the whole processes of the proposed autoencoder architecture. In addition, there are some other methods to perform the unsupervised learning on graphs, which do not rely on the reconstruction of the adjacency and/or feature matrix, e.g. the graph auto-encoder on directed acyclic graphs [79], pre-training GNNs via context prediction and attribute masking strategies [130], and deep graph Infomax using a noise-contrastive type objective with a standard binary cross-entropy loss between positive examples and negative examples [140].

## 4.3 Deep Graph Generative Models

The aforementioned work concentrates on embedding an input graph into a low-dimensional vector space so as to perform semi-supervised/supervised/unsupervised learning tasks on graphs. This subsection introduces deep graph generative models aiming to mimic real-world complex graphs. Generating complex graphs from latent representations is confronted with great challenges due to high nonlinearity and

arbitrary connectivity of graphs. Note that graph translation [145] is akin to graph generation. However, their difference lies in that the former takes two graphs, i.e. input graph and target graph, as input, and the latter only takes a single graph as input. The NetGAN [154] utilizes the generative adversarial network [149] to mimic the input real-world graphs. More specifically, it is composed of two components, i.e. a generator  $G$  and a discriminator  $D$ , as well. The discriminator  $D$  is modeled as a LSTM in order to distinguish real node sequences, which are yielded by the second-order random walks scheme node2vec [155], from faked ones. The generator  $G$  aims to generate faked node sequences via another LSTM, whose generating process is as follows.

$$\begin{aligned} v_0 = 0, z &\sim N_m(0, \mathbb{I}_m), (C_0, h_0) = g_{\theta'}(z) \\ (C_j, h_j, o_j) &= \text{LSTM}_G(C_{j-1}, h_{j-1}, v_{j-1}), \\ v_j &\sim \text{Cat}(\text{Softmax}(o_j)). \end{aligned}$$

The motif-targeted GAN [147] generalizes random walk based architecture of the NetGAN to characterize mesoscopic context of nodes. Different from [147, 154], the GraphVAE [146] adopts a probabilistic graph decoder to generate a probabilistic fully-connected graph, and then employs approximate graph matching to reconstruct the input graph. Its reconstruction loss is the cross entropy between the input and reconstructed graphs. The literature [150] defines a sequential decision-making process to add a node/edge via the graph network [85], readout operator and softmax function. The GraphRNN [152] is a deep autoregressive model, which generates graphs by training on a representative set of graphs and decomposes the graph generation process into a sequence of node and edge formations conditioned on the current generated graph. Specifically, let  $\pi$  be a permutation over  $N$  nodes and  $S_j^\pi \in \{0, 1\}^{j-1}$  be an adjacency vector representing the links between node  $\pi(v_i)$  and the previous nodes  $\{\pi(v_k), k = 1, \dots, j-1\}$ . The GraphRNN is denoted as  $P(S^\pi) = \prod_{j=1}^{N+1} P(S_j^\pi | S_{<j}^\pi)$ , where the conditional probability is computed by a RNN.

#### 4.4 Combinations of the PI and GNNs

The GNNs and PI are two different paradigms for learning from complicated real-world data. The former specializes in learning hierarchical representations based on local and global structural information, and the latter learning the dependencies between random variables. How to combine them becomes an intriguing research topic. This subsection provides a summarization of this topic.

**Conditional Random Field Layer Preserving Similarities between Nodes.** The literature [162] proposes a CRF layer for the GCNNs to enforce hidden layers to preserve similarities between nodes. Specifically, the input to each layer is viewed as a random vector around the output of the last layer. That is,  $X^{(l+1)}$  is a random vector around  $B^{(l+1)} = \text{GCNN}(X^{(l)}, A_G, X)$ . The objective function for the GCNN with a CRF layer can be reformulated as

$$J(W; X, A_G, Y) = \mathcal{L}(Y; B^{(L)}) + \sum_{l=1}^{L-1} \frac{\gamma}{2} \|X^{(l)} - B^{(l)}\|_F^2 + \mathcal{R}(X^{(l)}),$$

where the first term after "=" is the conventional loss function for semi-supervised node classification problem, and the last term is a regularization one implementing similarity constraint.  $\mathcal{R}(X^{(l)})$  is modeled as a CRF, i.e.

$$p(X^{(l)} | B^{(l)}) = \frac{1}{Z(B^{(l)})} \exp\left(-E(X^{(l)} | B^{(l)})\right)$$

where

$$E(X^{(l)} | B^{(l)}) = \sum_{v \in V} \varphi_v(X^{(l)}[v, :], B^{(l)}[v, :]) + \sum_{(u,v) \in E} \varphi_{u,v}(X^{(l)}[u, :], X^{(l)}[v, :], B^{(l)}[u, :], B^{(l)}[v, :]).$$

Let  $s_{u,v}$  denote the similarity between  $u$  and  $v$ . The unary energy component  $\varphi_v(X^{(l)}[v, :], B^{(l)}[v, :]) = \|X^{(l)}[v, :] - B^{(l)}[v, :]\|_2^2$ , and the pairwise energy component  $\varphi_{u,v}(X^{(l)}[u, :], X^{(l)}[v, :], B^{(l)}[u, :], B^{(l)}[v, :]) =$

$s_{u,v} \|X^{(l)}[u, :] - X^{(l)}[v, :]\|_2^2$  implementing the similarity constraint. The mean-field variational Bayesian inference is employed to approximate the posterior  $p(X^{(l)}|B^{(l)})$ . Consequently, the CRF layer is defined as

$$\left(X^{(l)}[v, :]\right)^{(k+1)} = \frac{\alpha B^{(l)}[v, :] + \beta \sum_{u \in N_G(v)} s_{u,v} \left(X^{(l)}[u, :]\right)^{(k)}}{\alpha + \beta \sum_{u \in N_G(v)} s_{u,v}}.$$

**Conditional GCNNs for Semi-supervised Node Classification.** The conditional GCNNs incorporate the Conditional Random Field (CRF) into the conventional GCNNs so that the semi-supervised node classification can be enhanced by both the powerful node representations and the dependencies of node labels. The GMNN [80] performs the semi-supervised node classification by incorporating the GCNN into the Statistical Relational Learning (SRL). Specifically, the SRL usually models the conditional probability  $p(Y_V|X_V)$  with the CRF, i.e.

$$p(Y_V|X_V) = \frac{1}{Z(X_V)} \prod_{(i,j) \in E} \varphi_{i,j}(y_i, y_j, X_V),$$

where  $y_* = Y_V[*, :]$ ,  $* = i, j$ . Note that  $Y_V$  is composed of the observed node labels  $Y_L$  and hidden node labels  $Y_U$ . The variational Bayesian inference is employed to estimate the posterior  $p(Y_U|Y_L, X_V)$ . The objective function is defined as the ELBO, i.e.

$$\text{ELBO}(q_{\theta_v}(Y_U|X_V)) = \mathbb{E}_{q_{\theta_v}(Y_U|X_V)} [\log(p_{\theta_l}(Y_L, Y_U|X_V)) - \log(q_{\theta_v}(Y_U|X_V))].$$

This objective can be optimized by the variational Bayesian EM algorithm [161], which iteratively updates the variational distribution  $q_{\theta_v}(Y_U|X_V)$  and the likelihood  $p_{\theta_l}(Y_U|Y_L, X_V)$ . In the VBE stage,  $q_{\theta_v}(Y_U|X_V) = \prod_{v \in U} q_{\theta_v}(y_v|X_V)$ , and  $q_{\theta_v}(y_v|X_V)$  is approximated by a GCNN. In the VBM stage, the pseudo-likelihood is employed to compute the following formula

$$\begin{aligned} & \mathbb{E}_{q_{\theta_v}(Y_U|X_V)} [\sum_{v \in V} \log p_{\theta_l}(y_n|Y_{V \setminus v}, X_V)] \\ &= \mathbb{E}_{q_{\theta_v}(Y_U|X_V)} [\sum_{v \in V} \log p_{\theta_l}(y_n|Y_{N_G(v)}, X_V)], \end{aligned}$$

and  $p_{\theta_l}(y_n|Y_{N_G(v)}, X_V)$  is approximated by another GCNN. The literature [159] adopts the similar idea to the GMNN. Its posterior is modeled as a CRF with unary energy components and pairwise energy components whose condition is the outputs of the prescribed GCNN. The maximum likelihood estimation employed to estimate the model parameters.

**GCNN-based Gaussian Process Inference.** A Gaussian Process (GP) defines a distribution over a function space and assumes any finite collection of marginal distributions follows a multivariate Gaussian distribution. A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  follows a Gaussian Process  $\text{GP}(m(\cdot), \kappa(\cdot, \cdot))$  if  $(f(X_1), \dots, f(X_N))^T$  follows a  $N$ -dimensional Gaussian distribution  $\mathcal{N}_N(\mu, \Sigma)$ , where  $\mu = (m(X_1), \dots, m(X_N))^T$  and  $\Sigma[j, k] = \kappa(X_j, X_k)$ , for any  $N$   $d$ -dimensional random vectors. For two  $d$ -dimensional random vectors  $X$  and  $X'$ , we have  $\mathbb{E}[f(X)] = m(X)$  and  $\text{Cov}(f(X), f(X')) = \kappa(X, X')$ . Given a collection of  $N$  samples  $\mathcal{D} = \{(X_j, y_j) : j = 1, \dots, N\}$ , the GP inference aims to calculate the probability  $p(y|X)$  for predictions, i.e.

$$f \sim \text{GP}(0(\cdot), \kappa(\cdot, \cdot)), y_j \sim \text{DIST}(\lambda(f(X_j))),$$

where  $\lambda(\cdot)$  is a link function and  $\text{DIST}(\cdot)$  denotes an arbitrary feasible noise distribution. To this end, the posterior  $p(f|\mathcal{D})$  needs to be calculated out firstly. The literature [156] employs amortized variational Bayesian inference to approximate  $p(f|\mathcal{D})$ , i.e.  $f = \mu + L\epsilon$ , and the GCNNs to estimate  $\mu$  and  $L$ .

**Other GCNN-based Probabilistic Inference.** The literature [157] combines the GCNNs and variational Bayesian inference to infer the input graph structure. The GCNNs are employed to estimate parameters of likelihood, and the latter approximates the posterior proportional to the likelihood and adjacency matrix prior. The literature [158] infers marginal probabilities in probabilistic graphical models by incorporating the GCNNs to the conventional message-passing inference algorithm. The literature [168] approximates the posterior in Markov logic networks with the GCNN-enhanced variational Bayesian inference. That is, the GCNNs can capture the dependencies of observed facts on hidden ones.

#### 4.5 Adversarial Attacks for the GNNs

For the classification tasks, it should be noted that in many circumstances where classifiers are deployed, adversaries deliberately contaminate data in order to fake the classifiers [149, 364]. This is the so-called adversarial attacks for the classification problems. As stated previously, the GNNs can solve semi-supervised node classification problems and supervised graph classification problems. Therefore, it is inevitable to study the adversarial attacks for GNNs and defense [365].

**Adversarial Attacks on Graphs.** The literature [182] firstly proposes a reinforcement learning based attack method which can learn a generalizable attack policy, while only requiring prediction labels from the target classifier. This paper provides a definite definition for a graph adversarial attacker. Given a sample  $(G, c, y) \in \{(G_j, c_j, y_j) : j = 1, \dots, M\}$  and a classifier  $f \in \mathcal{F}$ , the graph adversarial attacker  $g : \mathcal{F} \times \mathcal{G} \rightarrow \mathcal{G}$  asks to modify a graph  $G = (V, E)$  into  $\tilde{G} = (\tilde{V}, \tilde{E})$ , such that

$$\begin{aligned} & \max_g \mathbb{I}(f(\tilde{G}, c) \neq y) \\ & \text{s.t. } \tilde{G} = g(f, (G, c, y)), \\ & \mathcal{I}(G, \tilde{G}, c) = 1, \end{aligned}$$

where  $\mathcal{I} : \mathcal{G} \times \mathcal{G} \times V \rightarrow \{0, 1\}$ , named an equivalence indicator, checks whether two graph  $G$  and  $\tilde{G}$  are equivalent under the classification semantics. The equivalence indicator are usually defined in two fashions, namely explicit semantics and small modifications. The explicit semantics are defined as  $\mathcal{I}(G, \tilde{G}, c) = \mathbb{I}(f^*(G, c) = f^*(\tilde{G}, c))$  where  $f^*$  is a gold standard classifier, and the small modifications are defined as

$$\mathcal{I}(G, \tilde{G}, c) = \mathbb{I}\left(|(E - \tilde{E}) \cup (\tilde{E} - E)| < m\right) \cdot \mathbb{I}\left(\tilde{E} \subseteq N(G, b)\right),$$

where  $N(G, b) = \{(u, v) : u, v \in V, d_G(u, v) \leq b\}$ . The attack procedure is modeled as a finite horizon Markov Decision Process (MDP)  $\mathcal{M}_m(f, G, c, y)$ , and therefore can be optimized by Q-learning with a hierarchical Q-function so as to learn the attack policy. For the MDP  $\mathcal{M}_m(f, G, c, y)$ , its action at time step  $t$  is  $a_t \in \mathcal{A} \subseteq V \times V$  allowed to add or delete edges in the graph, its state at time step  $t$  is  $(\hat{G}_t, c)$  where  $\hat{G}_t$  is a partially modified graph with some of the edges added/deleted from  $G$ , and the reward function is defined as

$$R(\tilde{G}, c) = \begin{cases} 1 & f(\tilde{G}, c) \neq y \\ -1 & f(\tilde{G}, c) = y. \end{cases}$$

Note that the GCNNs are employed to parameterize the Q-function. The NETTACK, proposed in the literature [181], considers attacker nodes  $\mathcal{A}$  to satisfy the following constraints, i.e.

$$\begin{aligned} X'_{u,j} & \neq X_{u,j}^{(0)} \implies u \in \mathcal{A}, & \text{(Feature Attack)} \\ A'_{u,v} & \neq A_{u,v}^{(0)} \implies u \in \mathcal{A} \vee v \in \mathcal{A}, & \text{(Structure Attack)} \\ \sum_u \sum_j |X_{u,j}^{(0)} - X'_{u,j}| + \sum_{u < v} |A_{u,v}^{(0)} - A'_{u,v}| & \leq \Delta, & \text{(Equivalence Indicator)} \end{aligned} \tag{22}$$

where  $G'$  is derived by perturbing  $G^{(0)}$ . Let  $\mathcal{P}_{\Delta, \mathcal{A}}^{G^0}$  denote the set of all perturbed graphs  $G'$  fulfilling formulas (22). The goal is to find a perturbed graph  $G' = (A', X')$  that classifies a target node  $v_0$  as  $c_{\text{new}}$  and maximizes the log-probability/logit to  $c_{\text{old}}$ , i.e.

$$\max_{(A', X') \in \mathcal{P}_{\Delta, \mathcal{A}}^{G^0}} \max_{c_{\text{new}} \neq c_{\text{old}}} \log Z_{v_0, c_{\text{new}}}^* - \log Z_{v_0, c_{\text{old}}}^*,$$

where  $Z^* = f_{\theta^*}(A', X')$  with  $\theta^* = \arg \min_{\theta} \mathcal{L}(\theta; A', X')$ . It is noted that the NETTACK, accelerated by incremental computations, employs the GCNNs to model the classifier. The literature [183] adopts similar equivalence indicator to Eq. (22), and poisoning attacks are mathematically formulated as a



bilevel optimization problem, i.e.

$$\begin{aligned}
 & \min_{\tilde{G}} \mathcal{L}_{\text{attack}} \left( f_{\theta^*} \left( \tilde{G} \right) \right) \\
 & \text{s.t. } \tilde{G} \in \mathcal{P}_{\Delta, \mathcal{A}}^{G_0} \\
 & \theta^* = \arg \min_{\theta} \mathcal{L}_{\text{train}} \left( f_{\theta} \left( \tilde{G} \right) \right).
 \end{aligned} \tag{23}$$

This bilevel optimization problem in formula (23) is then tackled using meta-gradients, whose core idea is to treat the adjacency matrix of the input graph as a hyperparameter.

**Defense against the Adversarial Attacks.** A robust GCNN requires that it is invulnerable to perturbations of the input graph. Here, we introduce some works on defense against adversarial attacks. The robust GCN (RGCN) [180] can fortify the GCNs against adversarial attacks. More specifically, it adopts Gaussian distributions as the hidden representations of nodes, i.e.

$$X^{(l+1)}[j, :] \sim N \left( \mu_j^{(l+1)}, \text{diag}(\sigma_j^{(l+1)}) \right),$$

in each graph convolution layer so that the effects of adversarial attacks can be absorbed into the variances of the Gaussian distributions. The Gaussian based graph convolution is defined as

$$\begin{aligned}
 \mu_j^{(l+1)} &= \rho \left( \sum_{v_k \in N_G(v_j)} \frac{\mu_k^{(l)} \otimes \alpha_k^{(l)}}{\sqrt{\tilde{D}_{j,j} \tilde{D}_{k,k}}} W_{\mu}^{(l)} \right) \\
 \sigma_j^{(l+1)} &= \rho \left( \sum_{v_k \in N_G(v_j)} \frac{\sigma_k^{(l)} \otimes \alpha_k^{(l)} \otimes \alpha_k^{(l)}}{\tilde{D}_{j,j} \tilde{D}_{k,k}} W_{\sigma}^{(l)} \right),
 \end{aligned}$$

where  $\alpha_j^{(k)}$  are attention weights. Finally, the overall loss function is defined as regularized cross-entropy. The first regularization term is

$$\mathcal{L}_{\text{reg1}} = \sum_{j=1}^N D_{KL} \left( N \left( \mu_j^1, \text{diag} \left( \sigma_j^{(1)} \right) \right) \| N(0, 1) \right),$$

which is used to constrain the latent representations at the first layer. The second regularization term is to impose  $L_2$  regularization on parameters  $W_{\mu}^{(0)}$  and  $W_{\sigma}^{(0)}$  of the first layer. The literature [176] presents a batch virtual adversarial training method which appends a novel regularization term to the conventional objective function of the GCNNs. Therefore, the resultant loss is

$$\mathcal{L} = \mathcal{L}_0 + \alpha \cdot \frac{1}{N} \sum_{u \in V} E(p(y|X_u, W)) + \beta \cdot \mathcal{R}_{\text{vadv}}(V, W),$$

where  $\mathcal{L}_0$  is an average cross-entropy loss of all labelled nodes,  $E(\cdot)$  is the conditional entropy of a distribution, and  $\mathcal{R}_{\text{vadv}}(V, W)$  is the average Local Distributional Smoothness (LDS) loss for all nodes. That is,

$$\mathcal{R}_{\text{vadv}}(V, W) = \frac{1}{N} \sum_{u \in V} \text{LDS}(X[u, :], W, r_{\text{vadv}, u}),$$

where  $\text{LDS}(x, w, r_{\text{vadv}}) = D_{KL} \left( p(y|x, \widehat{W}) \| p(y|x + r_{\text{vadv}}, W) \right)$  and  $r_{\text{vadv}}$  is the virtual adversarial perturbation. Additionally, there are some works aiming at verifying certifiable (non-)robustness to structure and feature perturbations for the GCNNs and developing robust training algorithm [177, 178]. That is, once a node is certified, it is guaranteed for the node prediction to be robust to any structure or feature perturbations given an attack model. The literature [180] proposes to improve GCN generalization by minimizing the expected loss under small perturbations of the input graph. Its basic assumption is that the adjacency matrix  $A_G$  is perturbed by some random noises. Under this assumption, the objective function is defined as

$$\min_W \int q(\epsilon|\alpha) \mathcal{L}(X, Y, A_G(\epsilon), W) d\epsilon,$$

where  $A_G(\epsilon)$  denotes the perturbed adjacency matrix of  $G$  and  $q(\epsilon|\alpha)$  is a zero-centered density of the noise  $\epsilon$  so that the learned GCN is robust to these noises and generalizes well.

#### 4.6 Graph Neural Architecture Search

Neural Architecture Search (NAS) [171] has achieved tremendous success in discovering the optimal neural network architecture for image and language learning tasks. However, existing NAS algorithms cannot be directly generalized to find the optimal GNN architecture. Fortunately, there have been some studies to bridge this gap. The graph neural architecture search [172, 173] aims to search for an optimal GNN architecture within a designed search space. It usually exploits a reinforcement learning based controller, which is a RNN, to greedily validate the generated architecture, and then the validation results are fed back to the controller. The literature [175] proposes a Graph HyperNetwork (GHN) to amortize the search cost of training thousands of different networks, which is trained to minimize the training loss of the sampled network with the weights generated by a GCNN.

#### 4.7 Graph Reinforcement Learning

The GNNs can also be combined with the reinforcement learning so as to solve sequential decision-making problems on graphs. The literature [188] learns to walk over a graph from a source node towards a target node for a given query via reinforcement learning. The proposed agent M-Walk is composed of a deep RNN and Monte Carlo Tree Search (MCTS). The former maps a hidden vector representation  $h_t$ , yielded by a special RNN encoding the state  $s_t$  at time  $t$ , to a policy and Q-values, and the latter is employed to generate trajectories yielding more positive rewards. The NerveNet [187] propagates information over the underlying graph of an agent via a GCNN, and then predicts actions for different parts of the agent. The literature [186] combines the GNNs and Deep Reinforcement Learning (DRL), named DRL+GNN, to learn, operate and generalize over arbitrary network topologies. The DRL+GNN agent employs a GCNN to model the Q-value function.

#### 4.8 Applications

In this subsection, we introduce the applications of the GNNs. Due to the space limitation, the applications are illustrated with Fig. 7. In this figure, we only list the application fields and the corresponding references.

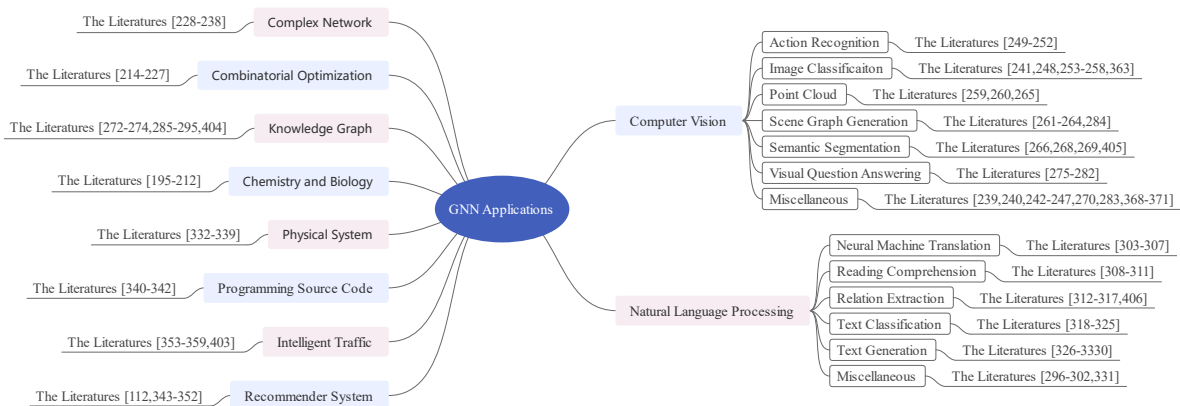


Figure 7 The GNN Application Summarizations

## 5 Benchmarks and Evaluation Pitfalls

In this section, we briefly introduce benchmarks, evaluation Pitfalls and applications of the GNNs. The benchmarks provide ground truth for various GNN architectures so that different GNNs can be compared fairly, the evaluation pitfalls empirically show that the existing evaluation criterion have potential pitfalls. Moreover, we summarize applications of the GNNs.

**Benchmarks.** Graph neural networks have become a powerful toolkit for mining complex graphs. It becomes more and more critical to evaluate the effectiveness of new GNN architectures and compare different GNN models under a standardized benchmark with consistent experimental settings and large datasets. A feasible benchmark for the GNNs should include appropriate graph datasets, robust coding interfaces and experimental settings so that different GNN architectures can be compared in the same settings. The literature [169] makes a pioneering effort to construct a reproducible GNN benchmarking framework in order to facilitate researchers to gauge the effectiveness of different GNN architectures. Specifically, it releases an open-source benchmark infrastructures for GNNs, hosted on GitHub based on PyTorch and DGL libraries [402], introduces medium-scale graph datasets with 12k-70k graphs of variable sizes 9-500 nodes, and identifies important building blocks of GNNs (graph convolutions, anisotropic diffusion, residual connections and normalization layers) with the proposed benchmark infrastructures. The literature [372] presents an Open Graph Benchmark (OGB) including challenging, real-world and large-scale benchmark graph datasets, encompassing multiple important graph deep learning tasks ranging from social and information networks to biological networks, molecular graphs and knowledge graphs, to facilitate scalable, robust and reproducible deep learning research on graphs. The OGB datasets, which provide a unified evaluation protocol using application-specific train/validation/test dataset splits and evaluation metrics, releases an end-to-end graph processing pipeline including graph data loading, experimental settings and model evaluations.

**Evaluation Pitfalls.** The literature [170] compares four typical GCNN architectures: GCN [6], MoNet [18], GAT [66] and GraphSAGE using three aggregation strategies [104] against 4 baseline models: logistic regression, multi-layer perceptron, label propagation and normalized laplacian label propagation, and uses a standardized training and hyper-parameter tuning procedure for all these models so as to perform a more fair comparison. The experimental results show that different train/validation/test splits of datasets lead to dramatically different rankings of models. In addition, their findings also demonstrate that simpler GCNN architectures can outperform more sophisticated ones only if the hyper-parameters and training procedures are tuned fairly for all models.

## 6 Future Research Directions

Although the GNNs have achieved tremendous success in many fields, there still exists some open problems. This section summarizes the future research directions of the GNNs.

**Highly Scalable GNNs.** The real-world graphs usually contain hundreds of millions of nodes and edges, and have dynamically evolving characteristics. It turns out that it is difficult for the existing GNN architectures to scale up to the huge real-world graphs. This motivates us to design highly scalable GNN architectures which can efficiently and effectively learn node/edge/graph representations for the huge dynamically-evolving graphs.

**Robust GNNs.** The existing GNN architectures are vulnerable to adversarial attacks. That is, the performance of the GNN models will sharply drop once the structure and/or initial features of the input graph are attacked by adversaries. Therefore, we should incorporate the attack-and-defense mechanism into the GNN architectures, i.e. constructing robust GNN architecture, so as to reinforce them against adversarial attacks.

**GNNs Going Beyond WL Test.** The capabilities of the spatial GCNNs are limited by the 1-WL test, and the higher-order WL test is computationally expensive. Consequently, two non-isomorphic graphs will produce the same node/edge/graph representations under appropriate conditions. This mo-

tivates us to develop a novel GNN framework going beyond WL test, or design an elegant higher-order GNN architectures corresponding to the higher-order WL test.

**Interpretable GNNs.** The existing GNNs work in a black box. We do not understand why they achieve state-of-the-art performance in terms of the node classification task, graph classification task and graph embedding task etc. Interpretability has become a major obstacle to apply the GNNs to real-world issues. Although there have been some studies to interpret some specific GNN models, they cannot interpret general GNN models. This motivates us to construct a unified interpretable framework for the GNNs.

## 7 Conclusions

This paper aims to provide a taxonomy, advances and trends of the GNNs. The GNN architectures are classified into 9 categories (graph convolutional neural networks, capability and interpretability, graph recurrent neural networks, deep graph representation learning, deep graph generative models, combinations of the PI and the GNNs, adversarial attacks for the GNNs, graph neural architecture search and graph reinforcement). Then, we review the benchmarks and evaluation pitfalls, applications and future research directions of the GNNs. It is expected that the relevant scholars can understand the computational principles of the GNNs, consolidate the foundations of the GNNs and apply them to more and more real-world issues, through reading this review.

**Acknowledgements** This work was supported by National Natural Science Foundation of China (Grant No. 62002255).

## References

- 1 Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, Pierre Vandergheynst. Geometric deep learning: going beyond Euclidean data, *IEEE Signal Processing Magazine*, 2017, 34(4): 18-42.
- 2 Ziwei Zhang, Peng Cui, Wenwu Zhu. Deep learning on graphs: a survey, 2018, <https://arxiv.org/abs/1812.04202>.
- 3 Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Maosong Sun. Graph neural networks: a review of methods and applications, 2018, <https://arxiv.org/abs/1812.08434>.
- 4 Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, Philip S. Yu. A comprehensive survey on graph neural networks, 2019, <https://arxiv.org/abs/1901.00596>.
- 5 David I. Shuman, Sunil K. Narang, Pascal Frossard, Antonio Ortega, Pierre Vandergheynst. The emerging field of signal processing on graphs: extending high-dimensional data analysis to networks and other irregular domains, *IEEE Signal Processing Magazine*, 2013, 30(3): 83-98.
- 6 Thomas N. Kipf, Max Welling. Semi-supervised classification with graph convolutional networks, *The International Conference on Learning Representations*, 2017, <https://openreview.net/forum?id=SJU4ayYgl>.
- 7 Joan Bruna, Wojciech Zaremba, Arthur Szlam, Yann LeCun. Spectral networks and locally connected networks on graphs, *The International Conference on Learning Representations*, 2014, <https://openreview.net/forum?id=DQNsQf-Us0DBa>.
- 8 Ruoyu Li, Sheng Wang, Feiyun Zhu, Junzhou Huang. Adaptive Graph Convolutional Neural Networks, *The AAAI Conference on Artificial Intelligence*, 2018, 3546-3553.
- 9 Wenbing Huang, Tong Zhang, Yu Rong, Junzhou Huang. Adaptive sampling towards fast graph representation learning, *The International Conference on Neural Information Processing Systems*, 2018, 4563-4572.
- 10 Jie Chen, Tengfei Ma, Cao Xiao. FastGCN: fast learning with graph convolutional networks via importance sampling, *The International Conference on Learning Representations*, 2018, <https://openreview.net/forum?id=rytstxWAW>.
- 11 Jianfei Chen, Jun Zhu, Le Song. Stochastic training of graph convolutional networks with variance reduction, *The International Conference on Machine Learning*, 2018, 942-950.
- 12 Ron Levie, Federico Monti, Xavier Bresson, Michael M. Bronstein. CayleyNets: Graph Convolutional Neural Networks with Complex Relational Spectral Filters, *IEEE Transactions on Signal Processing*, 2019, 67(1): 97-109.
- 13 Guohao Li, Matthias Müller, Ali Thabet, Bernard Ghanem. DeepGCNs: can GCNs go as deep as CNNs, *The IEEE International Conference on Computer Vision (ICCV)*, 2019, 9267-9276.
- 14 Renjie Liao, Zhizhen Zhao, Raquel Urtasun, Richard S. Zemel. LanczosNet: Multi-Scale Deep Graph Convolutional Networks, *The International Conference on Learning Representations*, 2019, <https://openreview.net/forum?id=BkedznAqKQ>.
- 15 Ines Chami, Rex Ying, Christopher Re, Jure Leskovec. Hyperbolic Graph Convolutional Neural Networks, *The International Conference on Neural Information Processing Systems*, 2019, 4868-4879.
- 16 Qimai Li, Zhichao Han, Xiaoming Wu. Deeper insights into graph convolutional networks for semi-supervised learning, *The AAAI Conference on Artificial Intelligence (AAAI)*, 2018, 3538-3545.

- 17 Hongbin Pei, Bingzhe Wei, Kevin C.C. Chang, Yu Lei, Bo Yang. Geom-GCN: geometric graph convolutional networks, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=S1e2agrFvS>.
- 18 Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodolà, Jan Svoboda, Michael M. Bronstein. Geometric deep learning on graphs and manifolds using mixture model CNNs, The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017, 5425-5434.
- 19 Chenyi Zhuang, Qiang Ma. Dual graph convolutional networks for graph-based semi-supervised classification, The World Wide Web Conference (WWW), 2018, 499-508.
- 20 James Atwood, Don Towsley. Diffusion-Convolutional Neural Networks. The International Conference on Neural Information Processing Systems (NIPS), 2016, 1993-2001.
- 21 Jun Wu, Jingrui He, Jiejun Xu. DEMO-Net: degree-specific graph neural networks for node and graph classification, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 406-415.
- 22 Felix Wu, Tianyi Zhang, Amauri Holanda de Souza Jr., Christopher Fifty, Tao Yu, Kilian Q. Weinberger. Simplifying Graph Convolutional Networks, The International Conference on Machine Learning (ICML), 2019, 6861-6871.
- 23 Weilin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, Cho-Jui Hsieh. Cluster-GCN: an efficient algorithm for training deep and large graph convolutional networks, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 257-266.
- 24 Xin Jiang, Kewei Cheng, Song Jiang, Yizhou Sun. Chordal-GCN: exploiting sparsity in training large-scale graph convolutional networks, 2020, <https://openreview.net/forum?id=rJI05AVtwB>.
- 25 Hongyang Gao, Zhengyang Wang, Shuiwang Ji. Large-scale learnable graph convolutional networks, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2018, 1415-1424.
- 26 Anh Viet Phan, Minh Le Nguyen, Yen Lam Hoang Nguyen, Lam Thu Bui. DGCNN: a convolutional neural network over large-scale labeled graphs, Neural Networks, 2018, 108(2018): 533-543.
- 27 Mikael Henaff, Joan Bruna, Yann LeCun. Deep convolutional networks on graph-structured data, 2018, <https://arxiv.org/abs/1506.05163>.
- 28 Tyler Derr, Yao Ma, Jiliang Tang. Signed Graph Convolutional Network, The IEEE International Conference on Data Mining, 2018, 929-934.
- 29 Trang Pham, Truyen Tran, Dinh Phung, Svetha Venkatesh. Column Networks for Collective Classification, The AAAI Conference on Artificial Intelligence, 2017, 2485-2491.
- 30 Ana Šušnjara, Nathanaël Perraudin, Daniel Kressner, Pierre Vandergheynst. Accelerated filtering on graphs using Lanczos method, 2015, <https://arxiv.org/abs/1509.04537>.
- 31 Yotam Hechtlinger, Purvasha Chakravarti, Jining Qin. A generalization of convolutional neural networks to graph-structured data, 2017, <https://arxiv.org/abs/1704.08165>.
- 32 Felipe Petroski Such, Shagan Sah, Miguel Dominguez, Suhas Pillai, Chao Zhang, Andrew Michael Nathan D. Cahill, Raymond Ptucha. Robust spatial filtering with graph convolutional neural networks, IEEE Journal of Selected Topics in Signal Processing, 2017, 11(6): 884-896.
- 33 Yawei Luo, Tao Guan, Junqing Yu, Ping Liu, Yi Yang. Every node counts: self-ensembling graph convolutional networks for semi-supervised learning, 2018, <https://arxiv.org/abs/1809.09925v1>.
- 34 Ziqi Liu, Chaochao Chen, Longfei Li, Jun Zhou, Xiaolong Li, Le Song, Yuan Qi. GeniePath: graph neural networks with adaptive receptive paths, The AAAI Conference on Artificial Intelligence, 2019, 4424-4431.
- 35 Renjie Liao, Marc Brockschmidt, Daniel Tarlow, Alexander L. Gaunt, Raquel Urtasun, Richard Zemel. Graph partition neural networks for semi-supervised classification, The International Conference on Learning Representations (ICLR Workshop), 2018, <https://openreview.net/forum?id=ByM6TrTUz>.
- 36 Yao Ma, Suhang Wang, Charu C. Aggarwal, Dawei Yin, Jiliang Tang. Multi-dimensional graph convolutional networks, 2018, <https://arxiv.org/abs/1808.06099>.
- 37 Tong Zhang, Wenming Zheng, Zhen Cui, Yang Li. Tensor graph convolutional neural network, 2018, <https://arxiv.org/abs/1803.10071v1>.
- 38 Jian Du, Shanghang Zhang, Guanhang Wu, José M.F. Moura, Soumya Kar. Topology adaptive graph convolutional networks, 2018, <https://arxiv.org/abs/1710.10370>.
- 39 Matthew Baron. Topology and prediction focused research on graph convolutional neural networks, 2018, <https://arxiv.org/abs/1808.07769v1>.
- 40 Lingxiao Ma, Zhi Yang, Youshan Miao, Jilong Xue, Ming Wu, Lidong Zhou, Yafei Dai. Towards efficient large-scale graph neural network computing, 2018, <https://arxiv.org/abs/1810.08403>.
- 41 Peter Meltzer, Marcelo Daniel Gutierrez Mallea, Peter J. Bentley. PiNet: a permutation invariant graph neural network for graph classification, 2019, <https://arxiv.org/abs/1905.03046>.
- 42 Fenyu Hu, Yanqiao Zhu, Shu Wu, Liang Wang, Tieniu Tan. Semi-supervised node classification via hierarchical graph convolutional networks, 2019, <https://arxiv.org/abs/1902.06667v2>.
- 43 Martin Simonovsky, Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs, The IEEE Conference on Computer Vision and Pattern Recognition, 2017, 29-38.
- 44 Yuzhou Chen, Yulia R. Gel, Konstantin Avrachenkov. Fractional graph convolutional networks (FGCN) for semi-supervised learning, 2020, <https://openreview.net/forum?id=BygacxrFwS>.
- 45 Shengding Hu, Meng Qu, Zhiyuan Liu, Jian Tang. Transfer active learning for graph neural networks, 2020, <https://openreview.net/forum?id=BklOXeBFDS>.
- 46 Jiaxue You, Rex Ying, Jure Leskovec. Position-aware graph neural networks, 2019, The International Conference on

- Machine Learning, 7134-7143.
- 47 Vincenzo Di. Massa, Gabriele Monfardini, Lorenzo Sarti, Franco Sarselli, Marco Maggini, Marco Gori. A comparison between recursive neural networks and graph neural networks, *The IEEE Joint Conference on Neural Networks*, 2006, 778-785.
  - 48 Mathias Niepert, Mohamed Ahmed, Konstantin Kutzkov. Learning convolutional neural networks for graphs, *The International Conference on Machine Learning*, 2016, 2014-2023.
  - 49 Simon S. Du, Kangcheng Hou, Barnabás Póczos, Ruslan Salakhutdinov. Graph neural tangent kernel: fusing graph neural networks with graph kernels, *The International Conference on Neural Information Processing Systems*, 2019, 5273-5733.
  - 50 Nicolas Keriven, Gabriel Peyré. Universal invariant and equivariant graph neural networks, *The International Conference on Neural Information Processing Systems*, 2019, 7092-7101.
  - 51 Zhihong Zhang, Dongdong Chen, Jianjia Wang, Lu Bai, Edwin R. Hancock. Quantum-based subgraph convolutional neural networks, *Pattern Recognition*, 2019, 88(2019): 38-49.
  - 52 Xinyi Zhang, Lihui Chen. Capsule Graph Neural Network, *The International Conference on Learning Representations*, 2019, <https://openreview.net/forum?id=Byl8BnRcYm>.
  - 53 Saurabh Verma, Zhili Zhang. Graph Capsule Convolutional Neural Networks, 2018, <https://arxiv.org/abs/1805.08090>.
  - 54 Marcelo Daniel Gutierrez Mallea, Peter Meltzer, Peter J. Bentley. Capsule neural networks for graph classification using explicit tensorial graph representations, 2019, <https://arxiv.org/abs/1902.08399>.
  - 55 Sijie Yan, Yuanjun Xiong, Dahua Lin. Spatial temporal graph convolutional networks for skeleton-based action recognition, *The AAAI Conference on Artificial Intelligence*, 2018, 7444-7452.
  - 56 Ashesh Jain, Amir R. Zamir, Silvio Savarese, Ashutosh Saxena. Structural-RNN: deep learning on spatio-temporal graphs, *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, 5308-5317.
  - 57 Yujun Cai, Lihao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, Nadia Magnenat Thalmann.
  - 58 Runhao Zeng, Wenbing Huang, Mingkui Tan, Yu Rong, Peilin Zhao, Junzhou Huang, Chuang Gan. Graph convolutional networks for temporal action localization, *The IEEE International Conference on Computer Vision (ICCV)*, 2019, 7094-7103.
  - 59 Bing Yu, Haoteng Yin, Zhanxing Zhu. Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting, *The International Joint Conference on Artificial Intelligence*, 2018, 3634-3640.
  - 60 Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, Hyunwoo J. Kim. Graph Transformer Networks, 2019, <https://arxiv.org/abs/1911.06455>.
  - 61 Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, Nitesh V. Chawla. Heterogeneous Graph Neural Network, *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2019, 793-803.
  - 62 Yizhou Zhang, Yun Xiong, Xiangnan Kong, Shanshan Li, Jinhong Mi, Yangyong Zhu. Deep Collective Classification in Heterogeneous Information Networks, *The World Wide Web Conference*, 2018, 399-408.
  - 63 Naganand Yadati, Tingran Gao, Shahab Asoodeh, Partha Talukdar, Anand Louis. Graph neural networks for soft semi-supervised learning on hypergraphs, 2020, <https://openreview.net/forum?id=rystJBKPB>.
  - 64 Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, Partha Talukdar. HyperGCN: a new method for training graph convolutional networks on hypergraphs, *The International Conference on Neural Information Processing Systems*, 2019, 1511-1522.
  - 65 Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, Yue Gao. Hypergraph Neural Networks, *The AAAI Conference on Artificial Intelligence (AAAI)*, 2019, 3558-3565.
  - 66 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, Yoshua Bengio. Graph Attention Networks, *The International Conference on Learning Representations (ICLR)*, 2018, <https://openreview.net/forum?id=rJXMpikCZ>.
  - 67 Kiran K. Thekumparampil, Chong Wang, Sweeong Oh, Lijia Li. Attention-based graph neural network for semi-supervised learning, 2018, <https://arxiv.org/abs/1803.03735>.
  - 68 Chao Shang, Qinqing Liu, Ko-Shin Chen, Jiangwen Sun, Jin Lu, Jinfeng Yi, Jinbo Bi. Edge attention-based multi-relational graph convolutional networks, 2018, <https://arxiv.org/abs/1802.04944>.
  - 69 Caglar Gulcehre, Misha Denil, Mateusz Malinowski, Ali Razavi, Razvan Pascanu, Karl Moritz Hermann, Peter Battaglia, Victor Bapst, David Raposo, Adam Santoro, Nando de Freitas. Hyperbolic Attention Networks, *The International Conference on Learning Representations*, 2019, <https://openreview.net/forum?id=rJxHsjRqFQ>.
  - 70 Xiaoran Xu, Songpeng Zu, Chengliang Gao, Yuan Zhang, Wei Feng. Modeling attention flow on graphs, 2018, <https://arxiv.org/abs/1811.00497>.
  - 71 Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, Huaiyu Wan. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting, *The AAAI Conference on Artificial Intelligence (AAAI)*, 2019, 922-929.
  - 72 Olaf Ronneberger, Philipp Fischer, Thomas Brox. U-Net: convolutional networks for biomedical image segmentation, *The International Conference on Medical Image Computing and Computer-assisted Intervention*, 2015, 234-241.
  - 73 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Lion Jones Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin. Attention is all you need, *The International Conference on Neural Information Processing (NIPS)*, 2017, 5998-6008.
  - 74 Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, Dit-Yan Yeung. GaAN: gated attention networks for learning on large and spatiotemporal graphs, *The International Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018, Article No. 139.
  - 75 Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Yanfang Ye, Peng Cui, Philip S. Yu. Heterogeneous Graph Attention

- Network, The World Wide Web Conference, 2018, 2022-2032.
- 76 Marco Gori, Gabriele Monfardini, Franco Scarselli. A new model for learning in graph domains, The International Joint Conference on Neural Networks (IJCNN), 2005, 729-734.
  - 77 Franco Scarselli, Marco Gori, Ah C. Tsoi, Markus Hagenbuchner, Gabriele Monfardini. Computational Capabilities of Graph Neural Networks, IEEE Transactions on Neural Networks, 2009, 20(1): 81-102.
  - 78 Michaël Defferrard, Xavier Bresson, Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering, Advances in Neural Information Processing Systems (NIPS), 2016, 3844-3852.
  - 79 Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, Yixin Chen. D-VAE: A Variational Autoencoder for Directed Acyclic Graphs, 2019, <https://arxiv.org/abs/1904.11088>.
  - 80 Meng Qu, Yoshua Bengio, Jian Tang. GMNN: Graph Markov Neural Networks, The International Conference on Machine Learning (ICML), 2019, 5241-5250.
  - 81 John Boaz Lee, Ryan Rossi, Xiangnan Kong. Graph classification using structural attention, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2018, 1666-1674.
  - 82 Yao Ma, Suhang Wang, Charu C. Aggarwal, Jiliang Tang. Graph convolutional networks with EigenPooling, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 723-731.
  - 83 Keyulu Xu, Weihua Hu, Jure Leskovec, Stefanie Jegelka. How powerful are graph neural networks, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=ryGs6iA5Km>.
  - 84 Amir H. Khasahmadi, Kaveh Hassani, Parsa Moradi, Leo Lee, Quaid Morris. Memory-based graph networks, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=r1laNeBYPB>.
  - 85 Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, et al. Relational inductive biases, deep learning, and graph networks, 2018, <https://arxiv.org/abs/1806.01261>.
  - 86 Junhyun Lee, Inyeop Lee, Jaewoo Kang. Self-Attention Graph Pooling, The International Conference on Machine Learning (ICML), 2019, 3734-3743.
  - 87 Thomas N. Kipf, Max Welling. Variational Graph Auto-Encoders, 2016, <https://arxiv.org/abs/1611.07308>.
  - 88 Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, Martin Grohe. Weisfeiler and Leman Go neural: higher-order graph neural networks, The AAAI Conference on Artificial Intelligence (AAAI), 2019, 4602-4609.
  - 89 Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, Stefanie Jegelka. Representation Learning on Graphs with Jumping Knowledge Networks, The International Conference on Machine Learning (ICML), 2018, 5453-5462.
  - 90 Sami Abu-EL-Haija, Bryan Perozzi, Amol Kapoor, Hrayr Harutyunyan, Nazanin Alipourfard, Kristina Lerman, Greg Ver Steeg, Aram Galstyan. MixHop: higher-order graph convolutional architectures via sparsified neighborhood mixing, The International Conference on Machine Learning (ICML), 2019, 21-29.
  - 91 Sungmin Rhee, Seokjun Seo, Sun Kim. Hybrid approach of relation network and localized graph convolutional filtering for breast cancer subtype classification, The International Joint Conference on Artificial Intelligence (IJCAI), 2018, 3527-3534.
  - 92 Muhan Zhang, Zhicheng Cui, Marion Neumann, Yixin Chen. An end-to-end deep learning architecture for graph classification, The AAAI Conference on Artificial Intelligence (AAAI), 2018, 4438-4445.
  - 93 Justin Gilmer Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, George E. Dahl. Neural Message Passing for Quantum Chemistry, The International Conference on Machine Learning (ICML), 2017, 1263-1272.
  - 94 Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, Jure Leskovec. Hierarchical graph representation learning with differentiable pooling, The International Conference on Neural Information Processing Systems (NeurIPS), 2018, 4805C4815.
  - 95 Hongyang Gao, Shuiwang Ji. Graph U-Nets, 2019, <https://arxiv.org/abs/1905.05178>.
  - 96 Čačalina Cangea, Petar Veličković, Nikola Jovanović, Thomas Kipf, Pietro Liò. Towards sparse hierarchical graph classifiers, 2018, <https://arxiv.org/abs/1811.01287>.
  - 97 Padraig Corcoran. Function space pooling for graph convolutional networks, 2019, <https://arxiv.org/abs/1905.06259>.
  - 98 Ryan Murphy, Balasubramaniam Srinivasan, Vinayak Rao, Bruno Ribeiro. Relational pooling for graph representations, The International Conference on Machine Learning (ICML), 2019, 4663-4673.
  - 99 Inderjit S. Dhillon, Yuqiang Guan, Brian Kulis. Weighted graph cuts without eigenvectors: a multilevel approach, IEEE Transactions on Pattern Analysis and Machine Intelligence, 2007, 29(11): 1944-1957.
  - 100 Jessica B. Hamrick, Kelsey R. Allen, Victor Bapst, Tina Zhu, Kevin R. McKee, Joshua B. Tenenbaum, Peter W. Battaglia. Relational inductive bias for physical construction in humans and machines, 2018, <https://arxiv.org/abs/1806.01203>.
  - 101 Xiaolong Wang, Ross Girshick, Abhinav Gupta, Kaiming He. Non-local neural networks, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 7794-7803.
  - 102 Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, Peter Battaglia. Graph networks as learnable physics engines for inference and control, The International Conference on Machine Learning (ICML), 2018, 4470-4479.
  - 103 David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gomez-Bombarelli, Timothy Hirzel, Alan Aspuru-Guzik, Ryan P. Adams. Convolutional networks on graphs for learning molecular fingerprints, The International Conference on Neural Information Processing Systems, 2015, 2224-2232.
  - 104 William L. Hamilton, Rex Ying, Jure Leskovec. Inductive representation learning on large graphs, The International Conference on Neural Information Processing Systems (NIPS), 2017, 1024-1034.

- 105 Franco Manessi, Alessandro Rozza, Mario Manzo. Dynamic graph convolutional networks, *Pattern Recognition*, 2020, 97(2020): Article No. 107000.
- 106 Kun Xu, Lingfei Wu, Zhiguo Wang, Yansong Feng, Michael Witbrock, Vadim Sheinin. Graph2Seq: graph to sequence learning with attention-based neural networks, 2018, <https://arxiv.org/abs/1804.00823>.
- 107 Yao Mao, Ziyi Guo, Zhaochun Ren, Eric Zhao, Jiliang Tang, Dawei Yin. Streaming graph neural networks, 2018, <https://arxiv.org/abs/1810.10627>.
- 108 Luana Ruiz, Fernando Gama, Alejandro Ribeiro. Gated Graph Convolutional Recurrent Neural Networks, 2019, <https://arxiv.org/abs/1903.01888>.
- 109 Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, Xavier Bresson. Structured sequence modeling with graph convolutional recurrent networks, *The International Conference on Neural Information Processing (ICONIP)*, 2018, 362-373.
- 110 Yuan Li, Xiaodan Liang, Zhiting Hu, Yinbo Chen, Eric P. Xing. Graph Transformer, 2019, <https://openreview.net/forum?id=HJei-2RcK7>.
- 111 Edouard Pineau, Nathan de Lara. Variational recurrent neural networks for graph classification, 2019, <https://arxiv.org/abs/1902.02721>.
- 112 Federico Monti, Michael M. Bronstein, Xavier Bresson. Geometric matrix completion with recurrent multi-graph neural network, *The International Conference on Neural Information Processing Systems*, 2017, 3697-3707.
- 113 Ehsan Hajiramezani, Arman Hasanzadeh, Krishna Narayanan, Nick Duffield, Mingyuan Zhou, Xiaoning Qian. Variational Graph Recurrent Neural Networks, *The International Conference on Neural Information Processing Systems*, 2019, 10701-10711.
- 114 Ziniu Hu, Yuxiao Dong, Kuansan Wang, Yizhou Sun. Heterogeneous Graph Transformer, 2020, <https://arxiv.org/abs/2003.01332>.
- 115 Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, Youshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation, *The Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, 1724-1734.
- 116 Xiao Huang, Qingquan Song, Yuening Li, Xia Hu. Graph recurrent networks with attributed random walks, *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2019, 732-740.
- 117 Victoria Zayats, Mari Ostendorf. Conversation modeling on Reddit using a graph-structured LSTM, *Transactions of the Association for Computational Linguistics*, 2018, 6(2018): 121-132.
- 118 Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, Wen-tau Yih. Cross-sentence N-ary relation extraction with graph LSTMs. *Transactions of the Association for Computational Linguistics*, 2017, 5(2017): 101-115.
- 119 Xavier Bresson, Thomas Laurent. Residual gated graph ConvNets, 2017, <https://arxiv.org/abs/1711.07553>.
- 120 Yue Zhang, Qi Liu, Linfeng Song. Sentence-State LSTM for text representation, *The Annual Meeting of the Association for Computational Linguistics*, 2018, 317-327.
- 121 Xiaodan Liang, Xiaohui Shen, Jiashi Feng, Liang Lin, Shuicheng Yan. Semantic Object Parsing with Graph LSTM, *The European Conference on Computer Vision (ECCV)*, 2016, 125-143.
- 122 Yujia Li, Daniel Tarlow, Marc Brockschmidt, Richard Zemel. Gated graph sequence neural networks, *The International Conference on Learning Representations*, 2016, <https://arxiv.org/abs/1511.05493>.
- 123 Kai Sheng Tai, Richard Socher, Christopher D. Manning. Improved semantic representations from tree-structured long short-term memory networks, *The Annual Meeting of the Association for Computational Linguistics*, 2015, 1556-1566.
- 124 Shaosheng Cao, Wei Lu, Qiongkai Xu. Deep neural networks for learning graph representations, *The AAAI Conference on Artificial Intelligence (AAAI)*, 2016, 1145C1152.
- 125 Thomas N. Kipf, Max Welling. Variational graph auto-encoders, 2016, <https://arxiv.org/abs/1611.07308>.
- 126 Wenting Zhao, Chunyan Xu, Zhen Cui, Tong Zhang, Jiatao Jiang, Zhenyu Zhang, Jian Yang. When work matters: transforming classical network structures to graph CNN, 2018, <https://arxiv.org/abs/1807.02653v1>.
- 127 Sébastien Lerique, Jacob Levy Abitol, Márton Karsai. Joint embedding of structure and features via graph convolutional networks, *Applied Network Science*, 2020, 5(2020): Article No. 5.
- 128 Jiwoong Park, Minsik Lee, Hyung Jin Chang, Kyuewang Lee, Jin Young Choi. Symmetric graph convolutional autoencoder for unsupervised graph representation learning, *The International Conference on Computer Vision*, 2019, 6519-6528.
- 129 Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, Chengqi Zhang. Adversarially regularized graph autoencoder for graph embedding, *The International Joint Conference on Artificial Intelligence (IJCAI)*, 2018, 2609-2615.
- 130 Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay S. Pande, Jure Leskovec. Strategies for pre-training graph neural networks, *The International Conference on Learning Representations (ICLR)*, 2020, <https://openreview.net/forum?id=HJIWWJSFDH>.
- 131 Zhihong Zhang, Dongdong Chen, Zeli Wang, Heng Li, Lu Bai, Edwin R. Hancock. Depth-based subgraph convolutional auto-encoder for network representation learning, *Pattern Recognition*, 2019, 90(2019): 363-376.
- 132 Daixin Wang, Peng Cui, Wenwu Zhu. Structural Deep Network Embedding, *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2016, 1225-1234.
- 133 Dingyuan Zhu, Peng Cui, Daixin Wang, Wenwu Zhu. Deep Variational Network Embedding in Wasserstein Space, *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2018, 2827-2836.
- 134 Ke Tu, Peng Cui, Xiao Wang, Philip S. Yu. Deep Recursive Network Embedding with Regular Equivalence, *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2018, 2357-2366.
- 135 Wenchao Yu, Cheng Zheng, Wei Cheng. Learning deep network representations with adversarially regularized autoen-



- coders, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2018, 2663-2671.
- 136 Ke Tu, Peng Cui, Xiao Wang, Fei Wang, Wenwu Zhu. Structural Deep Embedding for Hyper-Networks, The AAAI Conference on Artificial Intelligence (AAAI), 2018, 426-433.
- 137 Aditya Grover, Aaron Zweig, Stefano Ermon. Graphite: iterative generative modeling of graphs, The International Conference on Machine Learning (ICML), 2019, 2434-2444.
- 138 Zhen Zhang, Hongxia Yang, Jiajun Bu, Sheng Zhou, Pinggang Yu, Jianwei Zhang, Martin Ester, Can Wang. ANRL: attributed network representation learning via deep neural networks, The International Joint Conference on Artificial Intelligence (IJCAI), 2018, 3155-3161.
- 139 Wenwu Zhu, Xin Wang, Peng Cui. Deep learning for learning graph representations, 2020, <https://arxiv.org/abs/2001.00293v1>.
- 140 Petar Veličković, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, R Devon Hjelm. DEEP GRAPH INFOMAX, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=rklz9iAcKQ>.
- 141 Xiao Shen, Fulai Chung. Deep network embedding for graph representation learning in signed networks, IEEE Transactions on Cybernetics, 2020, 50(4): 1556-1568.
- 142 Palash Goyal, Nitin Kamra, Xinran He, Yan Liu. DynGEM: deep embedding method for dynamic graphs, 2018, <https://arxiv.org/abs/1805.11273v1>.
- 143 Yu Jin, Joseph F. JaJa. Learning graph-level representations with recurrent neural networks, 2018, <https://arxiv.org/abs/1805.07683>.
- 144 Wenchao Yu, Cheng Zheng, Wei Cheng, Charu C. Aggarwal, Dongjin Song, Bo Zong, Haifeng Chen, Wei Wang. Learning deep network representations with adversarially regularized autoencoders, The ACM International Conference on Knowledge Discovery and Data Mining, 2018, 2663-2671.
- 145 Xiaojie Guo, Lingfei Wu, Liang Zhao. Deep Graph Translation, 2018, <https://arxiv.org/abs/1805.09980>.
- 146 Martin Simonovsky, Nikos Komodakis. GraphVAE: towards generation of small graphs using variational autoencoders, The International Conference on Artificial Neural Networks (ICANN), 2018, 412-422.
- 147 Anuththari Gamage, Eli Chien, Jianhao Peng, Olgica Milenkovic. Multi-MotifGAN (MMGAN): motif-targeted graph generation and prediction, 2019, <https://arxiv.org/abs/1911.05469v1>.
- 148 Chun Wang, Shrui Pan, Guodong Long, Xingquan Zhu, Jing Jiang. MGAE: marginalized graph autoencoder for graph clustering. The International Conference on Information and Knowledge Management (CIKM), 2017, 889-898.
- 149 Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio. Generative Adversarial Nets, The International Conference on Neural Information Processing Systems (NIPS), 2014, 2672-2680.
- 150 Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, Peter Battaglia. Learning Deep Generative Models of Graphs, The International Conference on Learning Representations (ICLR Workshop), 2018, <https://openreview.net/forum?id=Hy1d-ebAb>.
- 151 David M. Blei, Alp Kucukelbir, Jon D. McAuliffe. Variational inference: a review for statisticians, Journal of the American Statistical Association, 2017, 112(518): 859-877.
- 152 Jiaxuan You, Rex Ying, Xiang Ren, William L. Hamilton, Jure Leskovec. GraphRNN: generating realistic graphs with deep auto-regressive models, The International Conference on Machine Learning (ICML), 2018, 5708-5717.
- 153 Aleksandar Bojchevski, Olexandr Shchur, Daniel Zügner, Stephan Günnemann. NetGAN: generating graphs via random walks, The International Conference on Machine Learning (ICML), 2018, 610-619.
- 154 Bryan Perozzi, Rami Al-Rfou, Steven Skiena. DeepWalk: online learning of social representations, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2014, 701-710.
- 155 Aditya Grover, Jure Leskovec. node2vec: scalable feature learning for networks, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2016, 855-864.
- 156 Linfeng Liu, Liping Liu. Amortized Variational Inference with Graph Convolutional Networks for Gaussian Processes, The International Conference on Artificial Intelligence and Statistics (AISTATS), 2019, 2291-2300.
- 157 Louis Tiao, Pantelis Elinas, Harrison Nguyen, Edwin V. Bonilla. Variational spectral graph convolutional networks, 2019, <https://arxiv.org/abs/1906.01852v1>.
- 158 KiJung Yoon, Renjie Liao, Yuwen Xiong, Lisa Zhang, Ethan Fetaya, Raquel Urtasun, Richard Zemel, Xaq Pitkow. Inference in probabilistic graphical models by graph neural networks, The International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=rkN1pF1vz>.
- 159 Tengfei Ma, Junyuan Shang, Jimeng Sun. CGNF: conditional graph neural fields, 2019, <https://openreview.net/forum?id=ryxMX2R9YQ>.
- 160 Sepp Hochreiter, Jürgen Schmidhuber. Long short-term memory, Neural Computation, 1997, 9(8): 1735-1780.
- 161 Radford M. Neal, Geoffrey E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. Learning in Graphical Models, 1998, 89: 355-368.
- 162 Hongchang Gao, Jian Pei, Heng Huang. Conditional Random Field Enhanced Graph Convolutional Neural Networks, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 276-284.
- 163 Andreas Loukas. What graph neural networks cannot learn-depth vs width, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=B1l2bp4YwS>.
- 164 Pablo Barceló, Egor V. Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, Juan Pablo Silva. The logical expressiveness of graph neural networks. The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=r1IZ7AEKvB>.

- 165 Federico Baldassarre, Hossein Azizpour. Explainability techniques for graph convolutional networks, 2019, <https://arxiv.org/abs/1905.13686>.
- 166 Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, Jure Leskovec. GNNExplainer: generating explanations for graph neural networks, The International Conference on Neural Information Processing Systems (NeurIPS), 2019, 9244-9255.
- 167 Nima Dehmamy, Albert-Laszlo Barabasi, Rose Yu. Understanding the Representation Power of Graph Neural Networks in Learning Graph Topology, The International Conference on Neural Information Processing Systems (NeurIPS), 2019, 15413-15423.
- 168 Yuyu Zhang, Xinshi Chen, Yuan Yang, Arun Ramamurthy, Bo Li, Yuan Qi, Le Song. Efficient probabilistic logic reasoning with graph neural networks, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=rJg76kStwH>.
- 169 Vijay Prakash Dwivedi, Chaitanya K. Joshi, Thomas Laurent, Yoshua Bengio, Xavier Bresson. Benchmarking graph neural networks, 2020, <https://arxiv.org/abs/2003.00982>.
- 170 Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, Stephan Günnemann. Pitfalls of graph neural network evaluation, 2019, <https://arxiv.org/abs/1811.05868>.
- 171 Wei Li, Shaogang Gong, Shaogang Gong. Neural graph embedding for neural architecture search, The AAAI Conference on Artificial Intelligence (AAAI), 2020.
- 172 Kaixiong Zhou, Qingquan Song, Xiao Huang, Xia Hu. Auto-GNN: neural architecture search of graph neural networks, 2019, <https://arxiv.org/abs/1909.03184>.
- 173 Yang Gao, Hong Yang, Peng Zhang, Chuan Zhou, Yue Hu. GraphNAS: graph neural architecture search with reinforcement learning, 2019, <https://arxiv.org/abs/1904.09981v1>.
- 174 Yihe Dong, Will Sawin, Yoshua Bengio. HNNH: hypergraph networks with hyperedge neurons, 2020, <https://arxiv.org/abs/2006.12278>.
- 175 Chris Zhang, Mengye Ren, Raquel Urtasun. Graph hypernetworks for neural architecture search, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=rkgW0oA9FX>.
- 176 Zhijie Deng, Yinpeng Dong, Jun Zhu. Batch virtual adversarial training for graph convolutional networks, The International Conference on Machine Learning (ICML Workshop), 2019, <https://graphreason.github.io/papers/3.pdf>.
- 177 Daniel Zügner, Stephan Günnemann. Certifiable robustness and robust training for graph convolutional networks, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 246-256.
- 178 Aleksandar Bojchevski, Stephan Günnemann. Certifiable robustness to graph perturbations, 2019, <https://arxiv.org/abs/1910.14356>.
- 179 Vikas Verma, Meng Qu, Alex Lamb, Yoshua Bengio, Juho Kannala, Jian Tang. GraphMix: regularized training of graph neural networks for semi-supervised learning, 2019, <https://arxiv.org/abs/1909.11715>.
- 180 Ke Sun, Piotr Koniusz, Zhen Wang. Fisher-Bures adversary graph convolutional networks, The International Conference on Uncertainty in Artificial Intelligence, 2019, Article No. 161.
- 181 Daniel Zügner, Amir Akbarnejad, Stephan Günnemann. Adversarial attacks on neural networks for graph data, The International Conference on Machine Learning (ICML), 2018, 2847-2856.
- 182 Hanjun Dai, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, Le Song. Adversarial Attack on Graph Structured Data, The International Conference on Machine Learning (ICML), 2018, 1115-1124.
- 183 Daniel Zügner, Stephan Günnemann. Adversarial attacks on graph neural networks via meta learning, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=Bylnx209YX>.
- 184 Dingyuan Zhu, Ziwei Zhang, Peng Cui, Wenwu Zhu. Robust graph convolutional networks against adversarial attacks, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 1399-1407.
- 185 Mikael Henaff, Joan Bruna, Yan LeCun. Deep convolutional networks on graph-structured data, 2015, <https://arxiv.org/abs/1506.05163>.
- 186 Paul Almasan, José Suárez-Varela, Arnau Badia-Sampera, Krzysztof Rusek, Pere Barlet-Ros, Albert Cabellos-Aparicio. Deep reinforcement learning meets graph neural networks: exploring a optical network routing use case, 2019, <https://arxiv.org/abs/1910.07421>.
- 187 Tingwu Wang, Renjie Liao, Jimmy Ba, Sanja Fidler. NerveNet: learning structured policy with graph neural networks, The International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=S1sqHMZCb>.
- 188 Yelong Shen, Jianshu Chen, Po-Sen Huang, Yuqing Guo, Jianfeng Gao. M-Walk: learning to walk over graphs using monte carlo tree search, The International Conference on Neural Information Processing Systems (NIPS), 2018, 6786-6797.
- 189 David K. Hammond, Pierre Vandergheynst, Rémi Gribonval. Wavelets on graphs via spectral graph theory, Applied and Computational Harmonic Analysis, 2011, 30(2): 129-150.
- 190 Dongmian Zou, Gilad Lerman. Graph convolutional neural networks via scattering, Applied and Computational Harmonic Analysis, 2019, <https://doi.org/10.1016/j.acha.2019.06.003>.
- 191 Bingbing Xu, Huawei Shen, Qi Cao, Yunqi Qiu, Xueqi Cheng. Graph Wavelet Neural Network, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=H1ewdiR5tQ>.
- 192 Matan Gavish, Boaz Nadler, Ronald R. Coifman. Multiscale wavelets on trees, graphs and high dimensional data: theory and applications to semi-supervised learning, The International Conference on Machine Learning (ICML), 2010, 367-374.
- 193 Richard C. Wilson, Edwin R. Hancock, Elżbieta Pekalska, Robert P.W. Duin. Spherical and Hyperbolic Embeddings

- of Data, *IEEE Transactions on Pattern Recognition and Machine Intelligence*, 2014, 36(11):2255-2269.
- 194 Claire Donnat, Marinka Zitnik, David Hallac, Jure Leskovec. Learning structural node embeddings via diffusion wavelets, *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2018, 1320-1329.
- 195 Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay S. Pande, Patrick F. Riley. Molecular graph convolutions: moving beyond fingerprints, *Journal of Computer-Aided Molecular Design*, 2016, 30(8): 595-608.
- 196 Seongok Ryu, Jaechang Lim, Seung Hwan Hong, Woo Youn Kim. Deeply learning molecular structure-property relationships using attention- and gate-agumented graph convolutional network, 2018, <https://arxiv.org/abs/1805.10988>.
- 197 Hyeoncheol Cho, Insung S. Choi. Three-dimensionally embedded graph convolutional network (3DGCN) for molecule interpretation, 2018, <https://arxiv.org/abs/1811.09794>.
- 198 Qi Liu, Miltiadis Allamanis, Marc Brockschmidt, Alexander L. Gaunt. Constrained graph variational autoencoders for molecule design, 2018, <https://arxiv.org/abs/1805.09076>.
- 199 Anees Kazi, Shayan Shekarforoush, S. Arvind Krishna, Hendrik Burwinkel, Gerome Vivar, karsten Kortüm, Seyed-Ahmad Ahmadi, Shadi Albarqouni, Nassir Navab. InceptionGCN: receptive field aware graph convolutional network for disease prediction, *The International Conference on Information Processing in Medical Imaging (IPMI)*, 2019, 73-85.
- 200 Wengong Jin, Regina Barzilay, Tommi Jaakkola. Junction tree variational autoencoder for molecular graph generation, 2018, <https://arxiv.org/abs/1802.04364>.
- 201 Marnka Zitnik, Monica Agrawal, Jure Leskovec. Modeling polypharmacy side effects with graph convolutional networks, *Bioinformatics*, 2018, 34(13): i457-i466.
- 202 Connor W. Coley, Wengong Jin, Luke Rogers, Timothy F. Jamison, Tommi S. Jaakkola, William H. Green, Regina Barzilay, Klavs F. Jensen. A graph-convolutional neural network model for the prediction of chemical reactivity, *Chemical Science*, 2019, 10(2): 370-377.
- 203 Nicola De Cao, Thomas Kipf. MolGAN: an implicit generative model for small molecular graphs, 2018, <https://arxiv.org/abs/1805.11973>.
- 204 Sarah Parisot, Sofia Ira Ktena, Enzo Ferrante, Matthew Lee, Ricardo Guerrero, Ben Glocker, Daniel Rueckert. Disease prediction using graph convolutional networks: applications to autism spectrum disorder and Alzheimer's disease, *Medical Image Analysis*, 2018, 48: 117-130.
- 205 Sofia Ira Ktena, Sarah Parisot, Enzo Ferrante, Martin Rajchl, Matthew Lee, Ben Glocker, Daniel Rueckert. Distance metric learning using graph convolutional networks: application to functional brain networks, *The International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2017, 469-477.
- 206 Jeremy Kawahara, Colin J. Brown, Steven P. Miller, Brian G. Booth, Vann Chau, Ruth E. Grunau, Jill G. Zwicker, Ghassan Hamarneh. BrainNetCNN: convolutional neural networks for brain networks; towards predicting neurodevelopment, *NeuroImage*, 2017, 146: 1038-1049.
- 207 Alex Fout, Jonathon Byrd, Basir Shariat, Asa Ben-Hur. Protein interface prediction using graph convolutional networks, *The International Conference on Neural Information Processing Systems*, 2017, 6530-6539.
- 208 Jiaxuan You, Bowen Liu, Rex Ying, Vijay Pande. Graph convolutional policy network for goal-directed molecular graph generation, *The International Conference on Neural Information Processing Systems*, 2018, 6412-6422.
- 209 Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F. Stewart, Jimeng Sun. GRAM: graph-based attention model for healthcare representation learning, *The International Conference on Knowledge Discovery and Data Mining*, 2017, 787-795.
- 210 Peng Han, Peng Yang, Peilin Zhao, Shuo Shang, Yong Liu, Jiayu Zhou, Xin Gao, Panos Kalnis. GCN-MF: disease-gene association identification by graph convolutional networks and matrix factorization, *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*, 2019, 705-713.
- 211 Emanuele Rossi, Federico Monti, Michael Bronstein, Pietro Liò. ncRNA classification with graph convolutional networks, 2019, <https://arxiv.org/abs/1905.06515>.
- 212 Sofia Ira Ktena, Sarah Parisot, Enzo Ferrante, Martin Rajchl, Matthew Lee, Ben Glocker, Daniel Rueckert. Metric learning with spectral graph convolutions on brain connectivity networks, 2018, 169: 431-442.
- 213 Risi Kondor, Truong Song Hy, Horace Pan, Brandon M. Anderson, Shubhendu Trivedi. Covariant compositional networks for learning graphs, *The International Conference on Learning Representations (ICLR Workshop)*, 2018, <https://openreview.net/forum?id=S1TgE7WR->.
- 214 Alex Nowak, Soledad Villar, Afonso S. Bandeira, Joan Bruna. A note on learning algorithms for quadratic assignment with graph neural networks, 2017, <https://arxiv.org/abs/1706.07450v1>.
- 215 Wouter Kool, Herke van Hoof, Max Welling. Attention, learn to solve routing problems, *The International Conference on Learning Representations (ICLR)*, 2019, <https://openreview.net/forum?id=ByxBFRqYm>.
- 216 Marcelo Prates, Pedro H.C. Avelar, Henrique Lemos, Luis C. Lamb, Moshe Y. Vardi. Learning to solve NP-Complete problems: a graph neural network for decision TSP, *The AAAI Conference on Artificial Intelligence (AAAI)*, 2019, 4731-4738.
- 217 Henrique Lemos, Marcelo Prates, Pedro Avelar, Luis Lamb. Graph colouring meets deep learning: effective graph neural network models for combinatorial problems, 2019, <https://arxiv.org/abs/1903.04598>.
- 218 Runzhong Wang, Junchi Yan, Xiaokang Yang. Learning combinatorial embedding networks for deep graph matching, *The IEEE International Conference on Computer Vision (ICCV)*, 2019, 3056-3065.
- 219 Irwan Bello, Hieu Pham, Quoc V. Le, Mohammad Norouzi, Samy Bengio. Neural combinatorial optimization with reinforcement learning, 2017, <https://arxiv.org/abs/1611.09940>.
- 220 Daniel Selsam, Matthew Lamm, Benedikt Bünz, Percy Liang, Leonardo de Moura, David L. Dill. Learning a

- SAT solver from single-bit supervision, The International Conference on Learning Representations (ICLR), 2019, [https://openreview.net/forum?id=HJMC\\_iA5tm](https://openreview.net/forum?id=HJMC_iA5tm).
- 221 Tianshu Yu, Runzhong Wang, Junchi Yan, Baoxin Li. Learning deep graph matching with channel-independent embedding and Hungarian attention, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=rJgBd2NYPH>.
- 222 Zhuwen Li, Qifeng Chen, Vladlen Koltun. Combinatorial optimization with graph convolutional networks and guided tree search, The International Conference on Neural Information Processing Systems (NIPS), 2018, 573-546.
- 223 Ryoma Sato, Makoto Yamada, Hisashi Kashima. Approximation Ratios of Graph Neural Networks for Combinatorial Problems, The International Conference on Neural Information Processing Systems (NeurIPS), 2019, 4081-4090.
- 224 Jiaxuan You, Haoze Wu, Clark Barrett, Raghuram Ramanujan, Jure Leskovec. G2SAT: learning to generate SAT formulas, The International Conference on Neural Information Processing Systems (NeurIPS), 2019, 10552-10563.
- 225 Elias B. Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, Le Song. Learning combinatorial optimization algorithms over graphs, The International Conference on Neural Information Processing Systems (NIPS), 2017, 6348-6358.
- 226 Tatsuro Kawamoto, Masashi Tsubaki, Tomoyuki Obuchi. Mean-field of graph neural networks in graph partitioning, The International Conference on Neural Information Processing Systems (NIPS), 2018, 4366-4376.
- 227 Hongteng Xu, Dixin Luo, Lawrence Carin. Scalable Gromov-Wasserstein learning for graph partitioning and matching, The International Conference on Neural Information Processing Systems (NIPS), 2019, 3046-3056.
- 228 Yunsheng Bai, Hao Ding, Song Bian, Ting Chen, Yizhou Sun, Wei Wang. SimGNN: a neural network approach to fast graph similarity computation, The ACM International Conference on Web Search and Data Mining (WSDM), 2019, 384-392.
- 229 Haoyan Xu, Runjian Chen, Yunsheng Bai, Jie Feng, Ziheng Duan, Ke Luo, Yizhou Sun, Wei Wang. Hierarchical and fast graph similarity computation via graph coarsening and deep graph learning, 2020, <https://arxiv.org/abs/2005.07115>.
- 230 Haoyan Xu, Ziheng Duan, Jie Feng, Runjian Chen, Yida Huang, Yueyang Wang. Graph partitioning and graph neural network based hierarchical graph matching for graph similarity computation, 2020, <https://arxiv.org/abs/2005.08008>.
- 231 Pedro H.C. Avelar, Henrique Lemos, Marcelo O.R. Prates, Luis Lamb. Multitask learning on graph neural networks: learning multiple graph centrality measures with a unified network, 2018, <https://arxiv.org/abs/1809.07695>.
- 232 Wen Zhang, Kai Shu, Huan Liu, Yalin Wang. Graph neural networks for user identity linkage, 2019, <https://arxiv.org/abs/1903.02174v1>.
- 233 Johannes Klicpera, Aleksandar Bojchevski, Stephan Günnemann. Predict the propagate: graph neural networks meet personalized pagerank, The International Conference on Learning Representations (ICLR), <https://openreview.net/forum?id=H1gL-2A9Ym>.
- 234 Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, Jie Tang. DeepInf: social influence prediction with deep learning, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2018, 2110-2119.
- 235 Jinyin Chen, Xuanheng Xu, Yangyang Wu, Haibin Zheng. GC-LSTM: graph convolution embedded LSTM for dynamic link prediction, 2018, <https://arxiv.org/abs/1812.04206>.
- 236 Franco Scarselli, Sweah Liang Yong, Marco Gori, Markus Hagenbuchner, Ah Chung Tsoi, Marco Maggini. Graph neural networks for ranking web pages, The IEEE/WIC/ACM International Conference on Web Intelligence (WI), 2005, 666-672.
- 237 Ben Gao, Vicent W. Zheng, Deng Cai, Xiaofei He, Yueting Zhuang. Link prediction via ranking metric dual-level attention network learning, The International Joint Conference on Artificial Intelligence (IJCAI), 2017, <https://doi.org/10.24963/ijcai.2017/493>.
- 238 Muhan Zhang, Yixin Chen. Link prediction based on graph neural networks, The International Conference on Neural Information Processing Systems (NIPS), 2018, 5171C5181.
- 239 Brian Chen, Bo Wu, Alireza Zareian, Hanwang Zhang, Shih-Fu Chang. General partial label learning via dual bipartite graph autoencoder, The AAAI Conference on Artificial Intelligence, 2020.
- 240 Stephen Phillips, Kostas Daniilidis. All graphs lead to rome: learning geometric and cycle-consistent representations with graph convolutional networks, 2019, <https://arxiv.org/abs/1901.02078>.
- 241 Zhaomin Chen, Xiushen Wei, Peng Wang, Yanwen Guo. Multi-label image recognition with graph convolutional networks, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 5177-5186.
- 242 Guillem Cucurull, Perouz Taslakian, David Vazquez. Context-aware visual compatibility prediction, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 12617-12626.
- 243 Yichao Yan, Qiang Zhang, Bingbing Ni, Wendong Zhang, Minghao Xu, Xiaokang Yang. Learning context graph for person search, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 2158-2167.
- 244 Lei Yang, Xiaohang Zhan, Dapeng Chen, Junjie Yan, Chen Change Loy, Dahua Lin. Learning to cluster faces on an affinity graph, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 2298-2306.
- 245 Zhongdao Wang, Liang Zheng, Yali Li, Shengjin Wang. Linkage based face clustering via graph convolution network, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 1117-1125.
- 246 Yanbin Liu, Juho Lee, Minseop Park, Saehoon Kim, Eunho Yang, Sung Ju Hwang, Yi Yang. Learning to propagate labels: transductive propagation network for few-shot learning, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=SyVuRiC5K7>.

- 247 Jongmin Kim, Taesup Kim, Sungwoong Kim, Chang D. Yoo. Edge-labeling graph neural network for few-shot learning, 2019, 11-20.
- 248 Yin Bi, Aaron Chadha, Alhabib Abbas, Eirina Bourtsoulatze, Yiannis Andreopoulos. Graph-based object classification for neuromorphic vision sensing, The IEEE International Conference on Computer Vision (ICCV), 2019, 491-501.
- 249 Cheyang Si, Wentao Chen, Wei Wang, Liang Wang, Tieniu Tan. An attention enhanced graph convolutional LSTM network for skeleton-based action recognition, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 1227-1236.
- 250 Jiaxing Zhong Nannan Li, Weijie Kong, Shan Liu, Thomas H. Li, Ge Li. Graph convolutional label noise cleaner: train a plug-and-play action classifier for anomaly detection, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 1237-1246.
- 251 Michele Guo, Edward Chou, De-An Huang, Shuran Song, Serena Yeung, Fei-Fei Li. Neural graph matching networks for fewshot 3d action recognition, The European Conference on Computer Vision (ECCV), 2018, 637-689.
- 252 Liang Lin, Lili Huang, Tianshui Chen, Yukang Gan, Hui Cheng. Knowledge-guided recurrent neural network learning for task-oriented action prediction, 2017, <https://arxiv.org/abs/1707.04677>.
- 253 Renchun You, Zhiyao Guo, Lei Cui, Xiang Long, Yingze Bao, Shilei Wen. Cross-modality attention with semantic graph embedding for multi-label classification, The AAAI Conference on Artificial Intelligence (AAAI), 2020.
- 254 Kenneth Marino, Ruslan Salakhutdinov, Abhinav Gupta. The more you know: using knowledge graphs for image classification, The IEEE International Conference on Computer Vision (CVPR), 2017, 20-28.
- 255 Chung-Wei Fang, Chih-Kuan Yeh, Yu-Chiang Frank Wang. Multi-label zero-shot learning with structured knowledge graphs, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 1576-1785.
- 256 Xiaolong Wang, Yufei Ye, Abhinav Gupta. Zero-shot recognition via semantic embedding and knowledge graphs, 2018, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 6857-6866.
- 257 Tianshui Chen, Muxin Xu, Xiaolu Hui, Hefeng Wu, Liang Lin. Learning semantic-specific graph representation for multi-label image recognition, The IEEE International Conference on Computer Vision (ICCV), 2019, 522-531.
- 258 Victor Garcia Satorras, Joan Bruna Estrach. Few-shot learning with graph neural networks, The International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=BJj6qGbRW>.
- 259 Loic Landrieu, Mohamed Boussaha. Point Cloud over Segmentation with graph-structured deep metric learning, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 7440-7449.
- 260 Dong Wook Shu, Sung Woo Park, Junseok Kwon. 3D point cloud generative adversarial network based on tree structured graph convolutions, The IEEE International Conference on Computer Vision (ICCV), 2019, 3859-3868.
- 261 Justin Johnson, Agrim Gupta, Fei-Fei Li. Image generation from scene graphs, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 1219-1228.
- 262 Yikang Li, Wanli Quyang, Bolei Zhou, Jianping Shi, Chao Zhang, Xiaogang Wang. Factorizable Net: an efficient subgraph-based framework for scene graph generation, The European Conference on Computer Vision (ECCV), 2018.
- 263 Danfei Xu, Yuke Zhu, Christopher B. Choy, Fei-Fei Li. Scene graph generation by iterative message passing, 2017, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 5410-5419.
- 264 Jianwei Yang, Jiasen Lu, Stefan Lee, Dhruv Batra, Devi Parikh. Graph R-CNN for scene graph generation, The European Conference on Computer Vision (ECCV), 2018.
- 265 Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein. Dynamic graph CNN for learning on point clouds. *ACM Transactions on Graphics*, 2019, 38(5): Article No. 146.
- 266 Gusi Te, Wei Hu, Amin Zheng, Zongming Guo. RGCNN: regularized graph cnn for point cloud segmentation, The ACM International Conference on Multimedia (MM), 2018, 746-754.
- 267 Xiaodan Liang, Liang Lin, Xiaohui Shen, Jiashi Feng, Shuicheng Yan, Eric P. Xing. Interpretable structure-evolving LSTM, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2017, 1010-1019.
- 268 Li Yi, Hao Su, Xingwen Guo, Leonidas J. Guibas. SyncSpecCNN: synchronized spectral CNN for 3D shape segmentation, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2017, 2282-2290.
- 269 Loic Landrieu, Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 4558-4567.
- 270 Yongfei Liu, Bo Wan, Xiaodan Zhu, Xuming He. Learning cross-modal context graph for visual grounding, The AAAI Conference on Artificial Intelligence (AAAI), 2020, <https://arxiv.org/abs/1911.09042>.
- 271 Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, Jie Tang. Cognitive graph for multi-hop reading comprehension at scale, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 2694-2703.
- 272 Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li, Weinan Zhang, Yong yu. Dynamically fused graph network for multi-hop reasoning, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 6140-6150.
- 273 Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, Maosong Sun. GEAR: graph-based evidence aggregating and reasoning for fact verification, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 892-901.
- 274 Zhenghao Liu, Chenyan Xiong, Maosong Sun. Kernel graph attention network for fact verification, 2019, <https://arxiv.org/abs/1910.09796v1>.
- 275 Nicola De Cao, Wilker Aziz, Ivan Titov. Question answering by reasoning across documents with graph convolutional networks, The Annual Conference of the North American Chapter of the Association for Computational Linguistics, 2019, 2306-2317.
- 276 Daesik Kim, Seonhoon Kim, Nojun Kwak. Textbook question answering with multi-modal context graph understanding and self-supervised open-set comprehension, The Annual Meeting of the Association for Computational Linguistics

- (ACL), 2019, 3568-3584.
- 277 Damien Teney, Lingqiao Liu, Anton van den Hengel. Graph-structured representations for visual question answering, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2017, 1-9.
- 278 Ronghang Hu, Anna Rohrbach, Trevor Darrell, Kate Saenko. Language-conditioned graph networks for relational reasoning, The IEEE International Conference on Computer Vision (ICCV), 2019, 10294-10303.
- 279 Linjie Li, Zhe Gan Yu Cheng, Jingjing Liu. Relation-aware graph attention network for visual question answering. The International Conference on Computer Vision (ICCV), 2019, 10313-10322.
- 280 Yangyang Cheng, Chun Yuan. Reasoning-aware graph convolutional network for visual question answering, 2020, <https://openreview.net/forum?id=SkpPvISYwS>.
- 281 Yu Cao, Meng Fang, Dacheng Tao. BAG: Bi-directional attention entity graph convolutional network for multi-hop reasoning question answering, The Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 2019, 357-362.
- 282 Medhini Narasimhan, Svetlana Lazebnik, Alexander G. Schwing. Out of the box: reasoning with graph convolution nets for factual visual question answering, The International Conference on Neural Information Processing Systems (NIPS), 2018, 2654-2665.
- 283 Xinlei Chen, Li-Jia Li, Fei-Fei Li, Abhinav Gupta. Iterative visual reasoning beyond convolutions, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2018, 7239-7248.
- 284 David Raposo, Adam Santoro, David Barrett, Razvan Pascanu, Timonhy Lillicrap, Peter Battaglia. Discovering objects and their relations from entangled scene representations, 2017, <https://arxiv.org/abs/1702.05068v1>.
- 285 Adam Santoro, David Raposo, David G.T. Barrett, Mateusz Malinowski, Razvan Pascanu, Peter W. Battaglia. A simple neural network module for relational reasoning, The International Conference on Neural Information Processing Systems (NIPS), 2017, 4967-4976.
- 286 Xu han, Zhiyuan Liu, Maosong Sun. Neural knowledge acquisition via mutual attention between knowledge graph and text, The AAAI Conference on Artificial Intelligence (AAAI), 2018, 4832-4839.
- 287 Huaxiu Yao, Chuxu Zhang, Ying Wei, Meng Jiang, Suhang Wang, Junzhou Huang, Nitesh V. Chawla, Zhenhui Li. Graph few-shot learning via knowledge transfer, The AAAI Conference on Artificial Intelligence (AAAI), 2020, <https://arxiv.org/abs/1910.03053>.
- 288 Daniel Oñoro-Rubio, Mathias Niepert, Alberto García-Durán, Roberto González, Roberto J. López-Sastre. Representation learning for visual-relational knowledge graphs, 2017, <https://arxiv.org/abs/1709.02314v5>.
- 289 Bill Yuchen Lin, Xinyue Chen, Jamin Chen, Xiang Ren. KagNet: knowledge-aware graph networks for commonsense reasoning, The International Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 2829-2839.
- 290 Mayu Otani, Yuta Nakashima, Esa Rahtu, Janne Heikkilä. Rethinking the evaluation of video summaries, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2019, 7596-7604.
- 291 Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne van den Berg, Ivan Titov, Max Welling. Modeling relational data with graph convolutional networks, The European Semantic Web Conference (ESWC), 2018, 593-607.
- 292 Takuo Hamaguchi, Hidekazu Oiwa, Masashi Shimbo, Yuji Matsumoto. Knowledge transfer for out-of-knowledge-base entities: a graph neural network approach, The International Joint Conference on Artificial Intelligence (IJCAI), 2017, 1802-1808.
- 293 Hao Zhou, Tom Young, Minlie Huang, Haizhou Zhao, Jingfang Xu, Xiaoyan Zhu. Commonsense knowledge aware conversation generation with graph attention, The International Joint Conference on Artificial Intelligence (IJCAI), 2019, 4623-4629.
- 294 Zhouxia Wang, Tianshui Chen, Jimmy Ren, Weihao Yu, Hui Cheng, Liang Lin. Deep reasoning with knowledge graph for social relationship understanding, The International Joint Conference on Artificial Intelligence (IJCAI), 2018, 1021-1028.
- 295 Namyong Park, Andrey Kan, Xin Luna Dong, Tong Zhao, Christos Faloutsos. Estimating node importance in knowledge graphs using graph neural networks, The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 596-606.
- 296 Shikhar Vashishta, Shib Sankar Dasgupta, Swayambhu Nath Ray, Partha Talukdar. Dating documents using graph convolution networks, The Annual Meeting of the Association for Computational Linguistics (ACL), 2018, 1605-1615.
- 297 Diego Marcheggiani, Ivan Titov. Encoding sentences with graph convolutional networks for semantic role labeling, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2017, 1506-1515.
- 298 Shikhar Vashishta, Manik Bhandari, Prateek Yadav, Piyush Rai, Chiranjib Bhattacharyya, Partha Talukdar. Incorporating syntactic and semantic information in word embeddings using graph convolutional networks, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 3308-3318.
- 299 Thien Huu Nguyen, Ralph Grishman. Graph convolutional networks with argument-aware pooling for event detection, The AAAI Conference on Artificial Intelligence (AAAI), 2018, 5900-5907.
- 300 Xiao Liu, Zhunchen Luo, Heyan Huang. Jointly multiple events extraction via attention-based graph information, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2018, 1247-1256.
- 301 Ao Li, Zhou Qin, Runshi Liu, Yiqun Yang, Dong Li. Spam review detection with graph convolutional networks, The International Conference on Information and Knowledge Management (CIKM), 2019, 2703-2711.
- 302 Pushkar Mishra, Marco Del Tredici, Helen Yannakoudakis, Ekaterina Shutova. Abusive Language Detection with Graph Convolutional Networks, The Conference of the North American Chapter of the Association for Computational Linguistics, 2019, 2145-2150.

- 303 Daniel Beck, Gholamreza Haffari, Trevor Cohn. Graph-to-Sequence learning using gated graph neural networks, The Annual Meeting of the Association for Computational Linguistics, 2018, 273-283.
- 304 Thang Luong, Hieu Pham, Christopher D. Manning. Effective approaches to attention-based neural machine translation, The Conference on Empirical Methods in Natural Language Processing, 2015, 1412-1421.
- 305 Joost Bastings, Ivan Titov, Wilker Aziz, Diego Marcheggiani, Khalil Simáan. Graph convolutional encoders for syntax-aware neural machine translation, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2017, 1957-1967.
- 306 Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio. Neural machine translation by jointly learning to align and translate, The International Conference on Learning Representations (ICLR), 2015, <https://arxiv.org/abs/1409.0473>.
- 307 Diego Marcheggiani, Joost Bastings, Ivan Titov. Exploiting semantics in neural machine translation with graph convolutional networks, The Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 2018, 486-492.
- 308 Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, Jie Tang. Cognitive graph for multi-hop reading comprehension at scale, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 2694-2703.
- 309 Ming Tu, Guangtao Wang, Jing Huang, Yun Tang, Xiaodong He, Bowen Zhou. Multi-hop reading comprehension across multiple documents by reasoning over heterogeneous graphs, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 2704-2713.
- 310 Linfeng Song, Zhiguo Wang, Mo Yu, Yue Zhang, Radu Florian, Daniel Gildea. Exploring graph-structured passage representation for multi-hop reading comprehension with graph neural networks, 2018, <https://arxiv.org/abs/1809.02040>.
- 311 Yu Chen, Lingfei Wu, Mohammed J. Zaki. GraphFlow: exploiting conversation flow with graph neural networks for conversational machine comprehension, 2020, <https://arxiv.org/abs/1908.00059>.
- 312 Zhijiang Guo, Yan Zhang, Wei Lu. Attention guided graph convolutional networks for relation extraction, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 241-251.
- 313 Hao Zhu, Yankai Lin, Zhiyuan Liu, Jie Fu, Tat-Seng Chua, Maosong Sun. Graph neural networks with generated parameters for relation extraction, The Annual Meeting of the Association for Computational Linguistics, 2019, 1331-1339.
- 314 Makoto Miwa, Mohit Bansal. End-to-End relation extraction using LSTMs on sequences and tree structures, The Annual Meeting of the Association for Computational Linguistics (ACL), 2016, 1105-1116.
- 315 Yuhao Zhang, Peng Qi, Christopher D. Manning. Graph convolutional over pruned dependency trees improves relation extraction, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2018, 2205C2215.
- 316 Linfeng Song, Yue Zhang, Zhiguo Wang, Daniel Gildea. N-ary relation extraction using graph-state LSTM, The Conference on Empirical Methods in Natural Language Processing, 2018, 2226-2235.
- 317 Ningyu Zhang, Shuming Deng, Zhanlin Sun, Guanying Wang, Xi Chen, Wei Zhang, HuaJun Chen. Long-tail relation extraction via knowledge graph embeddings and graph convolution networks, The Conference of the North American Chapter of the Association for Computational Linguistics (NAACL), 2019, 3016-3025.
- 318 Liang Yao, Chengsheng Mao, Yuan Luo. Graph convolutional networks for text classification, The AAAI Conference on Artificial Intelligence (AAAI), 2019, 7370-7377.
- 319 Chang Li, Dan Goldwasser. Encoding social information with graph convolutional networks for political perspective detection in News Media, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 2594-2604.
- 320 Xien Liu, Xinxin You, Xiao Zhang, Ji Wu, Ping Lv. Tensor graph convolutional networks for text classification, 2020, <https://arxiv.org/abs/2001.05313>.
- 321 Chen Zhang, Qiuchi Li, Dawei Song. Aspect-based sentiment classification with aspect-specific graph convolutional networks, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 4568-4578.
- 322 Deepanway Ghosal, Navonil Majumder, Soujanya Poria, Niyati Chhaya, Alexander Gelbukh. DialogueGCN: a graph convolutional neural network for emotion recognition in conversation, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 154-164.
- 323 Hu Linmei, Tianchi Yang, Chuan Shi, Houye Ji, Xiaoli Li. Heterogeneous graph attention networks for semi-supervised short text classification, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 4821-4830.
- 324 Binxuan Huang, Kathleen Carley. Syntax-aware aspect level sentiment classification with graph attention networks, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 5469-5477.
- 325 Hao Peng, Jianxin Li, Yu He, Yaopeng Liu, mengjiao Bao, Lihong Wang, Yangqiu Song. Large-scale hierarchical text classification with recursively regularized deep graph-CNN, The World Wide Web Conference (WWW), 2018, 1063-1072.
- 326 Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, Hannaneh Hajishirzi. Text generation from knowledge graphs with graph transformers, 2019, <https://arxiv.org/abs/1904.02342>.
- 327 Wei Li, Jingjing Xu, Yancheng He, Shengli Yan, Yunfang Wu, Xu Sun. Coherent comment generation for chinese articles with a graph-to-sequence model, The Annual Meeting of the Association for Computational Linguistics (ACL), 2019, 4843-4852.
- 328 Linfeng Song, Yue Zhang, Zhiguo Wang, Daniel Gildea. A graph-to-sequence model for AMR-to-text generation, The Annual Meeting of the Association for Computational Linguistics (ACL), 2018, 1616-1626.
- 329 Leonardo F.R. Ribeiro, Claire Gardent, Iryna Gurevych. Enhancing AMR-to-text generation with dual graph representations, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 3183-3194.

- 330 Yu Chen, Lingfei Wu, Mohammed J. Zaki. Reinforcement learning based graph-to-sequence model for natural question generation, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=HygnDhEtvR>.
- 331 Tao Gui, Yicheng Zou, Qi Zhang, Minlong Peng, Jinlan Fu, Zhongyu Wei, Xuanjing Huang. A lexicon-based graph neural network for Chinese NER, The Conference on Empirical Methods in Natural Language Processing (EMNLP), 2019, 1040-1050.
- 332 Jessica B. Hamrick, Kelsey R. Allen, Victor Bapst, Tina Zhu, Kevin R. McKee, Joshua B. Tenenbaum, Peter W. Battaglia. Relational inductive bias for physical construction in humans and machines, 2018, <https://arxiv.org/abs/1806.01203>.
- 333 Michael Chang, Tomer Ullman, Antonia Torralba, Joshua Tenenbaum. A compositional object-based approach to learning physical dynamics, The International Conference on Learning Representations (ICLR), 2017, <https://openreview.net/forum?id=Bkab5dqxe>.
- 334 Thomas N. Kipf, Elise van der Pol, Max Welling. Contrastive learning of structured world models, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=H1gax6VtDB>.
- 335 Thomas N. Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling. Neural relational inference for interacting systems, The International Conference on Machine Learning (ICML), 2018, 2688-2697.
- 336 Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, Koray Kavukcuoglu. Interaction networks for learning about objects, relations and physics, The International Conference on Neural Information Processing Systems (NIPS), 2016, 4502-4510.
- 337 Sainbayar Sukhbaatar, Arthur Szlam, Rob Fergus. Learning multiagent communication with backpropagation, The International Conference on Neural Information Processing Systems (NIPS), 2016, 2252-2260.
- 338 Yedid Hoshen. VAIN: attentional multi-agent predictive modeling, The International Conference on Neural Information Processing Systems (NIPS), 2017, 2698-2708.
- 339 Nicholas Watters, Andrea Tacchetti, Théophane Weber, Razvan Pascanu, Peter Battaglia, Daniel Zoran. Visual interaction networks: learning a physics simulator from video, The International Conference on Neural Information Processing Systems (NIPS), 2017, 4539-4547.
- 340 Miltiadis Allamanis, Marc Brockschmidt, Mahmoud Khademi. Learning to represent programs with graphs, The International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=BJOFETxR->.
- 341 Jiayi Wei, Maruth Goyal, Greg Durrett, Isil Dilling. LambdaNet: probabilistic type inference using graph neural networks, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=Hkx6hANtwH>.
- 342 Jessica Schrouff, Kai Wohlfahrt, Bruno Marnette, Liam Atkinson. Inferring Javascript types using graph neural networks, 2019, <https://arxiv.org/abs/1905.06707>.
- 343 Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, Tieniu Tan. Session-based recommendation with graph neural networks, The AAAI Conference on Artificial Intelligence (AAAI), 2019, 346-353.
- 344 Xiao Wang, Ruijia Wang, Chuan Shi, Guojie Song, Qingyong Li. Multi-component graph convolutional collaborative filtering, The AAAI Conference on Artificial Intelligence (AAAI), 2020.
- 345 Shaohua Fan, Junxiong Zhu, Xiaotian Han, Chuan Shi, Linmei Hu, Biyu Ma, Yongliang Li. Metapath-guided heterogeneous graph neural network for intent recommendation, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 2478-2486.
- 346 Yinwei Wei Xiang Wang, Liqiang Nie, Xiangnan He, Richang Hong, Tat-Seng Chua. MMGCN multi-modal graph convolution network for personalized recommendation of micro-video, The ACM International Conference on Multimedia (MM), 2019, 1437-1445.
- 347 Rianne van den Berg, Thomas N. Kipf, Max Welling. Graph convolutional matrix completion, 2017, <https://arxiv.org/abs/1706.02263>.
- 348 Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2018, 974-983.
- 349 Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, Tat-Seng Chua. KGAT: knowledge graph attention network for recommendation, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 950-958.
- 350 Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, Zhongyuan Wang. Knowledge-aware graph neural networks with label smoothness regularization for recommender systems, The International Conference on Knowledge Discovery and Data Mining (SIGKDD), 2019, 968-977.
- 351 Qitian Wu, Hengrui Zhang, Xiaofeng Gao, Peng He, Paul Weng, Han Gao, Guihai Chen. Dual graph attention networks for deep latent representation of multifaceted social effects in recommender systems, The World Wide Web Conference (WWW), 2019, 2091-2102.
- 352 Wenqi Fan, Yao Ma, Qing Li, Yuan He, Eric Zhao, Jiliang Tang, Dawei Yin. Graph neural networks for social recommendation, The World Wide Web Conference (WWW), 2019, 417-426.
- 353 Chuanpan Zheng, Xiaoliang Fan, Cheng Wang, Jianzhong Qi. GMAN: a graph multi-attention network for traffic prediction, The AAAI Conference on Artificial Intelligence (AAAI), 2020.
- 354 Cheonbok Park, Chunggi Le, Hyojin Bahng, Taeyun Won, Kihwan Kim, Seungmin Jin, Sungahn Ko, Jaegul Choo. STAGRAT: a spatio-temporal graph attention network for traffic forecasting, 2019, <https://arxiv.org/abs/1911.13181>.
- 355 Songtao he, Favyen Bastani, Satvat Jagwani, Edward Park, Sofiane Abbar, Mohammad Alizadeh, Hari Balakrishnan, Sanjay Chawla, Samuel Madden, Mohammad Amin Sadeghi. RoadTagger: robust road attribute inference with graph



- neural networks, The AAAI Conference on Artificial Intelligence (AAAI), 2020.
- 356 Zhiyong Cui, Kristian Henrickson, Ruimin Ke, Yinhai Wang. Traffic graph convolutional recurrent neural network: a deep learning framework for network-scale traffic learning and forecasting, *IEEE Transactions on Intelligent Transportation Systems (TITS)*, 2019, 1-12.
- 357 Frederik Diehl, Thomas Brunner, Michael Truong Le, Alois Knoll. Graph neural networks for modelling traffic participant interaction, 2019, <https://arxiv.org/abs/1903.01254>.
- 358 Yaguang Li, Rose Yu, Cyrus Shahabi, Yan Liu. Diffusion convolutional recurrent neural network: data-driven traffic forecasting, *The International Conference on Learning Representations (ICLR)*, 2018, <https://openreview.net/forum?id=SJiHXGWAZ>.
- 359 Jean-Baptiste Cordonnier, Andreas Loukas. Extrapolating paths with graph neural networks, *The International Joint Conference on Artificial Intelligence (IJCAI)*, 2019, 2187-2194.
- 360 Ryoma Sato. A survey on the expressive power of graph neural networks, 2020, <https://arxiv.org/abs/2003.04078>.
- 361 Dan Busbridge, Dane Sherburn, Pietro Cavallo, Nils Y. Yammerla. Relational graph attention networks, 2019, <https://arxiv.org/abs/1904.05811>.
- 362 Yu Rong, Tingyang Xu, Wenbing Huang, Somayeh Sojoudi, Junzhou Huang, Wenwu Zhu. Spectral graph attention network, 2020, <https://arxiv.org/abs/2003.07450>.
- 363 Tinghuai Wang, Guangming Wang, Kuan Eeik Tan, Donghui Tan. Spectral pyramid graph attention network for hyperspectral image classification, 2020, <https://arxiv.org/abs/2001.07108>.
- 364 Wei Liu, Sanjay Chawla. A game theoretical model for adversarial learning, *The IEEE International Conference on Data Mining Workshops (ICDM Workshop)*, 2009, 25-30.
- 365 Lichao Sun, Yingdong Dou, Carl Yang, Ji Wang, Philip S. Yu, Bo Li. Adversarial attack and defense on graph data: a survey, 2020, <https://arxiv.org/abs/1812.10528>.
- 366 Davide Bacciu, Federico Errica, Alessio Micheli, Marco Podda. A gentle introduction to deep learning for graphs, 2019, <https://arxiv.org/abs/1912.12693>.
- 367 Haopeng Zhang, Congying Xia, Li Sun. Graph-Bert: only attention is needed for learning graph representations, 2020, <https://arxiv.org/abs/2001.05140>.
- 368 Shizhe Chen, Yida Zhao, Qin Jin, Qi Wu. Fine-grained video-text retrieval with hierarchical graph reasoning, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, <https://arxiv.org/abs/2003.00392>.
- 369 Shizhe Chen, Qin Jin, Peng Wang, Qi Wu. Say as you wish: fine-grained control of image caption generation with abstract scene graphs, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, <https://arxiv.org/abs/2003.00387>.
- 370 Oytun Ulutan, ASM Iftekhar, B.S. Manjunath. VSGNet: spatial attention network for detecting human object interactions using graph convolutions, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, <https://arxiv.org/abs/2003.05541>.
- 371 Jiangke Lin, Yi Yuan, Tianjie Shao, Kun Zhou. Towards high-fidelity 3d face reconstruction from in-the-wild images using graph convolutional networks, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, <https://arxiv.org/abs/2003.05653>.
- 372 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, Jure Leskovec. Open graph benchmark: datasets for machine learning on graphs, 2020, <https://arxiv.org/abs/2005.00687>.
- 373 Lingxiao Zhao, Leman Akoglu. PairNorm: tackling oversmoothing in GNNs, *The International Conference on Learning Representations (ICLR)*, 2020, <https://openreview.net/forum?id=rkecl1rtwB>.
- 374 Fan R.K. Chung. Spectral Graph Theory, American Mathematical Society, 1992.
- 375 Sami Abu-EL-Haija, Amol Kapoor, Bryan Perozzi, Joonseok Lee. N-GCN: multi-scale graph convolution for semi-supervised node classification, *The Conference on Uncertainty in Artificial Intelligence (UAI)*, 2019: Article No. 310.
- 376 Antti Tarvainen, Harri Valpola. Mean teachers are better role models: weight-averaged consistency targets improve semi-supervised deep learning results, *The International Conference on Neural Information Processing Systems (NIPS)*, 2017: 1195-1204.
- 377 Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization, *The International Conference on Neural Information Processing Systems (NIPS)*, 2014, 2177-2185.
- 378 Peter D. Turney, Patrick Pantel. From frequency to meaning: vector space models of semantics. *Journal of Artificial Intelligence Research*, 2010, 37(2010): 141-188.
- 379 Yu Rong, Wenbing Huang, Tingyang Xu, Junzhou Huang. DropEdge: towards deep graph convolutional networks on node classification, *The International Conference on Learning Representations (ICLR)*, 2020, <https://openreview.net/group?id=ICLR.cc/2020/Conference>.
- 380 Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, Gabriele Monfardini. The graph neural network model, *IEEE Transactions on Neural Networks*, 2008, 20(1):61-80.
- 381 Antoni Buades, Bartomeu Coll, Jean-Michael Morel. A non-local algorithm for image denoising, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, 60-65.
- 382 Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun. Deep residual learning for image recognition, *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, 1063-6919.
- 383 Gao Huang, Zhuang Liu, Laurens Van Der Maaten, Kilian Q. Weinberger. Densely connected convolutional networks, *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, 1063-6919.
- 384 Fisher Yu, Vladlen Koltun. Multi-scale context aggregation by dilated convolutions, *The International Conference on*

- Learning Representations, 2016, <https://arxiv.org/abs/1511.07122>.
- 385 Bert De Brabandere, Xu Jia, Tinne Tuytelaars, Luc Van Gool. Dynamic filter networks, The International Conference on Neural Information Processing Systems (NIPS), 2016, 667-675.
- 386 Kilian Weinberger, Anirban Dasgupta, John Langford. Feature hashing for large scale multitask learning, The International Conference on Machine Learning, 2009, 1113-1120.
- 387 Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, David Lopez-Pax. Mixup: beyond empirical risk minimization, The International Conference on Learning Representations (ICLR), 2018, <https://openreview.net/forum?id=r1Ddp1-Rb>.
- 388 Vikas Verma, Alex Lamb, Christopher Beckham, Amir Najafi, Ioannis Mitliagkas, David Lopez-Pax, Yoshua Bengio. Manifold Mixup: better representations by interpolating hidden states, The International Conference on Machine Learning (ICML), 2019, 6438-6447.
- 389 Haggai Maron, Heli Ben-Hamu, Nadav Shamir, Yaron Lipman. Invariant and equivariant graph networks, The International Conference on Learning Representations (ICLR), 2019, <https://openreview.net/forum?id=Syx72jC9tm>.
- 390 Hao Yuan, Shuiwang Ji. StructPool: structured graph pooling via conditional random fields, The International Conference on Learning Representations (ICLR), 2020, [https://openreview.net/forum?id=BJxg\\_hVtwH](https://openreview.net/forum?id=BJxg_hVtwH).
- 391 John Lafferty, Andrew McCallum, Fernando CN Pereira. Conditional random field: probabilistic models for segmenting and labeling sequence data, The International Conference on Machine Learning (ICML), 2001, 282-289.
- 392 Dzmitry Bahdanau, kyungHyun Cho, Yoshua Bengio. Neural machine translation by jointly learning to align and translate, The International Conference on Learning Representations, 2015, <https://arxiv.org/abs/1409.0473>.
- 393 Volodymyr Minh, Nicolas Heess, Alex Graves, Koray Kavukcuoglu. Recurrent models of visual attention, The International Conference on Neural Information Processing Systems, 2014, 2204-2212.
- 394 John Boaz, Ryan A. Rossi, Sungchul Kim, Nesreen K. Ahmed, Eunye Koh. Attention models in graphs: a survey, ACM Transactions on Knowledge Discovery from Data, 2019, 13(6): Article No. 62.
- 395 Yu Zhou, Jianbin Huang, Heli Sun, Yizhou Sun, Shaojie Qiao, Stephen Wambura. Recurrent meta-structure for robust similarity measure in heterogeneous information networks, ACM Transactions on Knowledge Discovery from Data, 2019, 13(6): Article No. 64.
- 396 Chuan Shi, Yitong Li, Jiawei Zhang, Yizhou Sun, Philip S. Yu. A survey of heterogeneous information network analysis, IEEE Transactions on Knowledge and Data Engineering, 2017, 29(1): 17-37.
- 397 Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov Dumitru Erhan, Vincent Vanhoucke, Andrew Rabinovich. Going deeper with convolutions, The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2015, 1-9.
- 398 Afshin Rahimi, Trevor Cohn, Timothy Baldwin. Semi-supervised user geolocation via graph convolutional networks, The Annual Meeting of the Association for Computational Linguistics (ACL), 2015, 2009-2019.
- 399 Rupesh Kumar Srivastava, Klaus Greff, Jürgen Schmidhuber. Highway networks, <https://arxiv.org/abs/1505.00387>.
- 400 Sara Sabour, Nicholas Frosst, Geoffrey E. Hinton. Dynamic routing between capsules, The International Conference on Neural Information Processing Systems (NIPS), 2017, 3856-3866.
- 401 Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S. Du, Ken-ichi Kawarabayashi, Stefanie Jegelka. What can neural networks reason about, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=rJxbJeHFPS>.
- 402 Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang, Chao Ma, Ziyue Huang, Qipeng Guo, Hao Zhang, Haibin Lin, Junbo Zhao, Jinyang Li, Alexander Smola, Zheng Zhang. Deep Graph Library: towards efficient and scalable deep learning on graphs, The Workshop at the International Conference on Learning Representations (ICLR Workshop), 2019, <https://openreview.net/forum?id=q9RwOO-Ci5..>
- 403 Rui Dai, Shenkun Xu, Qian Gu, Chenguang Ji, Kaikui Liu. Hybrid spatio-temporal graph convolutional networks: improving traffic prediction with navigation data, The International Conference on Knowledge Discovery and Data Mining, 2020.
- 404 Xiaoran Xu, Wei Feng, Yunsheng Jiang, Xiaohui Xie, Zhiqing Sun, Zhihong Zhang. Dynamically pruned message passing networks for large-scale knowledge graph reasoning, The International Conference on Learning Representations (ICLR), 2020, <https://openreview.net/forum?id=rkeuAhVKvB&noteId=rkeuAhVKvB>.
- 405 Xiaojuan Qi, Renjie Liao, Jiaya Jia, Sanja Fidler, Raquel Urtasun. 3D graph neural networks for RGBD semantic segmentation, The International Conference on Computer Vision and Pattern Recognition (CVPR), 2017, 5199-5208.
- 406 Tsu-Jui Fu, Peng-Hsuan Li, Wei-Yun Ma. GraphRel: modeling text as relational graphs for joint entity and relation extraction. The Annual Meeting of the Association for Computational Linguistics, 2019, 1409-1418.