# Graph Neural Networks: Taxonomy, Advances and Trends

YU ZHOU*, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, China

HAIXIA ZHENG†, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, China

XIN HUANG, College of Data Science, Taiyuan University of Technology, China

SHUFENG HAO, College of Data Science, Taiyuan University of Technology, China

DENGAO LI, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, China

JUMIN ZHAO, College of Information and Computer/Shanxi Intelligent Perception Engineering Research Center, Taiyuan University of Technology, China

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 systematic and 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 250 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.

CCS Concepts: • **Computing methodologies → Neural networks**; **Learning latent representations**.

Additional Key Words and Phrases: Graph Convolutional Neural Network, Graph Recurrent Neural Network, Graph Pooling Operator, Graph Attention Mechanism, Graph Neural Network

---

*Corresponding Author.

†Corresponding Author.

---

Authors' addresses: Yu Zhou, zhouyu@tyut.edu.cn, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024; Haixia Zheng, zhenghaixia@tyut.edu.cn, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024; Xin Huang, huangxin@tyut.edu.cn, College of Data Science, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024; Shufeng Hao, haoshufeng@tyut.edu.cn, College of Data Science, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024; Dengao Li, lidengao@tyut.edu.cn, College of Data Science/Shanxi Spatial Information Network Engineering Technology Research Center, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024; Jumin Zhao, zhaojumin@tyut.edu.cn, College of Information and Computer/Shanxi Intelligent Perception Engineering Research Center, Taiyuan University of Technology, No. 79 Yingze West Street, WanBaiLin District, Taiyuan, Shanxi, China, 030024.

## 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 computa-
tional 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.

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[124]. Peng Cui et al[249] 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[81] 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[250] 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[75] 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
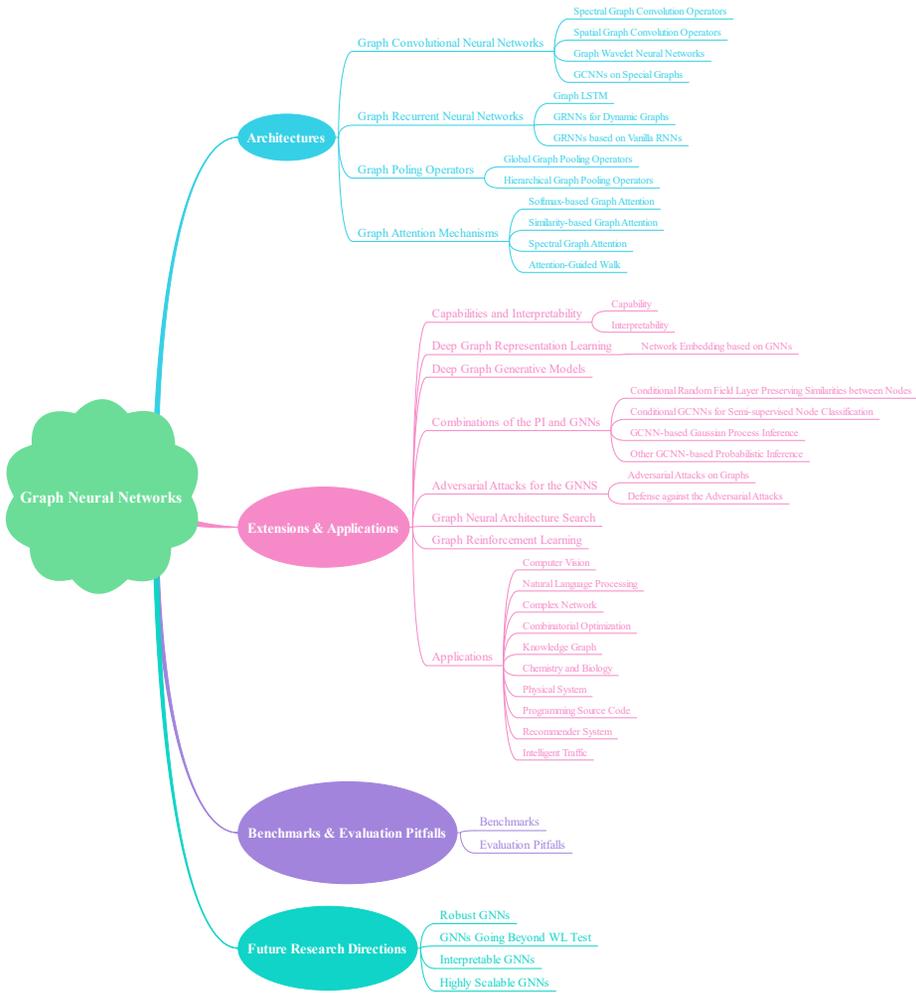
Fig. 1. The architecture of this paper.

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 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 $V = \{v_1, \cdots, v_N\}$ and $E = \{e_1, \cdots, e_M\}$. Each edge $e_j \in E$ can be denoted by $e_j = \left(v_{s_j}, v_{r_j}\right)$ 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, $\left(v_{s_j}, v_{r_j}\right) \neq \left(v_{r_j}, v_{s_j}\right)$ and therefore $A_G$ is asymmetric. A directed edge $e_j = \left(v_{s_j}, v_{r_j}\right)$ is also called an arch, i.e. $e_j = \langle v_{s_j}, v_{s_j} \rangle$. Otherwise $\left(v_{s_j}, v_{r_j}\right) = \left(v_{r_j}, v_{s_j}\right)$ 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 $a = (a_1, \cdots, a_N) \in \mathbb{R}^N$, $\text{diag}(a)$ (or $\text{diag}(a_1, \cdots, a_N)$) denotes a diagonal matrix consisting of the elements $a_n, n = 1, \cdots, N$.

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 $d$-dimensiaonl 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

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

The normalization of $L_G$ is defined by $\overline{L}_G = \mathbb{I}_N - D_G^{-\frac{1}{2}} A_G D_G^{-\frac{1}{2}}$. $\overline{L}_G$ is a real symmetric semi-positive definite matrix. So, it has $N$ ordered real non-negative eigenvalues $\{\lambda_n : n = 1, \cdots, N\}$ and corresponding orthonormal eigenvectors $\{u_n : n = 1, \cdots, N\}$, namely $\overline{L}_G = U \Lambda U^T$ where $\Lambda = \text{diag}(\lambda_1, \cdots, \lambda_N)$ and $U = (u_1, \cdots, u_N)$ denotes a orthonormomal matrix. Without loss of

generality, $0 = \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_N = \lambda_{\max}$. The eigenvectors $u_n, n = 1, \cdots, N$ 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[35] for a given graph signal $x$ can be denoted by

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

The inverse graph Fourier transform can be correspondingly denoted by

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

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

# 3 ARCHITECTURES

## 3.1 Graph Convolutional Neural Networks (GCNNs)

The GCNNs play pivotal roles on tackling the irregular data (e.g. graph and manifold). They 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 [120, 125, 228]. 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.



Fig. 2. Computational framework of the spectral GCNN.

*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

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

where $\circledast$ denotes the element-wise Hadamard product [45, 85, 125]. 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

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

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 $\overline{L}_G$ especially for large graphs. The article [193] proposes a regularization technique, namely GraphMix, to augment the vanilla GCNN with a parameter-sharing Fully Connected Network (FCN).

**Spectral Graph Filter.** Many studies [125] focus on designing different spectral graph filters. In order to circumvent the eigendecomposition, the spectral graph filter $f_\theta$ can formulated as a $K$-localized polynomial of the eigenvalues of the normalized graph Laplacian $\overline{L}_G$ [123, 160, 185], i.e.

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

In practice, the $K$-localized Chebyshev polynomial [123] is a favorable choice of formulating the spectral graph filter, i.e.

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

where the Chebyshev polynomial is defined as

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_k(x) = 2x T_{k-1}(x) - T_{k-2}(x) \tag{7}$$

and $\widetilde{\Lambda} = \frac{2}{\lambda_{\max}}\Lambda - \mathbb{I}_N$. The reason why $\widetilde{\Lambda} = \frac{2}{\lambda_{\max}}\Lambda - \mathbb{I}_N$ 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 $\overline{L}_G$ is less than or equal to 2 [24]) and $\theta = \theta_0 = -\theta_1$. Moreover, the renormalization trick is used here to mitigate the limitations of the vanishing/exploding gradient, namely substituting $\widetilde{D}_G^{-\frac{1}{2}} \widetilde{A}_G \widetilde{D}_G^{-\frac{1}{2}}$ for $\mathbb{I}_N + D_G^{-\frac{1}{2}} A_G D_G^{-\frac{1}{2}}$ where $\widetilde{A}_G = A_G + \mathbb{I}_N$ and $\widetilde{D}_G = \text{diag}\left( \sum_{k=1}^{N} \widetilde{A}[1, k], \cdots, \sum_{k=1}^{N} \widetilde{A}[N, k] \right)$. As a result, the Graph Convolutional Network (GCN) [165, 185] can be defined as

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

The Chebyshev spectral graph filter suffers from a drawback that the spectrum of $\overline{L}_G$ 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 [158] proposes the Cayley spectral graph filter via the order-$r$ Cayley polynomial $f_{c,h}(\lambda) = c_0 + 2\text{Re}\left( \sum_{j=0}^{r} c_h C(\lambda)^j \right)$ with the Cayley transform $C(\lambda) = \frac{\lambda - i}{\lambda + i}$. Moreover, there are many other spectral graph filters, e.g. [6, 18, 48, 72, 121,

153, 160, 166, 187, 239]. In addition, some studies employ the capsule network [167] to construct capsule-inspired GNNs [115, 168, 216].

**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. In general, two approaches, namely sampling [64, 73, 80, 201] and decomposition [200, 215], can be employed to mitigate the time and memory challenges for the spectral GCNNs.

**Depth Trap of Spectral GCNNs.** A bottleneck of GCNNs is that their performance maybe decease with ever-increasing number of layer. This decay is often attributed to three factors: (1) over-fitting resulting from the ever-increasing number of parameters; (2) gradient vanishing/explosion during training; (3) oversmoothing making vertices from different clusters more and more indistin-guishable. 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 [150]. There are some available approaches, e.g. [2, 112, 161, 221, 231], to circumvent the depth trap of the spectral GCNNs.

*3.1.2 Spatial Graph Convolution Operators.* Original spatial GCNNs [52, 53, 55, 194] constitutes a transition function, which must be a contraction map in order to ensure the uniqueness of states, and an update function. In the following, we firstly introduce a generic framework of the spatial GCNN, and then investigate its variants.

Graph networks (GNs) as generic architectures with relational inductive bias [149] 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 $\psi^{e \to v}(\cdot), \psi^{e \to u}(\cdot), \psi^{v \to u}(\cdot)$. The iterations are described as follows.

$$
\begin{aligned}
e_k' &= \phi^e\left(e_k, v_{r_k}, v_{s_k}, u\right), \quad \bar{e}_i' = \psi^{e \to v}\left(E_i'\right), \qquad v_i' = \phi^v\left(v_i, \bar{e}_i', u\right), \\
\bar{e}' &= \psi^{e \to u}\left(E'\right), \qquad\qquad u' = \phi^u\left(u, \bar{e}', \bar{v}'\right), \quad \bar{v}' = \psi^{v \to u}\left(V'\right)
\end{aligned}
\tag{9}
$$

where $e_k$ is an arch from $v_{s_k}$ to $v_{r_k}$, $E_i' = \left\{(e_k', s_k, r_k) : r_k = i, k = 1, \cdots, M\right\}$, $V' = \left\{v_i' : i = 1, \cdots, N\right\}$ and $E' = \left\{(e_k', s_k, r_k) : k = 1, \cdots, M\right\}$, 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, configuable within-block structure and flexible multi-block architectures. Below, we introduce three prevalent variants of the GNs, namely Message Passing Neural Networks (MPNNs) [92], Non-local Neural Networks (NLNNs) [4] and GraphSAGE [204].

**Variants of GNs——MPNNs.** MPNNs [92] 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 $\psi^{e \to v}(\cdot)$ 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,

$$
m_v^{(l+1)} = \sum_{u \in N_G(v)} M_l\left(x_v^{(l)}, x_u^{(l)}, e_{v,u}\right), \quad x_v^{(l+1)} = U_l\left(x_v^{(l)}, m_v^{(l+1)}\right)
$$

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. A lot of GCNNs can be regarded as special forms of the MPNN, e.g. [34, 147, 181, 235].
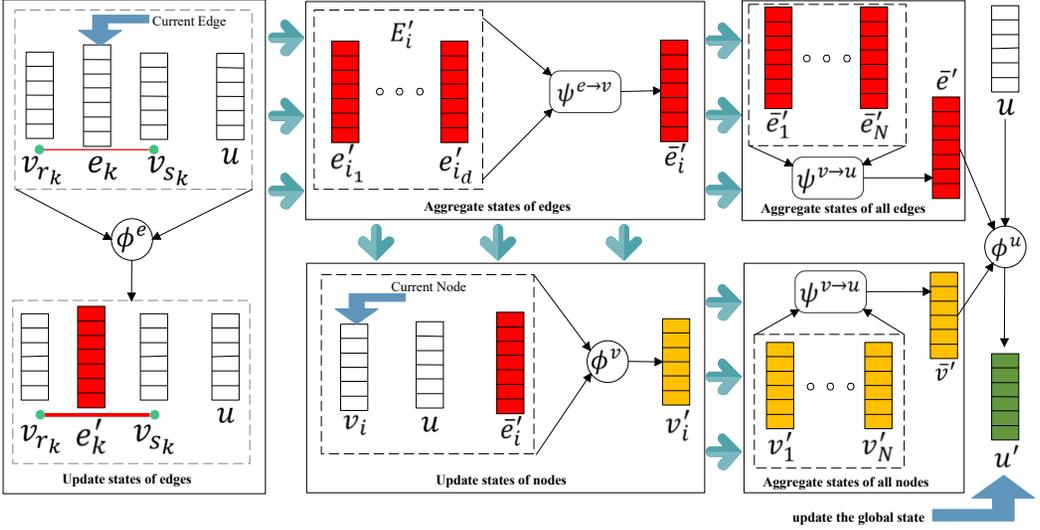
Fig. 3. Computational framework of the spatial GCNN.

**Variants of GNs——NLNNs.** NLNNs [4] give a general definition of non-local operations [8] which is a flexible building block and can be easily integrated into convolutional/recurrent layers. Specifically, the generic non-local operation is defined as

$$y_s = \frac{1}{C(x_s)} \sum_t f(x_s, x_t) g(x_t), \qquad (10)$$

where $f(\cdot, \cdot)$ denotes the affinity between $x_s$ and $x_t$, and $C(x_s) = \sum_t f(x_s, x_t)$ is a normalization factor. The affinity function $f(\cdot, \cdot)$ is of the following form

(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: $f(x_s, x_t) = \text{ReLU}(w_f^T [\theta(x_s), \eta(x_t)])$.

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 Eq. (10) plays the role of $\psi^{e \to v}(E_i')$.

**Variants of GNs——GraphSAGE.** GraphSAGE (SAMple and aggreGatE) [204] 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 AGGREGATE$^{(l)}(\cdot)$ and an update function UPDATE$^{(l)}$, i.e.

$$x_{N_G(v)}^{(l)} = \text{AGGREGATE}^{(l)} \left( \left\{ x_u^{(l-1)} : u \in N_G(v) \right\} \right)$$

$$x_v^{(l)} = \text{UPDATE}^{(l)} \left( \left\{ x_v^{(l-1)}, x_{N_G(v)}^{(l)} \right\} \right)$$

where $N_G(v)$ denotes a fixed-size set of neighbors of $v$ uniformly sampling from its whole neighbors. The aggregation function is of the following form

(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 [171] to aggregate the neighbors of $v$;

(3) Pooling Aggregator: $\text{AGGREGATE}^{(l)} = \max\left(\left\{\sigma\left(Wx_u^{(l)} + b\right) : u \in N_G(v)\right\}\right)$.

Note that the aggregation function and update function play the role of $\psi^{e\to v}(E_i')$ and $\phi^v(v_i, \bar{e}_i', u)$ in formula (9) 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 [16, 155], denoted as $\mathbb{H}_K^n$, is a unique, complete, simply connected $d$-dimensional Riemannian manifold with constant negative sectional curvature $-\frac{1}{K}$, i.e.

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

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 as $\mathcal{T}_x\mathbb{H}_K^n = \left\{v \in \mathbb{R}^{n+1} : \langle x, v\rangle_{\mathcal{M}} = 0\right\}$. 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 $\gamma_{x\to u}(\cdot)$ such that $\gamma_{x\to u}(0) = x$ and $\dot{\gamma}_{x\to u}(0) = u$ is denoted as $\gamma_{x\to u}(t) = \cosh\left(\frac{t}{\sqrt{K}}\right)x + \sqrt{K}\sinh\left(\frac{t}{\sqrt{K}}\right), u, t > 0$. 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 $\left(\mathbb{H}_K^n, d_{\mathcal{M}}^K(\cdot, \cdot)\right)$. In particular, $\mathbb{H}_1^n$, i.e. $K = 1$, is called the hyperboloid model of the hyperbolic space. Hyperbolic Graph Convolutional Networks (HGCN) [68] 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

$$\exp_x^K(v) = \cosh\left(\frac{\|v\|_{\mathcal{M}}}{\sqrt{K}}\right)x + \sqrt{K}\sinh\left(\frac{\|v\|_{\mathcal{M}}}{\sqrt{K}}\right)\frac{v}{\|v\|_{\mathcal{M}}}$$

$$\log_x^K(y) = d_{\mathcal{M}}^K(x, y)\frac{y + \frac{1}{K}\langle x, y\rangle_{\mathcal{M}} x}{\left\|y + \frac{1}{K}\langle x, y\rangle_{\mathcal{M}} x\right\|_{\mathcal{M}}},$$

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 HGCN architecture is composed of three components: a Hyperbolic Feature Transform (HFT), an Attention-Based Aggregation (ABA) and a Non-Linear Activation with Different Curvatures (NLADC). They are respectively defined as

$$h_j^{H,l} = \left(W^{(l)} \otimes^{K_{l-1}} x_j^{H,l-1}\right) \oplus^{K_{l-1}} b^{(l)} \qquad \text{(HFT)},$$

$$y_j^{H,l} = \text{AGGREGATE}^{K_{l-1}}(h^{H,l})_j \qquad \text{(ABA)},$$

$$x_j^{H,l} = \exp_o^{K_l}\left(\rho\left(\log_o^{K_{l-1}}\left(y_j^{H,l}\right)\right)\right) \qquad \text{(NLADC)},$$

where $o = \left(\sqrt{K}, 0, \cdots, 0\right) \in \mathbb{H}_K^n$, the subscript $j$ denotes the indices of nodes, the superscript $l$ denotes the layer of the HGCN. 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 $x^H \oplus^K b = \exp_{x^H}^K \left( P_{o \to x^H}^K (b) \right)$, where $P_{o \to x^H}^K (b)$ 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 $\text{Aggregate}^K (x^H)_j = \exp_{x^H}^K \left( \sum_{k \in N_G(j)} \omega_{j,k} \log_{x^H}^K \left( x_k^H \right) \right)$, where the attention weight $\omega_{j,k} = \text{Softmax}_{k \in N_G(j)} \left( \text{MLP} \left( \log_o^K \left( x_j^H \right) \bowtie \log_o^l \right) \right)$

**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 3 nodes. Here, we introduce the studies on the $k$-dimensional GCNNs [21]. 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 $\{\!\{ \cdots \}\!\}$ denote a multiset, $\text{Hash}(\cdot)$ a hashing function and $C_{l,k}^{(l)}(s)$ the node coloring (label) of $s = (s_1, \cdots, s_k) \in V^k$ at the $l$-th time. Moreover, let $N_G^j(s) = \left\{ (s_1, \cdots, s_{j-1}, r, s_{j+1}, \cdots, s_k) : r \in V \right\}$. The $k$-WL is computed by

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

where $c_j^{(l+1)} = \text{Hash} \left( \left\{\!\!\left\{ C_{l,k}^{(l)}(s') : s' \in N_G^j(s) \right\}\!\!\right\} \right)$. 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 often employed to learn the hierarchical representations of nodes in order to scale to larger graphs and mitigate the overfitting problem.

**Other Variants of GNs.** In addition to the aforementioned GNs and its variants, there are still many other spatial GCNNs which is defined from other perspectives, e.g. Diffusion-Convolutional Neural Network (DCNN) [69], Position-aware Graph Neural Network (P-GNN) [78], Memory-based Graph Neural Network (MemGNN) and Graph Memory Network (GMN) [5], Graph Partition Neural Network (GPNN) [152], Edge-Conditioned Convolution (ECC) [118], DEMO-Net [103], Column network [188], Graph-CNN [47].

**Invariance and Equivariance.** Permutation-invariance refers to that a function $f : \mathbb{R}^{n^k} \to \mathbb{R}$ (e.g. the aggregation function) is independent of any permutations of node/edge indices [57, 136], i.e. $f(P^T A_G P) = f(A_G)$ 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 $f : \mathbb{R}^{n^k} \to \mathbb{R}^{n^l}$ coincides with permutations of node/edge indices [57, 136], i.e. $f(P^T A_G P) = P^T f(A_G) P$ where $P$ and $A_G$ are defined as similarly as the permutation-invariance. For permutation-invariant aggregation functions, a straightforward choice is to take sum/max/average/concatenation as heuristic aggregation schemes [136]. 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 [248] proposes a scalable approach for learning adaptive receptive fields of GCNNs. It 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 $\alpha(x, y) = \text{Softmax}_y\left(\alpha^T \tanh\left(W_x^T x + W_y^T y\right)\right)$. The adaptive depth function is defined as a LSTM [171], i.e.

$$
\begin{aligned}
i_{v_j} &= \sigma\left(\left(W_i^{(t)}\right)^T h_{v_j}^{\text{temp}}\right) & f_{v_j} &= \sigma\left(\left(W_f^{(t)}\right)^T h_{v_j}^{\text{temp}}\right) \\
o_{v_j} &= \sigma\left(\left(W_o^{(t)}\right)^T h_{v_j}^{temp}\right) & \widetilde{C}_{v_j} &= \tanh\left(\left(W_c^{(t)}\right)^T h_{v_j}^{(temp)}\right) \\
C_{v_j}^{(t+1)} &= f_{v_j} \otimes C_{v_j}^{(t)} + i_{v_j} \otimes \widetilde{C}_{v_j} & h_{v_j}^{(t+1)} &= o_{v_j} \otimes \tanh\left(C_{v_j}^{(t+1)}\right).
\end{aligned}
\tag{11}
$$

GEOM-GCN [61] proposes a novel permutation-invariant geometric aggregation scheme consisting of three modules, namely node embedding, structural neighborhood, and bi-level aggregation. This aggregation scheme does not lose structural information of nodes and fail to capture long-range dependencies in disassortative graphs. For the permutation-invariant graph representations, PiNet [148] proposes an end-to-end spatial GCNN architecture that utilizes the permutation equivariance of graph convolutions. It is composed of a pair of double-stacked message passing layers, namely attention-oriented message passing layers and feature-oriented message passing layers.

**Depth Trap of Spatial GCNNs.** Similar to the spectral GCNNs, the spatial GCNNs is also confronted with the depth trap. As stated previously, the depth trap results from oversmoothing, overfitting and gradient vanishing/explosion. In order to escape from the depth trap, some studies propose some available strategies, e.g. DeepGCN [56] and Jumping Knowledge Network [99]. The jumping knowledge networks [99] adopt neighborhood aggregation with skip connections to integrate information from different layers. The DeepGCN [56] apply the residual/dense connections [54, 94] and dilated aggreagation [50] in the CNNs to construct the spatial GCNN architecture. They has three instantiations, namely ResGCN, DenseGCN and dilated graph convolution. ResGCN is inspired by the ResNet [94], 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 $\mathcal{F}(\cdot, \cdot)$ can be computed by spectral or spatial GCNNs. DenseGCN collectively exploit information from different GCNN layers like the DenseNet [54], which is defined to be

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

The dilated aggregation [50] 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 $(u_1, u_2, \cdots, u_{k \times d})$ are the first sorted $k \times d$ nearest neighbors, then $N_G^{(k,d)}(v) = \left\{u_1, u_{1+d}, \cdots, u_{1+(k-1)d}\right\}$. Thereby, we can construct a new graph $G^{(k,d)} = \left(V^{(k,d)}, E^{(k,d)}\right)$ where $V^{(k,d)} = V$ and $E^{(k,d)} = \left\{\langle v, u \rangle : v \in V, u \in N_G^{(k,d)}\right\}$. The dilated graph convolution layer can be obtained by running the spatial GCNNs over $G^{(k,d)}$.

*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) [36]. First of all, the SGWT is determined by a graph wavelet generating kernel $g : \mathbb{R}^+ \to \mathbb{R}^+$ with the property $g(0) = 0, g(+\infty) = \lim_{x \to \infty} g(x) = 0$. A 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.

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

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 $\Psi_g^s = U g(s\Lambda) U^T$ where $g(s\Lambda) = g\left(\text{diag}\left(s\lambda_1, \cdots, s\lambda_N\right)\right)$. A graph signal $x \in \mathbb{R}^N$ on $G$ can thereby be filtered by the spectral graph wavelet operator, i.e. $\mathcal{W}_{g,s}^x = \Psi_g^s x \in \mathbb{R}^N$. The literature [39] utilizes a special instance of the graph wavelet operator $\Psi_g^s$ to construct a graph scattering network, and proves its covariance and approximate invariance to permutations and stability to graph operations.

**Graph Wavelet Neural Networks.** The above spectral graph wavelet operator $\Psi_g^s$ can be employed to construct a Graph Wavelet Neural Network (GWNN) [14]. Let $\Psi_g^{-s} \triangleq \left(\Psi_g^s\right)^{-1}$. The graph wavelet based convolution is defined to be

$$x *_G y = \Psi_g^{-s} \left(\Psi_g^s x \circledast \Psi_g^s y\right).$$

The GWNN is composed of multiple layers of the graph wavelet based convolution. The structure of the $l$-th layer is defined as

$$X^{(l+1)}[:, j] = \rho \left( \sum_{k=1}^{d^{(l)}} \Psi_g^{-s} \Theta_{j,k}^{(l)} \Psi_g^s X^{(l)}[:, k] \right), \tag{12}$$

where $\Theta_{j,k}^{(l)}$ is a diagonal filter matrix learned in spectral domain. Eq. (12) can be rewritten as a matrix form, i.e. $X^{(l+1)} = \rho \left(\Psi_g^{-s} \Theta \Psi_g^s X^{(l)} W^{(l)}\right)$. The learnable filter matrix $\Theta$ can be replaced with the $K$-localized Chebyshev Polynomial so as to eschew the time-consuming eigendecomposition of $\overline{L}_G$.

*3.1.4   Summary.* The aforementioned GCNN architectures provide available ingredients of constructing the GNNs. In practice, we can construct our own GCNNs by assembling different modules introduced above. Additionally, some scholars also study the GCNNs from some novel perspectives, e.g. the parallel computing framework of the GCNNs [111], the hierarchical covariant compositional networks [157], the transfer active learning for GCNNs [174] and quantum walk based subgraph convolutional neural network [242]. They 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 the 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.

*3.2.1 Global Graph Pooling Operators.* Global graph pooling operators pool all of representations of nodes into a universal graph representation. Many literatures [90, 148, 160] 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) [162] provides a novel framework for graph representation with maximal representation power. Specifically, 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 undirected graphs is defined as

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

where $\Pi_{|V|}$ is the set of all distinct permutations on $V_G$ and $\overrightarrow{f}(\cdot, \cdot)$ is an arbitrary (possibly permutation-sensitive) vector-valued function. Specifically, $\overrightarrow{f}(\cdot, \cdot)$ can be denoted as Multi-Layer Perceptrons (MLPs), Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs) or Graph Neural Networks (GNNs). The literature [162] proves that $\bar{\bar{f}}(G)$ has the most expressive representation of $G$ under some mild conditions, and provides approximation approaches to making RP computationally tractable. In addition, there are some other available global graph pooling operators, e.g. SortPooling [130] and function space pooling [141].

*3.2.2 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 vanilla 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 [85] and Graclus [67]), learning a soft cluster assignment and selecting the first $k$ top-rank nodes.

**Invoking existing graph clustering algorithms.** The graph clustering aims to assign proximal nodes to the same cluster and in-proximal nodes to different clusters. The coarsened graph regards the resulting clusters as super-nodes and connections between two clusters as 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 [182] to compute the representations of the super-nodes. The literature [219] proposes the EigenPooling method and presents the relationship between the original and coarsened graph. In order to construct a coarsened graph of $G$, a graph clustering method is employed to partition $G$ into $K$ disjoint clusters, namely $\{G_k : k = 1, \cdots, K\}$. Suppose each cluster $G_k$ has $N_k$ nodes, namely $\{v_{k,1}, \cdots, v_{k,N_k}\}$, 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 $G_k, k = 1, \cdots, K$ 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 $A_{\text{coar}} = f\left(M_{\text{coar}}^T A_G M_{\text{coar}}\right)$ as well, where $f(\widetilde{a}_{i,j}) = 1$ if $\widetilde{a}_{i,j} > 0$ and $f(\widetilde{a}_{i,j}) = 0$ 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 [219] employs spectral clustering to obtain the coarsened graph, and then up-sample the Fourier basis of subgraphs $G_k, k = 1, \cdots, K$. 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 [49] proposes a novel Hierarchical Graph Convolutional Network (H-GCN) consisting of graph coarsening layers and graph refining layers. The former employs structural equivalence grouping and structural similarity grouping to construct the coarsened graph, and the latter restores the original topological structure of the corresponding graph.

**Learning a soft cluster assignment.** STRUCTPOOL [59], as a structured graph pooling technique, regards the graph pooling as a graph clustering problem so as to learn a cluster assignment matrix via the feature matrix $X$ and adjacency matrix $A_G$. Learning the cluster assignment matrix can formulated as a Conditional Random Field (CRF) [89] based probabilistic inference. Specifically, the input feature matrix $X$ is treated as global observation, and $Y = \{Y_1, \cdots, Y_N\}$ is a random field where $Y_i \in \{1, \cdots, K\}$ 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.

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

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 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 $\Psi_p = \left(\psi_p(y_s, y_t|X)\right)_{K \times K}$ 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 $A_{\text{coar}} = f\left(M^T A_G M\right)$ and $X_{\text{coar}} = M^T X$. DiffPool [154] 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. $\left(A_{G^{(l+1)}}, X^{(l+1)}\right) = \text{DiffPool}\left(A_{G^{(l)}}, Z^{(l)}\right)$. More specifically, $X^{(l+1)} = \left(M^{(l)}\right)^T Z^{(l)}$, $A_{G^{(l+1)}} = \left(M^{(l)}\right)^T A_{G^{(l)}} M^{(l)}$. 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. $Z^{(l)} = \text{GCNN}_{\text{embed}}\left(A_{G^{(l)}}, X^{(l)}\right)$, $M^{(l)} = \text{Softmax}\left(\text{GCNN}_{\text{pool}}\left(A_{G^{(l)}}, X^{(l)}\right)\right)$.

**Selecting the first $k$ top-rank nodes.** The literature [91] proposes a novel Self-Attention Graph Pooling operator (abbreviated as SAGPool). Specifically, SAGPool firstly employs the GCN [185] to calculate the self-attention scores, and then invokes the top-rank function to select the top $\lceil kN \rceil$ node indices, i.e.

$$Z = \rho\left(\widetilde{D}_G^{-\frac{1}{2}} \widetilde{A}_G \widetilde{D}_G^{-\frac{1}{2}} X\Theta\right), \quad \text{idx} = \text{top-rank}\left(Z, \lceil kN \rceil\right), \quad Z_{\text{mask}} = Z_{\text{idx}}$$
$$X' = X_{\text{idx}}, \qquad X_{\text{out}} = X' \circledast Z_{\text{mask}}, \qquad A_{\text{out}} = A_{\text{idx,idx}}.$$

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 [138], the literature [63] 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.

$$y = \frac{Xp}{\|p\|}, \qquad \text{idx} = \text{top-rank}(y, k), \quad \widetilde{y} = \tanh(y_{\text{idx}})$$
$$\widetilde{X}_{\text{coar}} = X_{\text{idx},:}, \quad A_{\text{coar}} = A_{\text{idx,idx}}, \qquad X_{\text{coar}} = \widetilde{X}_{\text{coar}} \circledast \left(\widetilde{y} \mathbf{1}_C^T\right),$$

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, let $X_{\text{refine}} = \text{Distribute}(0_{N \times d}, X_{\text{coar}}, \text{idx})$, 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 idx. Note that row vectors of $X_{\text{refine}}$ with indices in $idx$ are updated by the ones in $X_{\text{coar}}$, whereas other row vectors remain zero. It is worth noting that the literature [28] adopts the similar pooling strategy as gPool to learn the hierarchical representations of 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 [11, 40, 195]. Recently, applying the attention mechanisms to GCNNs has gained considerable attentions so that various attention techniques have been proposed. Below, we summarize the graph attention mechanisms on graphs from the next 4 perspectives [87], 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 $v_1, \cdots, v_{d_0}$, and their current feature vectors are respectively denoted as $x_0, x_1, \cdots, x_{d_0}$, where $d_0 = d_G(v_0)$.
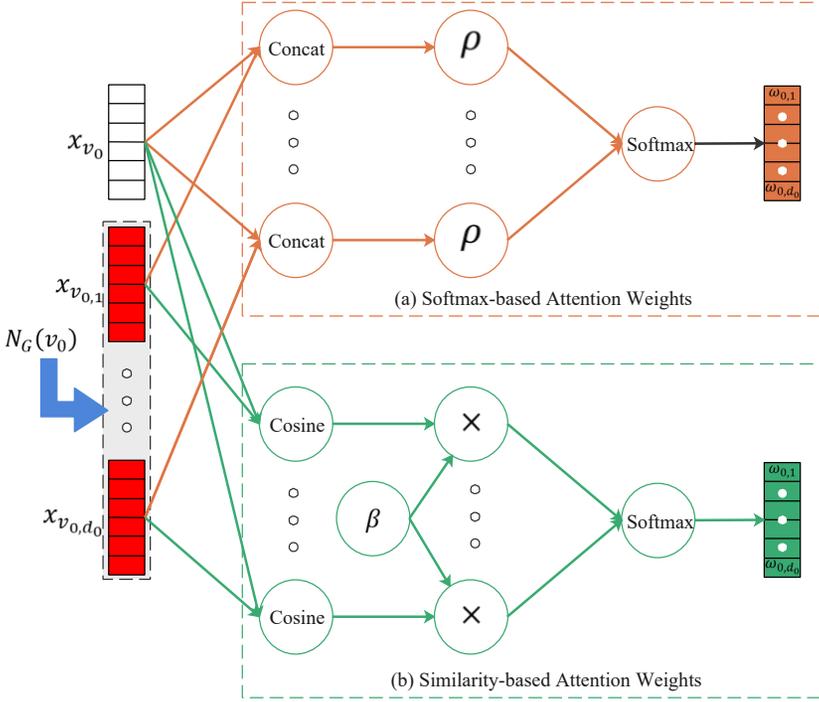


Fig. 4. Two kinds of graph attention mechanisms.

*3.3.1 Concatenation-based Graph Attention.* The softmax-based graph attention is typically implemented by employing a softmax with learnable weights [145, 233] to measure the relevance of $v_j, j = 1, \cdots, d_G(v_0)$ to $v_0$. More specifically, the softmax-based attention weights between $v_0$ and $v_j$ can be defined as

$$
\begin{aligned}
\omega_{0,j} &= \text{Softmax}\left(\left[e_{0,1}, \cdots, e_{0,d_G(v_0)}\right]\right) \\
&= \frac{\exp\left(\rho\left(a^T\left(Wx_0 \bowtie Wx_j\right)\right)\right)}{\sum_{k=1}^{d_G(v_0)} \exp\left(\rho\left(a^T\left(Wx_0 \bowtie W_k\right)\right)\right)},
\end{aligned}
\tag{13}
$$

where $e_{0,j} = \exp\left(\rho\left(a^T\left(Wx_0 \bowtie Wx_j\right)\right)\right)$, $a$ 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

$$
x_0' = \rho\left(\sum_{j=1}^{d_G(v_0)} \omega_{0,j} Wx_j\right).
\tag{14}
$$

In practice, multi-head attention mechanisms are usually employed to stabilize the learning process of the single-head attention [145]. For the multi-head attention, assume that the feature vector of each head is $x_0^{(h)} = \rho\left(\sum_{j=1}^{d_G(v_0)} \omega_{0,j}^{(h)} W^{(h)} x_j\right)$. The concatenation based multi-head attention is

computed by $x'_0 = \bowtie_{h=1}^{H} x_0^{(h)} = \bowtie_{h=1}^{H} \rho \left( \sum_{j=1}^{d_G(v_0)} \omega_{0,j}^{(h)} W^{(h)} x_j \right)$. 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 to assign different weights to heads, and gets the formulation of the gated multi-head attention mechanism. The Graph Transformer (GTR) [233] 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 [75] 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. $X^{(l+1)} = \text{Softmax} \left( \frac{QK^T}{\sqrt{d_h}} \right) V$, 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 [105] 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.

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

$$\omega_{0,j} = \frac{\exp \left( \beta \cdot \cos \left( W x_0, W x_j \right) \right)}{\sum_{k=1}^{d_G(v_0)} \exp \left( \beta \cdot \cos \left( W x_0, W x_k \right) \right)}, \tag{15}$$

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) [104] 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 (15).

### 3.3.3 Spectral Graph Attention.
The Spectral Graph Attention Networks (SpGAT) aims to learn representations for different frequency components [60]. The eigenvalues of the normalized graph Laplacian $\overline{L}_G$ can be treated as frequencies on the graph $G$. As stated in the Preliminary section, $0 = \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_N = \lambda_{\max}$. The SpGAT firstly extracts the low-frequency component $B_L = \{u_1, \cdots, u_n\}$ and the high-frequency component $B_H = \{u_{N-n+1}, \cdots, u_N\}$ from the graph Fourier bases $\{u_1, u_2, \cdots, u_N\}$. So, we have

$$\begin{aligned} X_L &= X^{(l)} \Theta_L, \quad X_H = X^{(l)} \Theta_H \\ X^{(l+1)} &= \rho \left( \text{Aggregate} \left( B_L F_L B_L^T X_L, B_H F_H B_H^T X_H \right) \right), \end{aligned} \tag{16}$$

where $F_L$ and $F_H$ respectively 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 $\Omega_L = \text{diag} \left( \omega_L, \cdots, \omega_L \right)$ and $\Omega_H = \text{diag} \left( \omega_H, \cdots, \omega_H \right)$ 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. $\omega_L =$

$\frac{\exp(\omega_L)}{\exp(\omega_L)+\exp(\omega_H)}$, $\omega_H = \frac{\exp(\omega_H)}{\exp(\omega_L)+\exp(\omega_H)}$. 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 [14, 26] 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. (16) for the spectral graph wavelet operator $\Psi^s_{L,g}$ and $\Psi^s_{H,g}$, i.e.

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

*3.3.4 Attention-guided Walk.* The two aforementioned kinds of attention mechanisms focus 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 [88], 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 $v_t$. The hidden representation $x^{(t)}$ of $v_t$ is computed by a recurrent neural network $f_x(\cdot)$ taking the step embedding $s^{(t)}$ and internal representation of the historical information from the previous step $x^{(t-1)}$ as input, i.e.

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

The step embedding $s^{(t)}$ is computed by a step network $f_s \left( r^{(t-1)}, x_{v_t}; \Theta_s \right)$ 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.

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

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

$$ r^{(t)} = f_r \left( x^{(t)}; \Theta_r \right), \quad \hat{l}^{(t)} = f_p \left( x^{(t)}; \Theta_p \right). $$

The ranking network $f_r \left( x^{(t)}; \Theta_r \right)$ determines which neighbors of $v_t$ should be prioritized in the next step, and the predicting network $f_p \left( x^{(t)}; \Theta_p \right)$ 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.

## 3.4 Graph Recurrent Neural Networks

The Graph Recurrent Neural Networks (GRNNs) generalize the Recurrent Neural Networks (RNNs) to process the graph-structured data. In general, the GRNN can be formulated as $h^{(l+1)}_{v_j} = \text{GRNN} \left( x^{(l)}_{v_j}, \left\{ h^{(l)}_{v_k} : v_k \in N_G(v_j) \cup \{v_j\} \right\} \right)$. Below, we introduce some available GRNN architectures.

*3.4.1 Graph LSTM.* The Graph Long Short Term Memroy (Graph LSTM) [93, 133, 190, 206, 213, 230, 234] generalizes the vanilla LSTM for the sequential data to the ones for general graph-structured data. Specifically, the graph LSTM updates the hidden states and cell states of nodes by the following
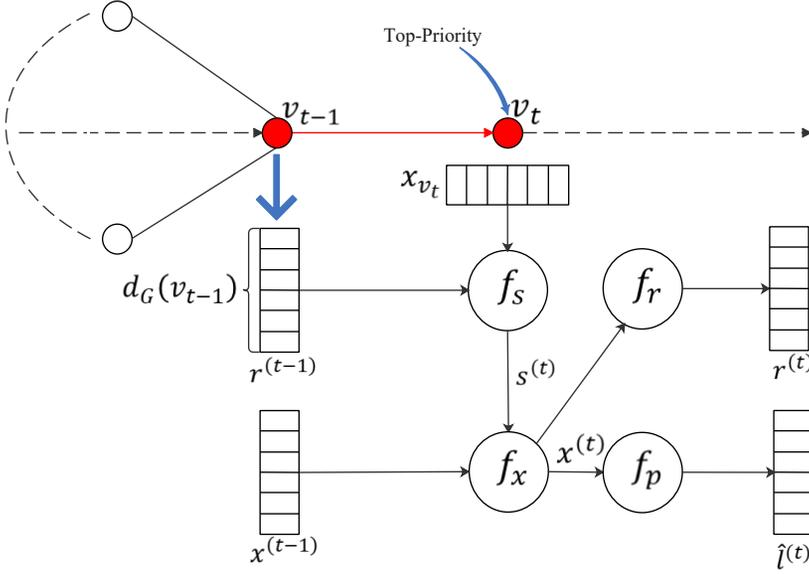
Fig. 5. The computational framework of the attention-guided walk.

formula,

$$i_{v_j}^{(l+1)} = \sigma \left( W_i x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_i h_{v_k}^{(l)} + b_i \right),$$

$$o_{v_j}^{(l+1)} = \sigma \left( W_o x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_o h_{v_k}^{(l)} + b_o \right)$$

$$\widetilde{c}_{v_j}^{(l+1)} = \tanh \left( W_c x_{v_j}^{(l)} + \sum_{v_k \in N_G(v_j) \cup \{v_j\}} U_c h_{v_k}^{(l)} + b_c \right),$$

$$f_{v_j,v_k}^{(l+1)} = \sigma \left( W_f x_{v_j}^{(l)} + U_f h_{v_k}^{(l)} + b_f \right)$$

$$c_{v_j}^{(l+1)} = i_{v_j}^{(l+1)} \circledast \widetilde{c}_{v_j}^{(l+1)} + \sum_{v_k \in N_G v_j \cup \{v_j\}} f_{v_j,v_k}^{(l+1)} \circledast c_{v_k}^{(l)},$$

$$h_{v_j}^{(l+1)} = o_{v_j}^{(l+1)} \circledast \tanh(c_{v_j}^{(l+1)}).$$

see the Fig. 6. The literature [212] 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) [235] employs the Gated Recurrent Unit (GRU) [106] to modify the vanilla GCNNs so that it can be extended to process the sequential data.

*3.4.2 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 features on nodes/edges evolve over time. The GRNNs are a straightforward approach to tackling dynamic graphs. Below, we introduce some studies on the GRNNs for dynamic graphs.

The literature [220] proposes a Dynamic Graph Neural Network (DGNN) concerning the dynamic graph only with 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
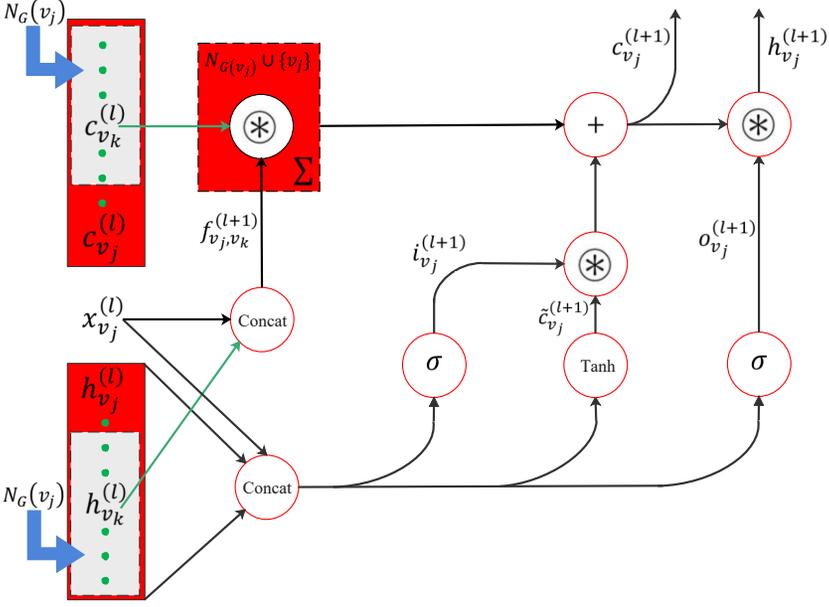
Fig. 6. Computational framework of the Graph LSTM.

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 [171] to respectively update the cell states and hidden states of the source and target, i.e.

$$\left(C_{v_s}^{(t)}, h_{v_s}^{(t)}\right) = \text{LSTM}_s\left(C_{v_s}^{(t-)}, h_{v_s}^{(t-)}, \Delta t_s\right), \quad \left(C_{v_g}^{(t)}, h_{v_g}^{(t)}\right) = \text{LSTM}_g\left(C_{v_g}^{(t-)}, h_{v_g}^{(t-)}, \Delta t_g\right).$$

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 [51] addresses the vertex- and graph-focused prediction tasks on dynamic graphs with a fixed node set by combining GCNs, LSTMs and fully connected layers. The Variational Graph Recurrent Neural Networks (VGRNNs) [42] 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. The semi-implicit variational inference is employed to approximate the posterior so as to generate the node embedding.

*3.4.3 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 a RNN, e.g. LSTM and GRU, to capture long short-term dependencies. The literature [208] performs joint random walks on attributed networks, and utilizes them to boost the deep node representation learning. The proposed GraphRNA in [208] consists of two key components, namely a collaborative walking mechanism

AttriWalk and a tailored deep embedding architecture for joint random walks GRN. Suppose $\mathcal{A}_G$ denotes the node-attribute matrix of size $N \times M$. The AtriWalk admits the transition matrix of size $\mathbb{R}_+^{(N+M) \times (N+M)}$ which is written as

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

After obtaining the sequences via the collaborative random walk, the bi-directional GRU [106] and pooling operator are employed to learn the global representations of sequences. The literature [41] 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 essentially provide ingredients of constructing the GNNs for us. Below, we investigate the extensions of the GNNs from the next 8 aspects: GCNNs on spectral graphs, 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, and briefly summarize the applications of the GNNs at last.

### 4.1 GCNNs on Special Graphs

The vanilla GCNNs aims at learning the representations of input graphs (directed or undirected, weighted or unweighted). The real-world graphs may have more additional characteristics, e.g. spatial-temporal graphs, heterogeneous graphs, hyper-graphs, signed graphs and so all. The GCNN for signed graphs [189] leverage the balance theory to aggregate and propagate information through positive and negative links.

*4.1.1 Heterogeneous Graphs.* Heterogeneous Graphs are composed of nodes and edges of different types, and each type of edges is called a relation between two types of nodes. 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 [232]. The heterogeneity and rich semantic information brings great challenges for designing heterogeneous graph convolutional neural networks. In general, a heterogeneous graph can be denoted as $H = (V, E, v, \zeta)$, where $v(v)$ 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 vanilla heterogeneous GCNNs, namely HetGNNs [25].

**Vanilla Heterogeneous GCNNs.** The Heterogeneous Graph Neural Networks (HetGNNs) [25] aims to resolve the issue of jointly considering heterogeneous structural information as well as heterogeneous content 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 $C_v$ denote the heterogeneous contents of

node $v$, which can be encoded as a fixed size embedding via a function $f_1(v)$, i.e.

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

where $\mathcal{F}C_{\theta_x}(\cdot)$ denotes feature transformer, e.g. identity or fully connected neural networks with parameter $\theta_x$, and $\overrightarrow{\text{LSTM}}(\cdot)$ and $\overleftarrow{\text{LSTM}}(.)$ is defined by the Eq. (11). The content embedding of the $t$-type neighbors of node $v$ can be aggregated as follows,

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

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 $\omega_{v,j} = \frac{\exp \left( \rho \left( u^t \left[ f_j(v) \bowtie f_1(v) \right] \right) \right)}{\sum_{f_j(v) \in \mathcal{F}(v)} \exp \left( \rho \left( u^t \left[ f_j(v) \bowtie f_1(v) \right] \right) \right)}$. In addition, the GraphInception [20] can be employed to learn the hierarchical relational features on heterogeneous graphs by converting the input graph into a multi-channel graph (each meta path as a channel) [227].

**Heterogeneous Graph Attention Mechanism.** The literature [210] 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 aims 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

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

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

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

where $\omega^{\Phi_p} = \frac{1}{|V|} \sum_{v_k \in V} q^T \tanh \left( W x_k^{\Phi_p} + b \right)$. As a result, the embedding matrix $X = \sum_{p=1}^P \beta^{\Phi_p} X^{\Phi_p}$. In addition, there are some available studies on the GCNNs for multi-relational graphs [17, 30, 218] and the transformer for dynamic heterogeneous graphs [170, 247].

*4.1.2 Spatio-Temporal Graphs.* Spatio-temporal graphs can be used to model traffic networks [13, 175] and skeleton networks [178, 237]. In general, a spatio-temporal graph is denoted as $G_{ST} = (V_{ST}, E_{ST})$ where $V_{ST} = \{v_{t,j} : t = 1, \cdots, T, j = 1, \cdots, N_{ST}\}$. 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 [13] proposes a novel deep learning framework, namely 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 $C_{\text{in}}$ channels, i.e. $X \in \mathbb{R}^{M \times C_{\text{in}}}$. The temporal gated convolution kernel $\Gamma \in \mathbb{R}^{K \times C_{\text{in}} \times 2C_{\text{out}}}$ is used to filter the input $Y$, i.e.

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

to yield an output $P \bowtie Q \in \mathbb{R}^{(M-K+1) \times (2C_{\text{out}})}$. 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. In addition, some other studies pay attention to the GCNNs on the spatio-temporal graphs from other perspectives, e.g. Structural-RNN [10] via a factor graph representation of the spatio-temporal graph and GCRNN [114, 229] combining the vanilla GCNN and RNN.

*4.1.3 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 convectional 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 $\mathcal{V} = \{v_1, \cdots, v_N\}$ like the conventional graph, $\mathcal{E} = \{e_1, \cdots, e_M\}$ 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 $\mathcal{D}_{\mathcal{V}} = \text{diag} \left( \deg_{\mathcal{V}}(v_1), \cdots, \deg_{\mathcal{V}}(v_N) \right), \mathcal{D}_{\mathcal{E}} = \text{diag} \left( \deg_{\mathcal{E}}(e_1), \cdots, \deg_{\mathcal{E}}(e_M) \right)$ and $\mathcal{W}_{\mathcal{G}} = \text{diag} \left( \omega(e_1), \cdots, \omega(e_M) \right)$. The hypergraph Laplacian [224] $\mathcal{L}_{\mathcal{G}}$ of $\mathcal{G}$ is defined to be

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

It can also be factorized by the eigendecomposition, i.e. $\mathcal{L}_{\mathcal{G}} = \mathcal{U} \Lambda \mathcal{U}^T$. The spectral hypergraph convolution operator, the Chebyshev hypergraph convolutional neural network and the hypergraph convolutional network can be defined in analogy to the Eqs (4,6,8). The HyperGraph Neural Network

(HGNN) architecture proposed in the literature [224] is composed of multiple layers of the hyperedge convolution, and each layer is defined as

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

In essence, the HGNN essentially views each hyperedge as a complete graph so that the hypergraph is converted into a conventional graph. Treating each hyperedge as a complete graph obviously incurs expensive computational cost. Hence, some studies [131, 132] propose various approaches to approximate the hyperedges. The HNHN [225] interleaves updating the node representations with the hyperedge representations by the following formulas,

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

## 4.2 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 their eyes on the capability and interpretability of the GCNNs.

*4.2.1 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 [21, 101, 163], 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

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

According to the literatures [21, 101], 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 [101] 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 [7] 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 [140], $p$-order graph moments [137], algorithmic alignment with the dynamic programming [100], generalization and representational limits of the GNNs [192].

*4.2.2    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 their eyes on the interpretability of the GNs [149]. The literature [44] 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 [246] 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 an optimization problem

$$\max_{G_S, X_S^F} I\left(Y, (G_S, X_S^F)\right) = H(Y) - H\left(Y | G = G_S, X = X_S^F\right), \tag{17}$$

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 $\left\{X^F[j,:] : v_j \in G_S\right\}$. The entropy term $H(Y)$ is constant because the spatial GCNN model $\Phi$ is fixed. In order to improve the tractability and computational efficiency of the GNNExplainer, the final optimization framework is reformulated as

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

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 $c$?".

## 4.3    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 [15]. Subsequently, many studies exploit deep learning techniques to learn low-dimensional representations of nodes [203]. 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 [29, 38, 97, 142, 173, 209] or RNNs to capture long and short-term dependencies of node sequences yielded by random walks [98, 202]. In the following, we introduce the network embedding approaches based on GNNs.

*4.3.1    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 [23, 169, 176, 184, 241, 243]. A representative approach among these ones is the Variational Graph Auto-Encoder (VGAE) [184] consisting of a GCN based

encoder and an inner product decoder. The GCN based encoder is defined to be

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

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

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

They adopt the evidence lower bound [37] as their objective function. The adversarially regularized (variational) graph autoencoder [176] extends the VGAE by adding an adversarial regularization term to the evidence lower bound. The literature [84] proposes a symmetric graph convolutional autoencoder which produces a low-dimensional latent nodes representations. Its encoder employs the Laplacian smoothing [150] 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 $X^{(l+1)} = (1 + \gamma)X^{(l)} - \gamma D^{-1}AX^{(l)} = X^{(l)} + \gamma(\mathbb{I}_N - D^{-1}A)X^{(l)}$, 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 [240], pre-training GNNs via context prediction and attribute masking strategies [198], and deep graph Infomax using a noise-contrastive type objective with a standard binary cross-entropy loss between positive examples and negative examples [146].

## 4.4 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 [214] 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 [15] utilizes the generative adversarial network [66] 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 [1], from faked ones. The generator $G$ aims to generate faked node sequences via another LSTM, whose generating process is as follows.

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

The motif-targeted GAN [9] generalizes random walk based architecture of the NetGAN to characterize mesoscopic context of nodes. Different from [9, 15], the GraphVAE [119] 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 [236] defines a sequential decision-making process to add a node/edge via the graph network [149], readout operator and softmax function. The GraphRNN [77] 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.

## 4.5 Combinations of the PI and GNNs

The GNNs and PI are two different learning paradigms for 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. This subsection provides a summarization of studies of combining these two paradigms.

*4.5.1 Conditional Random Field Layer Preserving Similarities between Nodes.* The literature [62] proposes a CRF layer for the GCNNs to enforce hidden layers to preserve similarities between nodes. Specifically, the input $X^{(l)}$ to $(l+1)$-th layer is a random vector around the output $B^{(l)} = $ GCNN$(X^{(l-1)}, A_G, X)$ of the $(l-1)$-th layer. The objective function for the GCNN with a CRF layer can be reformulated as

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

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. The similarity constraint $\mathcal{R}(X^{(l)})$ is modeled as a CRF, i.e. $p\left(X^{(l)}|B^{(l)}\right) = \frac{1}{Z(B^{(l)})} \exp\left(-E\left(X^{(l)}|B^{(l)}\right)\right)$ where the energy function $E\left(X^{(l)}|B^{(l)}\right)$ is formulated as

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

Let $s_{u,v}$ denote the similarity between $u$ and $v$. The unary energy component $\varphi_v(\cdot, \cdot)$ and pairwise energy component $\varphi_{u,v}(\cdot, \cdot, \cdot, \cdot)$ for implementing the similarity constraint are respectively formulated as

$$\varphi_v\left(X^{(l)}[v, :], B^{(l)}[v, :]\right) = \|X^{(l)}[v, :] - B^{(l)}[v, :]\|_2^2,$$
$$\varphi_{u,v}\left(X^{(l)}[u, :], X^{(l)}[v, :], B^{(l)}[u, :], B^{(l)}[v, :]\right) = s_{u,v}\|X^{(l)}[u, :] - X^{(l)}[v, :]\|_2^2.$$

The mean-field variational Bayesian inference is employed to approximate the posterior $p(B^{(l)}|X^{(l)})$. Consequently, the CRF layer is defined as

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

*4.5.2 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 [122] 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. $p(Y_V|X_V) = \frac{1}{Z(X_V)} \prod_{(i,j) \in E} \varphi_{i,j}\left(y_i, y_j, X_V\right)$, 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

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

This objective can be optimized by the variational Bayesian EM algorithm [151], 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 approximate

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

and $p_{\theta_l}\left(y_n|Y_{N_G(v)}, X_V\right)$ is approximated by another GCNN. The literature [183] 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.

*4.5.3  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 $f : \mathbb{R}^d \to \mathbb{R}$ follows a Gaussian Process GP$(m(\cdot), \kappa(\cdot, \cdot))$ iff $(f(X_1), \cdots, f(X_N))^T$ for any $N$ $d$-dimensional random vectors. follows a $N$-dimensional Gaussian distribution $\mathcal{N}_N(\mu, \Sigma)$, where $\mu = (m(X_1), \cdots, m(X_N))^T$ and $\Sigma = \left[\kappa(X_j, X_k)\right]_{N \times N}$, 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 $\mathcal{D} = \left\{(X_j, y_j) : j = 1, \cdots N\right\}$, 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 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 [110] employs amortized variational Bayesian inference to approximate $p(f|\mathcal{D})$, i.e. $f = \mu + L\epsilon$, and the GCNNs to estimate $\mu$ and $L$.

*4.5.4  Other GCNN-based Probabilistic Inference.* The literature [113] combines the GCNNs and variational Bayesian inference to infer the input graph structure. The literature [102] infers marginal probabilities in probabilistic graphical models by incorporating the GCNNs to the conventional message-passing inference algorithm. The literature [238] approximates the posterior in Markov logic networks with the GCNN-enhanced variational Bayesian inference.

## 4.6  Adversarial Attacks for the GNNs

In many circumstances where classifiers are deployed, adversaries deliberately contaminate data in order to fake the classifiers [66, 197]. 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 tasks. Therefore, it is inevitable to study the adversarial attacks for GNNs and defense [108].

*4.6.1  Adversarial Attacks on Graphs.* The literature [58] firstly proposes a reinforcement learning based attack method, which can learn a generalizable attack policy, on graphs. This paper provides a definition for a graph adversarial attacker. Given a sample $(G, c, y) \in \{(G_j, c_j, y_j) : j = 1, \cdots, M\}$ and a classifier $f \in \mathcal{F}$, the graph adversarial attacker $g : \mathcal{F} \times \mathcal{G} \to \mathcal{G}$ attempts to modify a graph $G = (V, E)$ into $\widetilde{G} = (\widetilde{V}, \widetilde{E})$, such that

$$\max_g \quad \mathbb{I}(f(\widetilde{G}, c) \neq y)$$
$$\text{s.t.} \quad \widetilde{G} = g(f, (G, c, y)),$$
$$\mathcal{I}(G, \widetilde{G}, c) = 1,$$

where $\mathcal{I} : \mathcal{G} \times \mathcal{G} \times V \rightarrow \{0, 1\}$, named an equivalence indicator, checks whether two graph $G$ and $\widetilde{G}$ 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 $\mathcal{I}\left(G, \widetilde{G}, c\right) = \mathbb{I}\left(f^*(G, c) = f^*(\widetilde{G}, c)\right)$ where $f^*$ is a gold standard classifier, and the small modifications are defined as $\mathcal{I}\left(G, \widetilde{G}, c\right) = \mathbb{I}\left(|(E - \widetilde{E}) \cup (\widetilde{E} - E)| < m\right) \cdot \mathbb{I}\left(\widetilde{E} \subseteq N(G, b)\right)$ where $N(G, b) = \{(u, v) : u, v \in V, d_G(u, v) <= b\}$. In order to learn an attack policy, the attack procedure is modeled as a finite horizon Markov Decision Process (MDP) $\mathcal{M}_m(f, G, c, y)$ and is therefore optimized by Q-learning with a hierarchical Q-function. For the MDP $\mathcal{M}_m(f, G, c, y)$, its action $a_t \in \mathcal{A} \subseteq V \times V$ at time step $t$ is defined to add or delete edges in the graph, its state $(\widehat{G}_t, c)$ at time step $t$ is a partially modified graph with some of the edges added/deleted from $G$, and the reward function is defined as

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

Note that the GCNNs are employed to parameterize the Q-function. The NETTACK [251] considers attacker nodes in $\mathcal{A}$ to satisfy a feature attack constraint $X'_{u,j} \neq X^{(0)}_{u,j} \implies u \in \mathcal{A}$, a structure attack constraint $A'_{u,v} \neq A^{(0)}_{u,v} \implies u \in \mathcal{A} \vee v \in \mathcal{A}$ and an equivalence indicator constraint $\sum_u \sum_j |X^{(0)}_{u,j} - X'_{u,j}| + \sum_{u<v} |A^{(0)}_{u,v} - A'_{u,v}| \leq \Delta$, where $G'$ is derived by perturbing $G^{(0)}$. Let $\mathcal{P}^{G0}_{\Delta, \mathcal{A}}$ denote the set of all perturbed graphs $G'$ satisfying these three constraints. 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. $\max_{(A', X') \in \mathcal{P}^{G0}_{\delta, \mathcal{A}}} \max_{c_{\text{new}} \neq c_{\text{old}}} \log Z^*_{v_0, c_{\text{new}}} - \log Z^*_{v_0, c_{\text{old}}}$ where $Z^* = f_{\theta^*}(A', X')$ with $\theta^* = \arg\min_\theta \mathcal{L}(\theta; A', X')$. The NETTACK employs the GCNNs to model the classifier. The literature [32] adopts the similar equivalence indicator, and poisoning attacks are mathematically formulated as a bilevel optimization problem, i.e.

$$\begin{aligned} \min_{\widetilde{G}} \quad & \mathcal{L}_{\text{attack}}\left(f_{\theta^*}\left(\widetilde{G}\right)\right) \\ s.t. \quad & \widetilde{G} \in \mathcal{P}^{G0}_{\Delta, \mathcal{A}} \\ & \theta^* = \arg\min_\theta \mathcal{L}_{\text{train}}\left(f_\theta\left(\widetilde{G}\right)\right). \end{aligned} \tag{18}$$

This bilevel optimization problem in formula (18) is then tackled using meta-gradients, whose core idea is to treat the adjacency matrix of the input graph as a hyperparameter.

*4.6.2 Defense against the Adversarial Attacks.* A robust GCNN requires that it is invulnerable to perturbations of the input graph. The robust GCN (RGCN) [96] 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

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

where $\alpha_j^{(k)}$ are attention weights. Finally, the overall loss function is defined as regularized cross-entropy. The literature [245] presents a batch virtual adversarial training method which appends a

novel regularization term to the conventional objective function of the GCNNs, i.e.

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

where $\mathcal{L}_0$ is an average cross-entropy loss of all labelled nodes, $E(\cdot)$ is the conditional entropy of a distribution, and $\mathcal{R}_{vadv}(V, W)$ is the average Local Distributional Smoothness (LDS) loss for all nodes. Specifically, $\mathcal{R}_{vadv}(V, W) = \frac{1}{N} \sum_{u \in V} \text{LDS}(X[u,:], W, r_{vadv,u})$ where $\text{LDS}(x, w, r_{vadv}) = D_{KL}\left(p(y|x, \widehat{W})||p(y|x + r_{vadv}, W)\right)$ and $r_{vadv}$ is the virtual adversarial perturbation. Additionally, there are some other studies aiming at verifying certifiable (non-)robustness to structure and feature perturbations for the GCNNs and developing robust training algorithm [3, 33]. The literature [96] 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.

## 4.7 Graph Neural Architecture Search

Neural Architecture Search (NAS) [196] 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 [95, 217] 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 [19] 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.

## 4.8 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 [223] 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 [186] propagates information over the underlying graph of an agent via a GCNN, and then predicts actions for different parts of the agent. The literature [143] 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.

## 4.9 Applications

In this subsection, we introduce the applications of the GNNs. Due to the space limitation, we only list the application fields, including complex network analysis [83, 86, 129], combinatorial optimization [76, 116, 164, 205], knowledge graph [12, 82, 109], bioinformatics [43, 117, 144], chemistry [27, 65, 135], brain network analysis [70, 179], physical system [4, 134, 222], source

code analysis [71, 79, 126], intelligent traffic [22, 159, 180], recommender systems [177, 207, 211], computer vision [46, 172, 177, 211, 226] and natural language processing [31, 107, 127, 156, 244].

## 5 BENCHMARKS AND EVALUATION PITFALLS

In this section, we briefly introduce benchmarks and evaluation Pitfalls. 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.

*5.0.1 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 [191] 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 [128], introduces medium-scale graph datasets with 12k-70k graphs of variable sizes 9-500 nodes, and identifies important building blocks of GNNs (graph convolutions, anistropic diffusion, residual connections and normalization layers) with the proposed benchmark infrastructures. The literature [199] 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.

*5.0.2 Evaluation Pitfalls.* The literature [139] compares four typical GCNN architectures: GCN [185], MoNet [45], GAT [145] and GraphSAGE using three aggregation strategies [204] 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.

*6.0.1 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.

*6.0.2 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.

*6.0.3   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 motivates 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.

*6.0.4   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 for the GNNs. We expand the content of this paper from 4 dimensions: architectures, extensions and applications, benchmarks and evaluation pitfalls, and future research directions. The GNN architectures are expanded from 4 perspectives: graph convolutional neural networks, graph pooling operators, graph attention mechanisms and graph recurrent neural networks. The extensions are expanded from 8 perspectives: GCNNS on special graphs, capability and interpretability, deep graph representation learning, deep graph generative models, combinations of the PI and GNNs adversarial attachks for the GNNs, graph neural architecture search and graph reinforcement. In the future directions, we propose 4 prospective topics on the GNNs: highly scalable GNNs, robust GNNs, GNNs going beyond WL test and interpretable GNNs. We expect 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.

## ACKNOWLEDGMENTS

## REFERENCES

[1]   Aditya Grover and Jure Leskovec. 2016.   node2vec: scalable feature learning for networks. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 855–864.

[2]   Afshin Rahimi, Trevor Cohn, and Timothy Baldwin. 2018. Semi-supervised user geolocation via graph convolutional networks. In *The Annual Meeting of the Association for Computational Linguistics (ACL)*. 2009–2019.

[3]   Aleksandar Bojchevski and Stephan Günnemann. 2019.   Certifiable robustness to graph perturbations.   (2019). https://arxiv.org/abs/1910.14356

[4]   Alvaro Sanchez-Gonzalez, Nicolas Heess, Jost Tobias Springenberg, Josh Merel, Martin Riedmiller, Raia Hadsell, and Peter Battaglia. 2018. Graph networks as learnable physics engines for inference and control. In *The International Conference on Machine Learning (ICML)*. 4470–4479.

[5]   Amir H. Khasahmadi, Kaveh Hassani, Parsa Moradi, Leo Lee, and Quaid Morris. 2020. Memory-based graph networks. In *The International Conference on Learning Representations (ICLR)*.

[6]   Ana {S}u{s}njara, Nathanaël Perraudin, Daniel Kressner, and Pierre Vandergheynst. 2015.   Accelerated filtering on graphs using Lanczos method. (2015).   https://arxiv.org/abs/1509.04537

[7]   Andreas Loukas. 2020. What graph neural networks cannot learn-depth vs width. In *The International Conference on Learning Representations (ICLR)*.

[8] Antoni Buades, Bartomeu Coll, and Jean-Michael Morel. 2005. A non-local algorithm for image denoising. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. 60–65.

[9] Anuththari Gamage, Eli Chien, Jianhao Peng, and Olgica Milenkovic. 2019. Multi-MotifGAN (MMGAN): motif-targeted graph generation and prediction. (2019). https://arxiv.org/abs/1911.05469v1

[10] Ashesh Jain, Amir R. Zamir, Silvio Savarese, and Ashutosh Saxena. 2016. Structural-RNN: deep learning on spatio-temporal graphs. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5308–5317.

[11] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomiz, ł{L}ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *The International Conference on Neural Information Processing (NIPS)*. 5998–6008.

[12] Bill Yuchen Lin, Xinyue Chen, Jamin Chen, and Xiang Ren. 2019. KagNet: knowledge-aware graph networks for commonsense reasoning. In *The International Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 2829–2839.

[13] Bing Yu, Haoteng Yin, and Zhanxing Zhu. 2018. Spatio-temporal graph convolutional networks: a deep learning framework for traffic forecasting. In *The International Joint Conference on Artificial Intelligence (IJCAI)*. 3634–3640.

[14] Bingbing Xu, Huawei Shen, Qi Cao, Yunqi Qiu, and Xueqi Cheng. 2019. Graph Wavelet Neural Network. In *The International Conference on Learning Representations (ICLR)*.

[15] Bryan Perozzi, Rami Ai-Rfou, and Steven Skiena. 2014. DeepWalk: online learning of social representations. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 701–710.

[16] Caglar Gulcehre, Misha Denil, Mateusz Malinowski, Ali Razavi, Razvan Pascanu, Karl Moritz Hermann, Peter Battaglia, Victor Bapst, David Raposo, Adam Santoro, and Nando de Freitas. 2019. Hyperbolic Attention Networks. In *The International Conference on Learning Representations*.

[17] Chao Shang, Qinqing Liu, Ko-Shin Chen, Jiangwen Sun, Jin Lu, Jinfeng Yi, and Jinbo Bi. 2018. Edge attention-based multi-relational graph convolutional networks. (2018). https://arxiv.org/abs/1802.04944

[18] Chenyi Zhuang and Qiang Ma. 2018. Dual graph convolutional networks for graph-based semi-supervised classification. In *The World Wide Web Conference (WWW)*. 499–508.

[19] Chris Zhang, Mengye Ren, and Raquel Urtasun. 2019. Graph hypernetworks for neural architecture search. In *The International Conference on Learning Representations (ICLR)*.

[20] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. 1–9.

[21] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. 2019. Weisfeiler and Leman Go neural: higher-order graph neural networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 4602–4609.

[22] Chuanpan Zheng, Xiaoliang Fan, Cheng Wang, and Jianzhong Qi. 2020. GMAN: a graph multi-attention network for traffic prediction. In *The AAAI Conference on Artificial Intelligence (AAAI)*.

[23] Chun Wang, Shirui Pan, Guodong Long, Xingquan Zhu, and Jing Jiang. 2017. MGAE: marginalized graph autoencoder for graph clustering. In *The International Conference on Information and Knowledge Management (CIKM)*. 889–898.

[24] Fan R.K. Chung. 1992. *Spectral graph theory*. American Mathematical Society.

[25] Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V. Chawla. 2019. Heterogeneous Graph Neural Network. In *Th ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 793–803.

[26] Claire Donnat, Marinka Zitnik, David Hallac, and Jure Leskovec. 2018. Learning structural node embeddings via diffusion wavelets. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 1320–1329.

[27] Connor W. Coley, Wengong Jin, Luke Rogers, Timothy F. Jamison, Tommi S. Jaakkola, William H. Green, Regina Barzilay, and Klavs F. Jensen. 2019. A graph-convolutional neural network model for the prediction of chemical reactivity. *Chemical Science* 10, 2 (2019), 370–377.

[28] C̆{a}t{a}lina Cangea, Petar Veli{c}kovi{c}, Nikola Jovanović, Thomas N. Kipf, and Pietro Liò. 2018. Towards sparse hierarchical graph classifiers. (2018). https://arxiv.org/abs/1811.01287

[29] Daixin Wang, Peng Cui, and Wenwu Zhu. 2016. Structural Deep Network Embedding. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 1225–1234.

[30] Dan Busbridge, Dane Sherburn, Pietro Cavallo, and Nils Y. Yhammerla. 2019. Relational graph attention networks. (2019). https://arxiv.org/abs/1904.05811

[31] Daniel Beck, Gholamreza Haffari, and Trevor Cohn. 2018. Graph-to-Sequence learning using gated graph neural networks. In *The Annual Meeting of the Association for Computational Linguistics*. 273–283.

[32] Daniel Zügner and Stephan Günnemann. 2019. Adversarial attacks on graph neural networks via meta learning. In *The International Conference on Learning Representations (ICLR)*.

[33] Daniel Zügner and Stephan Günnemann. 2019. Certifiable robustness and robust training for graph convolutional networks. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 246–256.

[34] David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gomez-Bombarelli, Timothy Hirzel, Alan Aspuru-Guzik, and Ryan P. Adams. 2015. Convolutional networks on graphs for learning molecular fingerprints. In *The International Conference on Neural Information Processing Systems*. 2224–2232.

[35] David I. Shuman, Sunil K. Narang, Pascal Frossard, Antonio Ortega, and Pierre Vandergheynst. 2013. The emerging field of signal processing on graphs: extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine* 30, 3 (2013), 83–98.

[36] David K. Hammond, Pierre Vandergheynst, and Rémi Gribonval. 2011. Wavelets on graphs via spectral graph theory. *Applied and Computational Harmonic Analysis* 30, 2 (2011), 129–150.

[37] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. 2017. Variational inference: a review for statisticians. *Journal of the Americian Statistical Association* 112, 518 (2017), 859–877.

[38] Dingyuan Zhu, Peng Cui, Daixin Wang, and Wenwu Zhu. 2018. Deep Variational Network Embedding in Wasserstein Space. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2827–2836.

[39] Dongmian Zou and Gilad Lerman. 2019. Graph convolutional neural networks via scattering. *Applied and Computational Harmonic Analysis* (2019).

[40] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *The International Conference on Learning Representations (ICLR)*.

[41] Edouard Pineau and Nathan de Lara. 2019. Variational recurrent neural networks for graph classification. (2019). https://arxiv.org/abs/1902.02721

[42] Ehsan Hajiramezanali, Arman Hasanzadeh, Krishna Narayanan, Nick Duffield, Mingyuan Zhou, and Xiaoning Qian. 2019. Variational Graph Recurrent Neural Networks. In *The International Conference on Neural Information Processing Systems*. 10701–10711.

[43] Emanuele Rossi, Federico Monti, Michael Bronstein, and Pietro Liò. 2019. ncRNA classification with graph convolutional networks. (2019). https://arxiv.org/abs/1905.06515

[44] Federico Baldassarre and Hossein Azizpour. 2019. Explainability techniques for graph convolutional networks. (2019). https://arxiv.org/abs/1905.13686

[45] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodolà, Jan Svoboda, and Michael M. Bronstein. 2017. Geometric deep learning on graphs and manifolds using mixture model CNNs. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5425–5434.

[46] Federico Monti, Michael M. Bronstein, and Xavier Bresson. 2017. Geometric matrix completion with recurrent multi-graph neural network. In *The International Conference on Neural Information Processing Systems*. 3697–3707.

[47] Felipe Petroski Such, Shagan Sah, Miguel Dominguez, Suhas Pillai, Chao Zhang, Andrew Michael, Nathan D. Cahill, and Raymond Ptucha. 2017. Robust spatial filtering with graph convolutional neural networks. *IEEE Journal of Selected Topics in Signal Processing* 11, 6 (2017), 884–896.

[48] Felix Wu, Tianyi Zhang, Amauri Holanda de Souza, Christopher Fifty, Tao Yu, and Kilian Q. Weinberger. 2019. Simplifying Graph Convolutional Networks. In *The International Conference on Machine Learning (ICML)*. 6861–6871.

[49] Fenyu Hu, Yanqiao Zhu, Shu Wu, Liang Wang, and Tieniu Tan. 2019. Semi-supervised node classification via hierarchical graph convolutional networks. (2019). https://arxiv.org/abs/1902.06667v2

[50] Fisher Yu and Vladlen Koltun. 2016. Multi-scale context aggregation by dilated convolutions. (2016). https://arxiv.org/abs/1511.07122v3

[51] Franco Manessi, Alessandro Rozza, and Mario Manzo. 2020. Dynamic graph convolutional networks. *Pattern Recognition* 97, 2020 (2020), No. 107000.

[52] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural network model. *IEEE Transactions on Neural Networks* 20, 1 (2008), 61–80.

[53] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2009. Computational Capabilities of Graph Neural Networks. *IEEE Transactions on Neural Networks* 20, 1 (2009), 81–102.

[54] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q. Weinberger. 2017. Densely connected convolutional networks. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. 2261–2269.

[55] Marco Gori, Gabriele Monfardini, and Franco Scarselli. 2005. A new model for learning in graph domains. In *The International Joint Conference on Neural Networks (IJCNN)*. 729–734.

[56] Guohao Li, Matthias Müller, Ali Thabet, and Bernard Ghanem. 2019. DeepGCNs: can GCNs go as deep as CNNs. In *The IEEE International Conference on Computer Vision (ICCV)*. 9267–9276.

[57] Haggai Maron, Heli Ben-Hamu, Nadav Shamir, and Yaron Lipman. 2019. Invariant and equivariant graph networks. In *The International Conference on Learning Representations (ICLR)*.

[58] Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. 2018. Adversarial Attack on Graph Structured Data. In *The International Conference on Machine Learning (ICML)*. 1115–1124.

[59] Hao Yuan and Shuiwang Ji. 2020. StructPool: structured graph pooling via conditional random fields. In *The International Conference on Learning Representations (ICLR)*.

[60] Heng Chang, Yu Rong, Tingyang Xu, Wenbing Huang, Somayeh Sojoudi, Junzhou Huang, and Wenwu Zhu. 2020. Spectral graph attention network. (2020). https://arxiv.org/abs/2003.07450

[61] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. 2020. Geom-GCN: geometric graph convolutional networks. In *The International Conference on Learning Representations (ICLR)*.

[62] Hongchang Gao, Jian Pei, and Heng Huang. 2019. Conditional Random Field Enhanced Graph Convolutional Neural Networks. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 276–284.

[63] Hongyang Gao and Shuiwang Ji. 2019. Graph U-Nets. (2019). https://arxiv.org/abs/1905.05178

[64] Hongyang Gao, Zhengyang Wang, and Shuiwang Ji. 2018. Large-scale learnable graph convolutional networks. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 1415–1424.

[65] Hyeoncheol Cho and Insung S. Choi. 2018. Three-dimensionally embedded graph convolutional network (3DGCN) for molecule interpretation. (2018). https://arxiv.org/abs/1811.09794

[66] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Azron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *The International Conference on Neural Inforamtion Processing Systems (NIPS)*. 2672–2680.

[67] Inderjit S. Dhillon, Yuqiang Guan, and Brian Kulis. 2007. Weighted graph cuts without eigenvectors: a multilevel approach. *IEEE Transactions on Pattern Analysis and Mchine Intelligence* 29, 11 (2007), 1944–1957.

[68] Ines Chami, Rex Ying, Christopher Re, and Jure Leskovec. 2019. Hyperbolic Graph Convolutional Neural Networks. In *The International Conference on Neural Information Processing Systems (NeurPS)*. 4868–4879.

[69] James Atwood and Don Towsley. 2016. Diffusion-Convolutional Neural Network. In *The International Conference on Neural Information Processing Systems (NIPS)*. 1993–2001.

[70] Jeremy Kawahara, Colin J. Brown, Steven P. Miller, Brian G. Booth, Vann Chau, Ruth E. Grunau, Jill G. Zwicker, and Ghassan Hamarneh. 2017. BrainNetCNN: convolutional neural networks for brain networks; towards predicting neurodevelopment. *NeuroImage* 146 (2017), 1038–1049.

[71] Jessica Schrouff, Kai Wohlfahrt, Bruno Marnette, and Liam Atkinson. 2019. Inferring Javascript types using graph neural networks. (2019). https://arxiv.org/abs/1905.06707

[72] Jian Du, Shanghang Zhang, Guanhang Wu, José M.F. Moura, and Soummya Kar. 2018. Topology adaptive graph convolutional networks. (2018). https://arxiv.org/abs/1710.10370

[73] Jianfei Chen, Jun Zhu, and Le Song. 2018. Stochastic training of graph convolutional networks with variance reduction. In *The International Conference on Machine Learning (ICML)*. 942–950.

[74] Jiani Zhang, Xingjian Shi, Junyuan Xie, Hao Ma, Irwin King, and Dit-Yan Yeung. 2018. GaAN: gated attention networks for learning on large and spatiotemporal graphs. In *The International Conference on Uncertainty in Artificial Intelligence (UAI)*. No. 139.

[75] Jiawei Zhang, Haopeng Zhang, Congying Xia, and Li Sun. 2020. Graph-Bert: only attention is needed for learning graph representations. (2020). https://arxiv.org/abs/2001.05140

[76] Jiaxuan You, Haoze Wu, Clark Barrett, Raghuram Ramanujan, and Jure Leskovec. 2019. G2SAT: learning to generate SAT formulas. In *The International Conference on Neural Information Processing Systems (NeruPS)*. 10552–10563.

[77] Jiaxuan You, Rex Ying, Xiang Ren, William L. Hamilton, and Jure Leskovec. 2018. GraphRNN: generating realistic graphs with deep auto-regressive models. In *The Internatinonal Conference on Machine Learing (ICML)*. 5708–5717.

[78] Jiaxue You, Rex Ying, and Jure Leskovec. 2019. Position-aware graph neural networks. In *The International Conference on Machine Learning (ICML)*. 7134–7143.

[79] Jiayi Wei, Maruth Goyal, Greg Durrett, and Isil Dilling. 2020. LambdaNet: probabilistic type inference using graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[80] Jie Chen, Tengfei Ma, and Cao Xiao. 2018. FastGCN: fast learning with graph convolutional networks via importance sampling. In *The International Conference on Learning Representations (ICLR)*.

[81] Jie Zhou, Gangqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, and Maosong Sun. 2018. Graph neural networks: a review of methods and applications. (2018). https://arxiv.org/abs/1812.08434

[82] Jie Zhou, Xu Han, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. GEAR: graph-based evidence aggregating and reasoning for fact verification. In *The Annual Meeting of the Association for Computational Linguistics (ACL)*. 892–901.

[83] Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. 2018. DeepInf: social influence prediction with deep learning. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2110–2119.

[84] Jiwoong Park, Minsik Lee, Hyung Jin Chang, Kyuewang Lee, and Jin Young Choi. 2019. Symmetric graph convolutional autoencoder for unsupervised graph representation learning. In *The International Conference on Computer Vision*. 6519–6528.

[85] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. 2014. Spectral networks and locally connected networks on graphs. In *The International Conference on Learning Representations (ICLR)*.

[86] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. 2019. Predict the propagate: graph neural networks meet personalized pagerank. In *The International Conference on Learning Representations (ICLR)*.

[87] John B. Lee, Ryan A. Rossi, Sungchul Kim, Nesreen K. Ahmed, and Eunyee Koh. 2019. Attention models in graphs: a survey. *ACM Transactions on Knowledge Discovery from Data* 13, 6 (2019), No. 62.

[88] John Boaz Lee, Ryan Rossi, and Xiangnan Kong. 2018. Graph classification using structural attention. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 1666–1674.

[89] John Lafferty, Andrew McCallum, and Fernando CN Pereira. 2001. Conditional random field: probabilistic models for segmenting and labeling sequence data. In *The International Conference on Machine Learning (ICML)*. 282–289.

[90] Jun Wu, Jingrui He, and Jiejun Xu. 2019. DEMO-Net: degree-specific graph neural networks for node and graph classification. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 406–415.

[91] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. 2019. Self-Attention Graph Pooling. In *The International Conference on Machine Learning (ICML)*. 3734–3743.

[92] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. 2017. Neural Message Passing for Quantum Chemistry. In *The International Conference on Machine Learning (ICML)*. 1263–1272.

[93] Kai Sheng Tai, Richard Socher, and Christopher D. Manning. 2015. Improved semantic representations from tree-structured long short-term memory networks. In *The Annual Meeting of the Association for Computational Linguistics*. 1556–1566.

[94] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. 770–778.

[95] Kaixiong Zhou, Qingquan Song, Xiao Huang, and Xia Hu. 2019. Auto-GNN: neural architecture search of graph neural networks. (2019). https://arxiv.org/abs/1909.03184

[96] Ke Sun, Piotr Koniusz, and Zhen Wang. 2019. Fisher-Bures adversary graph convolutional networks. In *The International Conference on Uncertainty in Artificial Intelligence*. No. 161.

[97] Ke Tu, Peng Cui, Xiao Wang, Fei Wang, and Wenwu Zhu. 2018. Structural Deep Embedding for Hyper-Networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 426–433.

[98] Ke Tu, Peng Cui, Xiao Wang, and Philip S. Yu. 2018. Deep Recursive Network Embedding with Regular Equivalence. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2357–2366.

[99] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation Learning on Graphs with Jumping Knowledge Networks. In *The International Conference on Machine Learning (ICML)*. 5453–5462.

[100] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S. Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2020. What can neural networks reason about. In *The International Conference on Learning Representations (ICLR)*.

[101] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How powerful are graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[102] KiJung Yoon, Renjie Liao, Yuwen Xiong, Lisa Zhang, Ethan Fetaya, Raquel Urtasun, Richard Zemel, and Xaq Pitkow. 2018. Inference in probabilistic graphical models by graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[103] Kilian Weinberger, Anirban Dasgupta, and John Langford. 2009. Feature hashing for large scale multitask learning. In *The International Conference on Machine Learning (ICML)*. 1113–1120.

[104] Kiran K. Thekumparampil, Chong Wang, Sewoong Oh, and Lijia Li. 2018. Attention-based graph neural network for semi-supervised learning. (2018). https://arxiv.org/abs/1803.03735

[105] Kun Xu, Lingfei Wu, Zhiguo Wang, Yansong Feng, Michael Witbrock, and Vadim Sheinin. 2018. Graph2Seq: graph to sequence learning with attention-based neural networks. (2018). https://arxiv.org/abs/1804.00823

[106] Kyunghyun Cho, Bart van Merrienboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Youshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *The Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 1724–1734.

[107] Liang Yao, Chengsheng Mao, and Yuan Luo. 2019. Graph convolutional networks for text classification. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 7370–7377.

[108] Lichao Sun, Yingtong Dou, Carl Yang, Ji Wang, Philip S. Yu, and Bo Li. 2020. Adversarial attack and defense on graph data: a survey. (2020). https://arxiv.org/abs/1812.10528

[109] Lin Qiu, Yunxuan Xiao, Yanru Qu, Hao Zhou, Lei Li, Weinan Zhang, and Yong Yu. 2019. Dynamically fused graph network for multi-hop reasoning. In *The Annual Meeting of the Association for Computational Linguistics (ACL)*. 6140–6150.

[110] Linfeng Liu and Liping Liu. 2019. Amortized Variational Inference with Graph Convolutional Networks for Gaussian Processes. In *The International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2291–2300.

[111] Lingxiao Ma, Zhi Yang, Youshan Miao, Jilong Xue, Ming Wu, Lidong Zhou, and Yafei Dai. 2018. Towards efficient large-scale graph neural network computing. (2018). https://arxiv.org/abs/1810.08403

[112] Lingxiao Zhao and Leman Akoglu. 2020. PairNorm: tackling oversmoothing in GNNs. In *The International Conference on Learning Representations (ICLR)*.

[113] Louis Tiao, Pantelis Elinas, Harrison Nguyen, and Edwin V. Bonilla. 2019. Variational spectral graph convolutional networks. (2019). https://arxiv.org/abs/1906.01852v1

[114] Luana Ruiz, Fernando Gama, and Alejandro Ribeiro. 2019. Gated Graph Convolutional Recurrent Neural Networks. (2019). https://arxiv.org/abs/1903.01888

[115] Marcelo Daniel Gutierrez Mallea, Peter Meltzer, and Peter J. Bentley. 2019. Capsule neural networks for graph classification using explicit tensorial graph representations. (2019). https://arxiv.org/abs/1902.08399

[116] Marcelo Prates, Pedro H.C. Avelar, Henrique Lemos, Luis C. Lamb, and Moshe Y. Vardi. 2019. Leaning to solve NP-Complete problems: a graph neural network for decision TSP. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 4731–4738.

[117] Marnka Zitnik, Monica Agrawal, and Jure Leskovec. 2018. Modeling polypharmacy side effects with graph convolutional networks. *Bioinformatics* 34, 13 (2018), i457–i466.

[118] Martin Simonovsky and Nikos Komodakis. 2017. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 29–38.

[119] Martin Simonovsky and Nikos Komodakis. 2018. GraphVAE: towards generation of small graphs using variational autoencoders. In *The International Conference on Artificial Neural Networks (ICANN)*. 412–422.

[120] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. 2016. Learning convolutional neural networks for graphs. In *The International Conference on Machine Learning (ICLR)*. 2014–2023.

[121] Matthew Baron. 2018. Topology and prediction focused research on graph convolutional neural networks. (2018). https://arxiv.org/abs/1808.07769v1

[122] Meng Qu, Yoshua Bengio, and Jian Tang. 2019. GMNN: Graph Markov Neural Networks. In *The International Conference on Machine Learning (ICML)*. 5241–5250.

[123] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems (NIPS)*. 3844–3852.

[124] Michael M. Bronstein, Joan Bruna, Yan LeCun, Arthur Szlam, and Pierre Vandergheynst. 2017. Geometric deep learning: going beyond Euclidean data. *IEEE Signal Processing Magazine* 34, 4 (2017), 18–42.

[125] Mikael Henaff, Joan Bruna, and Yan LeCun. 2015. Deep convolutional networks on graph-structured data. (2015). https://arxiv.org/abs/1506.05163

[126] Miltiadis Allamanis, Marc Brockschmidt, and Mahmoud Khademi. 2018. Learning to represent programs with graphs. In *The International Conference on Learning Representations (ICLR)*.

[127] Ming Ding, Chang Zhou, Qibin Chen, Hongxia Yang, and Jie Tang. 2019. Cognitive graph for multi-hop reading comprehension at scale. In *The Annual Meeting of the Association for Computational Linguistics (ACL)*. 2694–2703.

[128] 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, and Zheng Zhang. 2019. Deep Graph Library: towards efficient and scalable deep learning on graphs. In *The Workshop at the International Conference on Learning Representations (ICLR Workshop)*.

[129] Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. In *The International Conference on Neural Information Processing Systems (NIPS)*. 5171–5181.

[130] Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. 2018. An end-to-end deep learning architecture for graph classification. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 4438–4445.

[131] Naganand Yadati, Madhav Nimishakavi, Prateek Yadav, Vikram Nitin, Anand Louis, and Partha Talukdar. 2019. HyperGCN: a new method for training graph convolutional networks on hypergraphs. In *The International Conference on Neural Information Processing Systems (NeurPS)*. 1511–1522.

[132] Naganand Yadati, Tingran Gao, Shahab Asoodeh, Partha Talukdar, and Anand Louis. 2020. Graph neural networks for soft semi-supervised learning on hypergraphs. (2020).

[133] Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. 2017. Cross-sentence N-ary relation extraction with graph LSTMs. *Transactions of the Association for Computational Linguistics* 5, 2017 (2017), 101–115.

[134] Nicholas Watters, Andrea Tacchetti, Théophane Weber, Razvan Pascanu, Peter Battaglia, and Daniel Zoran. 2017. Visual interaction networks: learning a physics simulator from video. In *The International Conference on Neural Inforamtion Processing Systems (NIPS)*. 4539–4547.

[135] Nicola De Cao and Thomas N. Kipf. 2018. MolGAN: an implicit generative model for small molecular graphs. (2018). https://arxiv.org/abs/1805.11973

[136] Nicolas Keriven and Gabriel Peyré. 2019. Universal invariant and equivariant graph neural networks. In *The International Conference on Neural Information Processing Systems (NeurPS)*. 7092–7101.

[137] Nima Dehmamy, Albert-Laszlo Barabasi, and Rose Yu. 2019. Understanding the Representation Power of Graph Neural Networks in Learning Graph Topology. In *The International Conference on Neural Information Processing*

*Systems (NeurPS)*. 15413–15423.

[138] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-Net: convolutional networks for biomedical image segmentation. In *The International Conference on Medical Image Computing and Computer-assisted Intervention (MCCAI)*. 234–241.

[139] Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. 2019. Pitfalls of graph neural network evaluation. (2019). https://arxiv.org/abs/1811.05868

[140] Pablo Barceló, Egor V. Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, and Juan Pablo Silva. 2020. The logical expressiveness of graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[141] Padraig Corcoran. 2019. Function space pooling for graph convolutional networks. (2019). https://arxiv.org/abs/1905.06259

[142] Palash Goyal, Nitin Kamra, Xinran He, and Yan Liu. 2018. DynGEM: deep embedding method for dynamic graphs. (2018). https://arxiv.org/abs/1805.11273v1

[143] Paul Almasan, José Suárez-Varela, Arnau Badia-Sampera, Krzysztof Rusek, Pere Barlet-Ros, and Albert Cabellos-Aparicio. 2019. Deep reinforcement learning meets graph neural networks: exploring a routing optimization use case. (2019). https://arxiv.org/abs/1910.07421

[144] Peng Han, Peng Yang, Peilin Zhao, Shuo Shang, Yong Liu, Jiayu Zhou, Xin Gao, and Panos Kalnis. 2019. GCN-MF: disease-gene association identification by graph convolutional networks and matrix factorization. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 705–713.

[145] Petar Veli{c}ković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *The International Conference on Learning Representations (ICLR)*.

[146] Petar Veli{c}ković, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R Devon Hjelm. 2019. Deep graph infomax. In *The International Conference on Learning Representations (ICLR)*.

[147] Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, and Koray Kavukcuoglu. 2016. Interaction networks for learning about objects, relations and physics. In *The International Conference on Neural Information Processing Systems (NIPS)*. 4502–4510.

[148] Peter Meltzer, Marcelo Daniel Gutierrez Mallea, and Peter J. Bentley. 2019. PiNet: a permutation invariant graph neural network for graph classification. (2019). https://arxiv.org/abs/1905.03046

[149] Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, Caglar Gulcehre, Francis Song, Andrew Ballard, Justin Gilmer, George E. Dahl, Ashish Vaswani, Kelsey Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan Wierstra, Pushmeet Kohli, Matt Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu. 2018. Relational inductive biases, deep learning, and graph networks. (2018). https://arxiv.org/abs/1806.01261

[150] Qimai Li, Zhichao Han, and Xiaoming Wu. 2018. Deeper insights into graph convolutional networks for semi-supervised learning. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 3538–3545.

[151] Radford M. Neal and Geoffrey E. Hinton. [n.d.]. A view of the EM algorithm that justifies incremental, sparse and other variants. ([n. d.]).

[152] Renjie Liao, Marc Brockschmidt, Daniel Tarlow, Alexander L. Gaunt, Raquel Urtasun, and Richard Zemel. 2018. Graph partition neural networks for semi-supervised classification. In *The International Conference on Learning Representations (ICLR Workshop)*.

[153] Renjie Liao, Zhizhen Zhao, Raquel Urtasun, and Richard Zemel. 2019. LanczosNet: Multi-scale deep graph convolutional networks. In *The International Conference on Learning Representations (ICLR)*.

[154] Rex Ying, Jiaxuan You, Christopher Morris, Xiang Ren, William L. Hamilton, and Jure Leskovec. 2018. Hierarchical graph representation learning with differentiable pooling. In *The International Conference on Nueral Information Processing Systems (NeurPS)*. 4805–4815.

[155] Richard C. Wilson, Edwin R. Hancock, Elżbieta Pekalska, and Robert P.W. Duin. 2014. Spherical and Hyperbolic Embeddings of Data. *IEEE Transactions on Pattern Recognition and Machine Intelligence* 36, 11 (2014), 2255–2269.

[156] Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. Text generation from knowledge graphs with graph transformers. (2019). https://arxiv.org/abs/1904.02342

[157] Risi Kondor, Truong Son Hy, Horace Pan, Brandon M. Anderson, and Shubhendu Trivedi. 2018. Covariant compositional networks for learning graphs. In *The International Conference on Learning Representations (ICLR Workshop)*.

[158] Ron Levie, Federico Monti, Xavier Bresson, and Michael M. Bronstein. 2019. CayleyNets: Graph Convolutional Neural Networks with Complex Relational Spectral Filters. *IEEE Transactions on Signal Processing* 67, 1 (2019), 97–109.

[159] Rui Dai, Shenkun Xu, Qian Gu, Chenguang Ji, and Kaikui Liu. 2020. Hybrid spatio-temporal graph convolutional network: improving traffic prediction with navigation data. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 3074–3082.

[160] Ruoyu Li, Sheng Wang, Feiyun Zhu, and Junzhou Huang. 2018. Adaptive Graph Convolutional Neural Networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 3546–3553.

[161] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. 2015. Highways networks. (2015). https://arxiv.org/abs/1505.00387

[162] Ryan Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. 2019. Relational pooling for graph representations. In *The International Conference on Machine Learning (ICML)*. 4663–4673.

[163] Ryoma Sato. 2020. A survey on the expressive power of graph neural networks. (2020). https://arxiv.org/abs/2003.04078

[164] Ryoma Sato, Makoto Yamada, and Hisashi Kashima. 2019. Approximation Ratios of Graph Neural Networks for Combinatorial Problems. In *The International Conference on Neural Information Processing Systems (NeurPS)*. 4081–4090.

[165] Sami Abu-EL-Haija, Amol Kapoor, Bryan Perozzi, and Joonseok Lee. 2019. N-GCN: multi-scale graph convolution for semi-supervised node classification. In *The Conference on Uncertainty in Artificial Intelligence (UAI)*. No. 310.

[166] Sami Abu-EL-Haija, Bryan Perozzi, Amol Kapoor, Nazanin Alipourfard, Kristina Lerman, Hrayr Harutyunyan, Greg Ver Steeg, and Aram Galstyan. 2019. MixHop: higher-order graph convolutional architectures via sparsified neighborhood mixing. In *The International Conference on Machine Learning (ICML)*. 21–29.

[167] Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. 2017. Dynamic routing between capsules. In *The International Conference on Neural Information Processing Systems (NIPS)*. 3856–3866.

[168] Saurabh Verma and Zhili Zhang. 2018. Graph Capsule Convolutional Neural Networks. (2018). https://arxiv.org/abs/1805.08090

[169] S{e}bastien Lerique, Jacob Levy Abitol, and Márton Karsai. 2020. Joint embedding of structure and features via graph convolutional networks. *Applied Network Science* 5, 2020 (2020), No. 5.

[170] Seongjun Yun, Minbyul Jeong, Raehyun Kim, Jaewoo Kang, and Hyunwoo J. Kim. 2019. Graph Transformer Networks. (2019). https://arxiv.org/abs/1911.06455

[171] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Computation* 9, 8 (1997), 1735–1780.

[172] Shaohua Fan, Junxiong Zhu, Xiaotian Han, Chuan Shi, Linmei Hu, Biyu Ma, and Yongliang Li. 2019. Metapath-guided heterogeneous graph neural network for intent recommendation. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2478–2486.

[173] Shaosheng Cao, Wei Lu, and Qiongkai Xu. 2016. Deep neural networks for learning graph representations. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 1145–1152.

[174] Shengding Hu, Meng Qu, Zhiyuan Liu, and Jian Tang. 2020. Transfer active learning for graph neural networks. (2020). https://openreview.net/forum?i

[175] Shengnan Guo, Youfang Lin, Ning Feng, Chao Song, and Huaiyu Wan. 2019. Attention based spatial-temporal graph convolutional networks for traffic flow forecasting. In *The AAAI Conference on Artificial Intelligence (AAAI)*. pages = 922–929,.

[176] Shirui Pan, Ruiqi Hu, Guodong Long, Jing Jiang, Lina Yao, and Chengqi Zhang. 2018. Adversarially regularized graph autoencoder for graph embedding. In *The International Joint Conference on Artificial Intelligence (IJCAI)*. 2609–2615.

[177] Shu Wu, Yuyuan Tang, Yanqiao Zhu, Liang Wang, Xing Xie, and Tieniu Tan. 2019. Session-based recommendation with graph neural networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 346–353.

[178] Sijie Yan, Yuanjun Xiong, and Dahua Lin. 2018. Spatial temporal graph convolutional networks for skeleton-based action recognition. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 7444–7452.

[179] Sofia Ira Ktena, Sarah Parisot, Enzo Ferrante, Martin Rajchl, Matthew Lee, Ben Glocker, and Daniel Ruekert. 2017. Distance metric learning using graph convolutional networks: application to functional brain networks. In *The International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. 469–477.

[180] Songtao He, Favyen Bastani, Satvat Jagwani, Edward Park, Sofiane Abbar, Mohammad Alizadeh, Hari Balakrishnan, Sanjay Chawla, Samuel Madden, and Mohammad Amin Sadeghi. 2020. RoadTagger: robust road attribute inference with graph neural networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*.

[181] Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay S. Pande, and Patrick F. Riley. 2016. Molecular graph convolutions: moving beyond fingerprints. *Journal of Computer-Aided Molecular Design* 30, 8 (2016), 595–608.

[182] Sungmin Rhee, Seokjun Seo, and Sun Kim. 2018. Hybrid approach of relation network and localized graph convolutional filtering for breast cancer subtype classification. In *Th International Joint Conference on Artificial Intelligence (IJCAI)*. 3527–3534.

[183] Tengfei Ma, Junyuan Shang, and Jimeng Sun. 2019. CGNF: conditional graph neural fields. (2019). https://openreview.net/forum?i

[184] Thomas N. Kipf and Max Welling. 2016. Variational Graph Auto-Encoders. (2016). https://arxiv.org/abs/1611.07308

[185] Thomas N. Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *The International Conference on Learning Representations (ICLR)*.

[186] Tingwu Wang, Renjie Liao, Jimmy Ba, and Sanja Fidler. 2018. NerveNet: learning structured policy with graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[187] Tong Zhang, Wenming Zheng, Zhen Cui, and Yang Li. 2018. Tensor graph convolutional neural network. (2018).
      https://arxiv.org/abs/1803.10071v1

[188] Trang Pham, Truyen Tran, Dinh Phung, and Svetha Venkatesh. 2017. Column Networks for Collective Classification.
      In *The AAAI Conference on Artificial Intelligence (AAAI)*. 2485–2491.

[189] Tyler Derr, Yao Ma, and Jiliang Tang. 2018. Signed Graph Convolutional Network. In *The IEEE International Conference on Data Mining (ICDM)*. 929–934.

[190] Victoria Zayats and Mari Ostendorf. 2018. Conversation modeling on Reddit using a graph-structured LSTM.
      *Transactions of the Association for Computational Linguistics* 6, 2018 (2018), 121–132.

[191] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. 2020. Benchmarking
      graph neural networks. (2020). https://arxiv.org/abs/2003.00982

[192] Vikas K. Garg, Stefanie Jegelka, and Tommi Jaakkola. 2020. Generalization and representational limits of graph neural
      networks. (2020). https://arxiv.org/abs/2002.06157

[193] Vikas Verma, Meng Qu, Alex Lamb, Yoshua Bengio, Juho Kannala, and Jian Tang. 2019. GraphMix: regularized
      training of graph neural networks for semi-supervised learning. (2019). https://arxiv.org/abs/1909.11715

[194] Vincenzo Di Massa, Cabriele Monfardini, Lorenzo Sarti, Franco Scarselli, Marco Maggini, and Marco Gori. 2006. A
      comparison between recursive neural networks and graph neural networks. In *The IEEE Joint Conference on Neural Networks (IJCNN)*. 778–785.

[195] Volodymyr Minh, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. 2014. Recurrent models of visual attention.
      In *The International Conference on Neural Information Processing Systems (NIPS)*. 2204–2212.

[196] Wei Li, Shaogang Gong, and Shaogang Gong. 2020. Neural graph embedding for neural architecture search. In *The AAAI Conference on Artificial Intelligence (AAAI)*.

[197] Wei Liu and Sanjay Chawla. 2009. A game theoretical model for adversarial learning. In *The IEEE International Conference on Data Mining Workshops (ICDM Workshop)*. 25–30.

[198] Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay S. Pande, and Jure Leskovec. 2020. Strategies
      for pre-training graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[199] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec.
      2020. Open graph benchmark: datasets for machine learning on graphs. (2020). https://arxiv.org/abs/2005.00687

[200] Weilin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-GCN: an efficient algorithm
      for training deep and large graph convolutional networks. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 257–266.

[201] Wenbing Huang, Tong Zhang, Yu Rong, and Junzhou Huang. 2018. Adaptive sampling towards fast graph representation learning. In *The International Conference on Neural Information Processing Systems (NeurPS)*. 4563–4572.

[202] Wenchao Yu, Cheng Zheng, and Wei Cheng. 2018. Learning deep network representations with adversarially
      regularized autoencoders. In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2663–2671.

[203] Wenwu Zhu, Xin Wang, and Peng Cui. 2020. Deep learning for learning graph representations. (2020). https://arxiv.org/abs/2001.00293v1

[204] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *The International Conference on Neural Information Processing Systems (NIPS)*. 1024–1034.

[205] Wouter Kool, Herke van Hoof, and Max Welling. 2019. Attention, learn to solve routing problems. In *The International Conference on Learning Representations (ICLR)*.

[206] Xavier Bresson and Thomas Laurent. 2017. Residual gated graph ConvNets. (2017). https://arxiv.org/abs/1711.07553

[207] Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. 2019. KGAT: knowledge graph attention
      network for recommendation. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 950–958.

[208] Xiao Huang, Qingquan Song, Yuening Li, and Xia Hu. 2019. Graph recurrent networks with attributed random walks.
      In *The ACM International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 732–740.

[209] Xiao Shen and Fulai Chung. 2020. Deep network embedding for graph representation learning in signed networks.
      *IEEE Transactions on Cybernetics* 50, 4 (2020), 1556–1568.

[210] Xiao Wang, Houye Ji, Chuan Shi, Bai Wang, Peng Cui, Philip S. Yu, and Yanfang Ye. 2018. Heterogeneous Graph
      Attention Network. In *The World Wide Web Conference*. 2022–2032.

[211] Xiao Wang, Ruijia Wang, Chuan Shi, Guojie Song, and Qingyong Li. 2020. Multi-component graph convolutional
      collaborative filtering. In *The AAAI Conference on Artificial Intelligence (AAAI)*.

[212] Xiaodan Liang, Liang Lin, Xiaohui Shen, Jiashi Feng, Shuicheng Yan, and Eric P. Xing. 2017. Interpretable structure-evolving LSTM. In *The IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*. 1010–1019.

[213] Xiaodan Liang, Xiaohui Shen, Jiashi Feng, Liang Lin, and Shuicheng Yan. 2016. Semantic Object Parsing with Graph
      LSTM. In *The European Conference on Computer Vision (ECCV)*. 125–143.

[214] Xiaojie Guo, Lingfei Wu, and Liang Zhao. 2018. Deep Graph Translation. (2018). https://arxiv.org/abs/1805.09980

[215] Xin Jiang, Kewei Cheng, Song Jiang, and Yizhou Sun. 2020. Chordal-GCN: exploiting sparsity in training large-scale graph convolutional networks. (2020).

[216] Xinyi Zhang and Lihui Chen. 2019. Capsule Graph Neural Network. In *The International Conference on Learning Representations (ICLR)*.

[217] Yang Gao, Hong Yang, Peng Zhang, Chuan Zhou, and Yue Hu. 2019. GraphNAS: graph neural architecture search with reinforcement learning. (2019). https://arxiv.org/abs/1904.09981v1

[218] Yao Ma, Suhang Wang, Charu C. Aggarwal, Dawei Yin, and Jiliang Tang. 2018. Multi-dimensional graph convolutional networks. (2018). https://arxiv.org/abs/1808.06099

[219] Yao Ma, Suhang Wang, Charu C. Aggarwal, and Jiliang Tang. 2019. Graph convolutional networks with EigenPooling. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 723–731.

[220] Yao Mao, Ziyi Guo, Zhaochun Ren, Eric Zhao, Jiliang Tang, and Dawei Yin. 2018. Streaming graph neural networks. (2018). https://arxiv.org/abs/1810.10627

[221] Yawei Luo, Tao Guan, Junqing Yu, Ping Liu, and Yi Yang. 2018. Every node counts: self-ensembling graph convolutional networks for semi-supervised learning. (2018). https://arxiv.org/abs/1809.09925v1

[222] Yedid Hoshen. 2017. VAIN: attentional multi-agent predictive modeling. In *The International Conference on Neural Information Processing Systems (NIPS)*. 2698–2708.

[223] Yelong Shen, Jianshu Chen, Po-Sen Huang, Yuqing Guo, and Jianfeng Gao. 2018. M-Walk: learning to walk over graphs using monte carlo tree search. In *The International Conference on Neural Information Processing Systems (NIPS)*. 6786–6797.

[224] Yifan Feng, Haoxuan You, Zizhao Zhang, Rongrong Ji, and Yue Gao. 2019. Hypergraph Neural Networks. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 3558–3565.

[225] Yihe Dong, Will Sawin, and Yoshua Bengio. 2020. HNHN: hypergraph networks with hyperedge neurons. (2020). https://arxiv.org/abs/2006.12278

[226] Yinwei Wei, Xiang Wang, Liqiang Nie, Xiangnan He, Richang Hong, and Tat-Seng Chua. 2019. MMGCN: multi-modal graph convolution network for personalized recommendation of micro-video. In *The ACM International Conference on Multimedia (MM)*. 1437–1445.

[227] Yizhou Zhang, Yun Xiong, Xiangnan Kong, Shanshan Li, Jinhong Mi, and Yangyong Zhu. 2018. Deep Collective Classification in Heterogeneous Information Networks. In *The World Wide Web Conference (WWW)*. 399–408.

[228] Yotam Hechtlinger, Purvasha Chakravarti, and Jining Qin. 2017. A generalization of convolutional neural networks to graph-structured data. (2017). https://arxiv.org/abs/1704.08165

[229] Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, and Xavier Bresson. 2018. Structured sequence modeling with graph convolutional recurrent networks. In *The International Conference on Neural Information Processing (ICONIP)*. 362–373.

[230] Yu Jin and Joseph F. JaJa. 2018. Learning graph-level representations with recurrent neural networks. (2018). https://arxiv.org/abs/1805.07683

[231] Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. 2020. DropEdge: towards deep graph convolutional networks on node classification. In *The International Conference on Learning Representations (ICLR)*.

[232] Yu Zhou, Jianbin Huang, Heli Sun, Yizhou Sun, Shaojie Qiao, and Stephen Wambura. 2019. Recurrent meta-structure for robust similarity measure in heterogeneous information networks. *ACM Transactions on Knowledge Discovery from Data* (2019), No. 64.

[233] Yuan Li, Xiaodan Liang, Zhiting Hu, Yinbo Chen, and Eric P. Xing. 2019. Graph Transformer. (2019). https://openreview.net/forum?id=HJei-2RcK7

[234] Yue Zhang, Qi Liu, and Linfeng Song. 2018. Sentence-State LSTM for text representation. In *The Annual Meeting of the Association for Computational Linguistics*. 317–327.

[235] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. [n.d.]. Gated graph sequence neural networks.

[236] Yujia Li, Oriol Vinyals, Chris Dyer, Razvan Pascanu, and Peter Battaglia. 2018. Learning Deep Generative Models of Graphs. In *The International Conference on Learning Representations (ICLR Workshop)*.

[237] Yujun Cai, Liuhao Ge, Jun Liu, Jianfei Cai, Tat-Jen Cham, Junsong Yuan, and Nadia Magnenat Thalmann. 2019. Exploiting spatial-temporal relationships for 3d pose estimation via graph convolutional networks. In *The IEEE International Conference on Computer Vision (ICCV)*. 2272–2281.

[238] Yuyu Zhang, Xinshi Chen, Yuan Yang, Arun Ramamurthy, Bo Li, Yuan Qi, and Le Song. 2020. Efficient probabilistic logic reasoning with graph neural networks. In *The International Conference on Learning Representations (ICLR)*.

[239] Yuzhou Chen, Yulia R. Gel, and Konstanin Avrachenkov. 2020. Fractional graph convolutional networks (FGCN) for semi-supervised learning. (2020). https://openreview.net/forum?i

[240] Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, and Yixin Chen. 2019. D-VAE: A Variational Autoencoder for Directed Acyclic Graphs. (2019). https://arxiv.org/abs/1904.11088

[241] Zhen Zhang, Hongxia Yang, Jiajun Bu, Sheng Zhou, Pinggang Yu, Jianwei Zhang, Martin Ester, and Can Wang. 2018.
      ANRL: attributed network representation learning via deep neural networks. In *The International Joint Conference on
      Artificial Intelligence (IJCAI)*. 3155–3161.
[242] Zhihong Zhang, Dongdong Chen, Jianjia Wang, Lu Bai, and Edwin R. Hancock. 2019. Quantum-based subgraph
      convolutional neural networks. *Pattern Recognition* 88, 2019 (2019), 38–49.
[243] Zhihong Zhang, Dongdong Chen, Zeli Wang, Heng Li, Lu Bai, and Edwin R. Hancock. 2019. Depth-based subgraph
      convolutional auto-encoder for network representation learning. *Pattern Recognition* 90, 2019 (2019), 363–376.
[244] Zhijiang Guo, Yan Zhang, and Wei Lu. 2019. Attention guided graph convolutional networks for relation extraction.
      In *The Annual Meeting of the Association for Computational Linguistics (ACL)*. 241–251.
[245] Zhijie Deng, Yinpeng Dong, and Jun Zhu. 2019. Batch virtual adversarial training for graph convolutional networks.
      In *The International Conference on Machine Learning (ICML Workshop)*.
[246] Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. GNNExplainer: generating
      explanations for graph neural networks. In *The International Conference on Neural Information Processing Systems
      (NeurPS)*. 9244–9255.
[247] Ziniu Hu, Yuxiao Dong, Kuansan Wang, and Yizhou Sun. 2020. Heterogeneous Graph Transformer. (2020). https:
      //arxiv.org/abs/2003.01332
[248] Ziqi Liu, Chaochao Chen, Longfei Li, Jun Zhou, Xiaolong Li, Le Song, and Yuan Qi. 2019. GeniePath: graph neural
      networks with adaptive receptive paths. In *The AAAI Conference on Artificial Intelligence (AAAI)*. 4424–4431.
[249] Ziwei Zhang, Peng Cui, and Wenwu Zhu. 2018. Deep learning on graphs: a survey. (2018). https://arxiv.org/abs/1812.
      04202
[250] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. 2019. A comprehensive
      survey on graph neural networks. (2019). https://arxiv.org/abs/1901.00596
[251] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. 2018. Adversarial attacks on neural networks for graph
      data. In *The International Conference on Knowledge Discovery and Data Mining (SIGKDD)*. 2847–2856.