# Reinforcement learning on graphs: A survey

Mingshuo Nie , Dongming Chen , and Dongqi Wang

*Abstract*—Graph mining tasks arise from many different application domains, ranging from social networks, transportation to E-commerce, etc., which have been receiving great attention from the theoretical and algorithmic design communities in recent years, and there has been some pioneering work employing the research-rich Reinforcement Learning (RL) techniques to address graph data mining tasks. However, these graph mining methods and RL models are dispersed in different research areas, which makes it hard to compare them. In this survey, we provide a comprehensive overview of RL and graph mining methods and generalize these methods to Graph Reinforcement Learning (GRL) as a unified formulation. We further discuss the applications of GRL methods across various domains and summarize the method descriptions, open-source codes, and benchmark datasets of GRL methods. Furthermore, we propose important directions and challenges to be solved in the future. As far as we know, this is the latest work on a comprehensive survey of GRL, this work provides a global view and a learning resource for scholars. In addition, we create an online open-source for both interested scholars who want to enter this rapidly developing domain and experts who would like to compare GRL methods.

*Index Terms*—Graph reinforcement learning, Graph mining, Reinforcement learning, Graph neural networks

## I. INTRODUCTION

THE recent success of reinforcement learning (RL) has solved challenges in different domains such as robotics [1] , games [2], Natural Language Processing (NLP) [3], etc. RL addresses the problem of how agents should learn to take actions to maximize cumulative reward through interactions with the environment [4]. The rapid development of RL in cross-disciplinary domains has motivated scholars to explore novel RL models to address real-world applications, e.g., financial and economic [5]–[7], biomedical [8], [9], and transportation [10], [11]. On the other hand, many real world data can be represented with graphs, and data mining for graph structure has received extensive research, such as link prediction [12], [13], node classification [14], and graph classification [15]. Various strategies have been proposed by scholars to address graph data, including novel information aggregation functions [16], graph structure pooling [17], and

Mingshuo Nie, Dongming Chen, Dongqi Wang are with the Software College, Northeastern University, Shenyang 110169, Liaoning, China. (e-mail: niemingshuo@stumail.neu.edu.cn; chendm@mail.neu.edu.cn; wangdq@swc.neu.edu.cn)

neighborhood sampling methods [12]. With the increasing graph scale and the continuous development of RL methods in recent years, scholars are focusing on combining graph mining with RL, and there is increasing interest in addressing decision problems arising in graph mining tasks with powerful RL methods [18]–[20]. The collaborative research on graph mining algorithms and RL models is gradually increasing, we show the trends of published papers on Graph Reinforcement Learning (GRL) ranging from January 2017 to April 2022 in Figure 1, and we will continue to update more research articles in our published open-source repository.
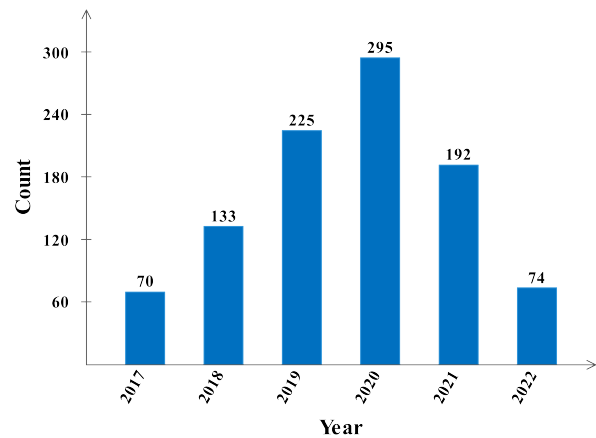


Fig. 1. The trend of GRL (January 2017-April 2022).

The traditional methods and deep learning-based models for graph mining tasks have major differences in terms of model design and training process from RL-based methods, and scholars are facing many challenges in employing RL methods to analyze graph data. (1) Machine learning models commonly used in computer vision and NLP effectively capture hidden patterns of Euclidean data (e.g., images and speech), but these models will not be available directly for graph data due to the data complexity. (2) RL methods are difficult to accurately establish the environment or state space, and design action space or reward functions for specific graph mining tasks [21]. Specifically, graphs are massive, high-dimensional, discrete objects, and are challenging to work with in a RL modeling context. (3) The specific objectives of graph mining tasks require RL methods to design effective architectures for the specific objectives that are adapted to graph data. (4) The emergence of large-scale graph data creates limitations for the RL methods to be employed. The historical data-driven RL methods need to continuously learn from the environment to update the parameters, and the increasing of the data scale causes inefficiencies in RL methods. (5) Many

interdisciplinary data in the real world can be modeled as graph structures, and interdisciplinary data analysis integrating domain knowledge will drive models toward expansion and complexity, which leads to significant difficulties in designing RL methods. We believe that the scalable and efficient RL methods to solve graph mining tasks are the key research problem.

To address the above challenges, scholars have been working extensively in the domain of RL and graph mining, and there are attempts in many fields [22], [23], including rumor detection [24], recommendation systems [25], [26], and automated machine learning (AutoML) [27]–[30]. We believe that existing methods that combine RL techniques with graph data mining for co-learning are divided into two main categories: (1) Solving RL problems by exploiting graph structures [31]–[37]. (2) Solving graph mining tasks with RL methods [38]–[40].

We define GRL as solutions and measures for solving graph mining tasks by analyzing critical components such as nodes, links, and subgraphs in graphs with RL methods to achieve exploration of topological structure and attribute information of the graphs. The success of GRL in many domains is partially attributed to: (1) **Generalization and Adaptability.** RL methods enable learning specific objectives in networks with different topologies without considering the effect of network structure on prediction performance since RL methods allow for adaptive learning by representing graph mining tasks as the MDP with sequential decision characteristics. (2) **Data-driven and Efficient.** Existing graph data mining methods introduce abundant expert knowledge or require some rules to be developed manually [18], [41], [42], while RL methods allow fast learning without expert knowledge, and these methods achieve learning objectives more efficiently.

There are a limited number of comprehensive reviews available to summarize the works using RL methods to solve graph data mining, thus, it is necessary to provide this systematic review of methods and frameworks for GRL. We provide an overview of the latest research contents and trends in the domains discussed above and summarize the extensive literature for some specific real world applications to prove the importance of GRL. We aim to provide an intuitive understanding and high-level insight into the different methods so that scholars can choose the suitable direction to explore. To the best of our knowledge, this is the first work to provide a comprehensive survey of various methods in GRL.

Our paper makes notable contributions summarized as follows:

- **Comprehensive review.** We provide a systematic and comprehensive overview of modern GRL methods for graph data, and we categorize these methods according to research domains and characteristics of RL methods, which focus on solving data mining problems on graphs with RL.
- **Recent literature.** We provide a summary of recent work about GRL and an investigation of real world applications. By doing the survey, we hope to provide a useful resource and global perspective for the graph mining and RL communities.

- **Future directions.** We discuss six possible future directions in terms of automated RL, Hierarchical RL, multi-agent RL, subgraph pattern mining, explainability, and evaluation metrics.
- **Abundant resources.** We collect abundant resources on GRL, including state-of-the-art models, benchmark datasets, and open-source codes. We create an open-source repository for GRL, which contains numerous papers on GRL in recent years. In this open-source repository, we provide links to the papers and code (optional) to allow scholars to get a faster overview of the research contents and methods, and to provide more inspiration for research on GRL.

The rest of this survey is organized as follows. Section II provides a brief introduction to the graph data mining concepts discussed in this paper and a brief review of the key algorithms commonly used in GRL. In Section III, we review some current research on the topic of GRL and list the state, action, reward, termination, RL algorithms, and evaluation metrics. In Section IV, we discuss and propose some future research directions and challenges. Finally, we summarize the paper in Section V.

## II. PRELIMINARIES

Graphs are the natural abstraction of complex correlations found in numerous domains, different types of graphs have been the main topic in many scientific disciplines such as computer science, mathematics, engineering, sociology, and economics [43]. We summarize the common graph data mining problems in GRL research as follows, and the scholars try to solve these problems with RL-based methods.

### A. Graph Neural Networks

Recently, Graph Neural Networks (GNNs) have been employed to model and compute graph data to improve the ability to reason and predict relationships in graphs. These models have a strong relational induction to effectively learn the relationships between nodes an links in graphs and establish rules for graph data, ranging from social media to biological network analysis and network modeling and optimization [44]–[46]. GNN learns the embedding of nodes and graphs by aggregating the features of target nodes and their neighboring nodes. The basic principle of the GNN model can be concluded as follows: Given a graph G with $m$ nodes, the G can be represented by the adjacency matrix $A \in \{0,1\}^{m \times m}$ and the feature matrix $X \in R^{m \times d}$, where $d$ indicates that each node corresponds to a $d$-dimensional feature vector, the information aggregation operation of GNN is defined as Equation 1.

$$X_{i+1} = \sigma \left( D^{-\frac{1}{2}} \hat{A} D^{-\frac{1}{2}} X_i W_i \right) \tag{1}$$

where $X_i$ denotes the input feature matrix of the $i$-th layer GCN. $X_0 = X$ denotes the original feature vector of the input graph, and $X_i \in R^{m \times d_i}$ is transformed into $X_{i+1} \in R^{m \times d_{i+1}}$ by the message passing mechanism. $\hat{A} = A + I$ denotes that self-loops are added to the adjacency matrix of the input graph. $D$ denotes the diagonal node degree matrix after

the normalization operation is executed on $\hat{A}$. In addition, $W_i \in R^{d_i \times d_{i+1}}$ is the learnable weight matrix and $\sigma(\cdot)$ denotes the nonlinear activation function. Motivated by Weisfeiler-Lehman (WL) graph isomorphism test, learning a GNN includes three main components [47], [48]: (1) **Initialization.** Initialize the feature vectors of each node by the attribute of the vertices, (2) **Neighborhood Detection.** Determiner the local neighborhood for each node to further gather the information and (3) **Information Aggregation.** Update the node feature vectors by aggregating and compressing feature vectors of the detected neighbors.

Many different message passing strategies have been proposed to enable graph data mining, ranging from the novel information aggregation functions [16] to graph structure pooling [17] and neighborhood sampling approaches [12]. Different GNN models provide high-level variability in the design of neighborhood detection and information aggregation, and these various models define the different node neighborhood information and the aggregation compression methods to achieve graph learning and they commonly consider that the features of the target nodes are obtained by aggregating and combining the features of their neighbors.

### B. Network Representation Learning

Models and algorithms for network representation learning are pervasive in society, and impact human behavior via social networks, search engines, and recommendation systems. Network representation learning has been recently proposed as a new learning paradigm to embed network nodes into a low-dimensional vector space, by preserving network topology structure, node content, and other side information. The low-dimensional embeddings of nodes obtained enable the application of vector classification models or vector-based machine learning models for downstream tasks, e.g., link prediction, node classification, graph clustering, etc. Graph mining methods based on network representation learning avoid the problem of missing information caused by the direct application of the original network structure. In addition, the deep network representation learning-based methods follow the message passing strategy and the graph embeddings obtained to preserve the proximity information in the low-dimensional vector space.

The network representation learning is defined as: Given a network $G = (V, E, X, Y)$, where $E$ is the set of edges in the network, $V$ is the set of nodes in the network, $X$ is the attribute feature of the nodes, and $Y$ denotes the label of the nodes, the network structure with node labels is commonly employed for node classification tasks. The network representation learning methods allow the graph to be embedded based on the mapping function $\mathcal{F}: v \rightarrow e_v \in R^d$, where $e$ is the embedding of the learned node $v$, the mapping function $f$ is employed to map the network structure into the low-dimensional vector space and preserve the topological and attribute information of the graph, to ensure that nodes with similar structure and attributes have higher proximity in the vector space. Zhang et al. [49] propose that the learned node representations should satisfy the three conditions: (1) **Low-dimensional.** The dimension of learned node representations

should be much smaller than the dimension of the original adjacency matrix representation. (2) **Informative.** The learned node representations should preserve node proximity reflected by network structure, node attributes and node labels. (3) **Continuous.** The learned node representations should have continuous real values to support subsequent network analytic tasks and have smooth decision boundaries to ensure the robustness of these tasks.

In GRL, scholars have introduced RL into the design principles of network representation learning, and have formulated the network representation learning problem as the Markov Decision Process (MDP), they have also employed RL methods in multi-relational networks and heterogeneous information networks for constructing input data for network representation learning models or optimizing the selection scheme of neighbors.

### C. Adversarial Attacks

Deep neural networks are very sensitive to adversarial attacks, which can significantly change the prediction results by slightly perturbing the input data [50]. Sun et al. [51] categorize the network adversarial attack problem as (1) **Model-Objective Attack.** The strategies of these attacks are attacking the specified model with attack methods to make the models become non-functional working in multiple scenarios, including evasion attacks and poisoning attacks. In addition, RL-based adversarial attack methods [19], [50], [52] have attracted much interest. (2) **Data-Objective Attack.** Data-Objective attacks do not attack a specific model. Such attacks happen when the attacker only has access to the data but does not have enough information about the model, including statistical information and model poisoning [53], [54].

### D. Knowledge Graphs

Knowledge graphs, which preserve rich human knowledge and facts, are commonly employed in downstream AI applications, and large-scale knowledge graphs such as DBpedia [55], Freebase [56], and Yago [57] are proved to be the infrastructure for tasks such as recommendation systems, dialogue generation [58], and Question and Answer (Q&A) tasks. The knowledge graph can be defined as $G = (h, r, t)$ to preserve the multi-relationship graph with a large a number of facts, where $h$ denotes the head entity, $t$ denotes the tail entity, and $r$ denotes the relationship between $h$ and $t$. Scholars have proposed many methods to infer the missing knowledge graphs with explicit information for knowledge graph inference and completion. Knowledge graph embedding and multi-hop path reasoning are the main methods to tackle knowledge graph completion, and knowledge graph completion methods are divided into three categories [59]: (1) **Path ranking-based Methods** [60]. Such methods are to a path to connect two entities as a feature to predict the relationship between two entities. (2) **Representation learning-based Methods** [61], [62]. Such methods aim to represent the semantic information of the research object as a dense low-dimensional real-value vector and reason via vector operations. (3) **RL-based methods** [3]. These methods define the knowledge graph reasoning task as MDP.

## E. Reinforcement Learning Methods

RL methods have recently achieved huge success in a variety of applications [63], which automatically tackles sequential decision problems in real-world environments via goal-directed learning and decision making, and such methods achieve great success in many real match games [64], [65].

The RL is formulated as an MDP, which is a sequential decision mathematical model in which actions affect the current short-term rewards, the subsequent states, and future rewards. MDP serves to simulate the random strategies and rewards of agents. In MDP, the prediction and control problem can be implemented by dynamic programming. The formal definition of MDP is the tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, p(s_0), \gamma\}$, where $\mathcal{S}$ denotes the set of all possible states, which is the generalization of the environment, $\mathcal{A}$ denotes the set of actions that can be adopted in the states, which is the all possible actions of the agent, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathcal{R}$ denotes the reward function, which is the rewards that the environment returns to the agent after executing an action, $\mathcal{T} : \mathcal{S} \times \mathcal{A} \to p(\mathcal{S})$ is identified by the state transition function, $\gamma \in [0, 1]$ denotes the discount factor, which can be treated as a hyperparameter of the agent and serves to promote the faster availability of short-term rewards to the agent. The process of interaction between agent and environment is shown in Figure 2.
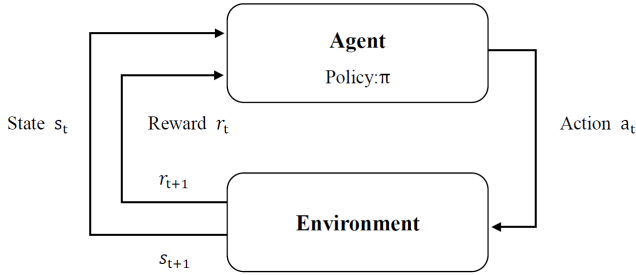


Fig. 2. The interaction process between agent and environment. The interaction process between agent and environment is described as: The first state is sampled from the initial state distribution $p(s_0)$. At timestep $t$, the agent makes a calculated decision based on the state $s_t \in \mathcal{S}$ returned from the environment and chooses the corresponding action $a_t \in \mathcal{A}$, which will be applied to the environment, resulting in the changes in the environment. The agent obtains the state representation $s_{t+1} \sim \mathcal{T}(\cdot \mid s_t, a_t)$ and a reward $r_t = \mathcal{R}(s_t, a_t, s_{t+1})$ according to the changes in the environment. The agent acts in the environment according to a policy $\pi : \mathcal{S} \to p(\mathcal{A})$. By repeatedly selecting actions and transitioning to a next state, we can sample a trace $\{s_0, a_0, r_0, s_1, a_1, r_1, ...\}$ through the environment.

Our goal is to find a policy $\pi$ that maximizes our expected return $Q(s, a)$, the target policy is defined as Equation 2 [66]:

$$\pi^* = \arg \max_{\pi} Q(s, a)$$
$$= \arg \max_{\pi} E_{\pi, \mathbb{T}}[\sum_{k=0}^{K} \gamma^k r_{t+k} \mid s_t = s, a_t = a] \quad (2)$$

RL algorithms are categorized into Model-based and Model-free algorithms according to the accessibility of the agent to the environment [67]. Model-based algorithms have to first construct the state representation that the model should predict, and there are potential problems between representation learning, model learning, and planning. In addition, model-based algorithms often suffer from the challenge that the agents cannot effectively model the real environment. Most of the GRL methods are based on the Model-free algorithms.

Deep Reinforcement Learning (DRL) is poised to revolutionize the field of artificial intelligence and represents a step toward building autonomous systems with a higher-level understanding of the visual world [68]. Deep learning enables RL to scale to decision-making problems that were previously intractable, i.e., settings with high-dimensional state and action spaces.

In GRL, it is common for scholars to employ model-free RL algorithms for graph data mining. This survey is aimed at solving graph mining tasks with RL methods. Therefore, we do not summarize all the RL methods, and instead, we focus on the current research results in the field of GRL. In this section, we present these RL methods for graph data mining.

*1) Q-learning:* Q-learning [69] is an off-policy learner and Temporal Difference (TD) learning method, which is an important result of early research on RL. The target policy in Q-learning updates the values directly on the Q-table to achieve the selection of the optimal policy, whereas, the behavior policy employs the $\varepsilon$-greedy policy for semi-random exploration of the environment. The learning goal of the action value function $Q$ to be learned in Q-learning is that the optimal action value function $q^*$ is learned by direct approximation, so that the Q-learning algorithm can be formulated as Equation 3.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha[R_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)] \quad (3)$$

where the equation denotes that state $s_t$ explores the environment with a behavior policy based on the values in the Q-table at timestep $t$, i.e., it performs action $a$, the reward $R$ and a new state $s_{t+1}$ is obtained based on the feedback from the environment. The equation is executed to update to obtain the latest the Q-table, it will continue to update the new state $s_{t+1}$ after completing the above operation until the termination time $t$.

*2) REINFORCE:* REINFORCE [70] does not optimize directly on the policy space but learns the parameterized policy without the intermediate value estimation function, which uses the Monte Carlo method for learning the policy parameters with the estimated returns and the full trace. The method creates a neural network-based policy that takes states as inputs and generates probability distributions in the operation space as outputs. The goal of the policy is to maximize the expected reward. This discounted reward is defined as the sum of all rewards obtained by the agent in the future. The policy $\pi$ is parameterized with a set of weights $\theta$ so that $\pi(s; \theta) \equiv \pi(s)$, which is the action probability distribution on the state, and REINFORCE is updated with Equation 4:

$$\Delta \omega_{i,j} = \alpha_{i,j}(r - b_{i,j}) \frac{\vartheta}{\vartheta \omega_{i,j}} \ln g_i \quad (4)$$

where $a_{i,j}$ is a non-negative learning factor, $r$ denotes the discounted reward value, and $b_{i,j}$ is a representation function of the state for reducing the variance of the gradient estimate. Williams [70] points that $b_{i,j}$ could have a profound effect on the convergence speed of the algorithm. $g_i$ is the probability density function for randomly generating unit activation-based actions.

*3) Actor - Critic:* Actor-Critic algorithm [71] combines the basic ideas of magnetic gradient and approximate dynamic programming, which employs a parameterized policy and a value function, and the value function to provide a better $\hat{A}(s, a)$ estimate for the computation of the policy gradient. The Actor-Critic algorithm studies both policy and state value functions, combining the advantages of value function-based and policy gradient-based algorithms. In this method, the Actor denotes the policy function for learning the policy that can obtain as many rewards as possible, and the Critic denotes the estimated value function for evaluating the estimated value of the current policy. Figure 3 is the framework of Actor-Critic algorithm.
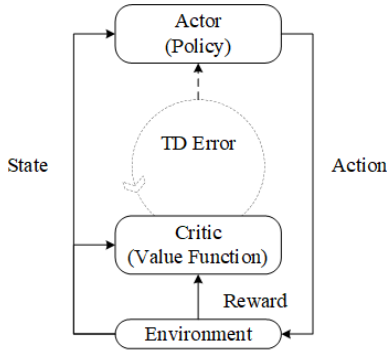


Fig. 3. The framework of Actor-Critic algorithm. Actor receives a state from the environment and selects an action to perform. Meanwhile, Critic receives the current state and the state generated by the previous interaction, and calculates the TD error to update Critic and Actor.

The Actor-Critic architecture is widely employed as the basic framework for RL algorithms Asynchronous Advantage Actor-Critic (A3C) [72], Advantage Actor-Critic (A2C) [73], Deterministic Policy Gradient (DPG) [74], Deep Deterministic Policy Gradient (DDPG) [75]. The A3C develops for single and distributed machine setups combines advantage updates with an actor-critic formulation that relies on asynchronous update policy and value function networks trained in parallel on several processing threads [68].

*4) Deep Q-network:* Q-table updates with the Q-learning algorithm in large-scale graph data suffer from the problem of high number of intermediate state values, leading to dimensional disasters. To address the above challenges, the main method proposed by scholars is the Deep Q-Network (DQN) [2] which combines value function approximation and neural networks through function approximation and state/action space reduction. DQN is widely exploited to learn policies by leveraging deep neural networks.

The scholars propose to employ Q-learning algorithms to represent states with graph embeddings according to the property of graph data in order to minimize the impact of non-Euclidean structures on the Q-table scale. Overall, DQN combines deep learning with Q-learning and approximates the action value function with deep neural networks, and obtains the trace $\{s_t, a_t, r_t, s_{t+1}\}$ from the interaction of the agent with the environment. This method can learn successful policies directly from high-dimensional sensory inputs using

end-to-end RL [2]. The loss function $L(\theta_t)$ of this method is given by Equation 5.

$$L(\theta_t) = \mathbb{E}[(r + \gamma \max_{a+1} Q_{\theta_{t-1}}(s_{t+1}, a_{t+1}) - Q_{\theta_t}(s_t, a_t))^2] \tag{5}$$

DQN introduces two techniques to stabilize the training process: (1) A replay buffer to reuse past experiences. (2) A separate target network that is periodically updated. Since the success of DQN, a large number of improved algorithms have been proposed. Double Deep Q-Network (DDQN) [76] reduces the risk of high estimation bias in Q-learning by decoupling selection and evaluation. the DDQN addresses the problem of not sufficiently considering the importance of different samples in the DQN algorithm by calculating the priority of each sample in the experience pool and increasing the probability of valuable training samples. Moreover, scholars have proposed more improved versions of DQN [77], [78] to address problems such as the lack of long-term memory capability of DQN.

## III. REINFORCEMENT LEARNING ON GRAPHS

Existing methods to solve graph data mining problems with RL methods focus on network representation learning, adversarial attacks, relational reasoning. In addition, many real-world applications study the GRL problem from different perspectives. We provide a detailed overview of the research contents in GRL in this section. Furthermore, we provide a list of abbreviations and commonly used notations in papers in Appendix A and Appendix B.

### A. Datasets & Open-source

*1) Datasets:* We provide a summary of the experimental datasets employed in these studies cited in this survey and a list of the statistical information of datasets, sources, tasks, and citations, as shown in Table I and II. In addition, we summarize the experimental results of node classification on Cora [79], CiteSeer [79], and PubMed [79] and the results of knowledge graph reasoning on NELL-995 [3], WN18RR [80], and FB15K-237 [81], where the classification accuracy is the evaluation metric for the node classification task and Hit@1, Hit@10, and MRR are the evaluation metrics for the knowledge graph reasoning, as shown in III and Table IV.

*2) Open-source:* To facilitate research into GRL, we develop an open-source repository. These papers in this repository are categorized by year and cover the fields of solving GRL problem, such as node classification, graph classification, and model explainability, which are available in public[1]. We believe that this repository could provide a comprehensive and efficient summary of methods in the field of GRL. In addition, we summarize the open-source implementations of GRL methods reviewed in the survey and provide the links of the codes of the models in Appendix C.

[1]https://github.com/neunms/Reinforcement-learning-on-graphs-A-survey

TABLE I
COMMONLY USED DATASETS FOR GRL

| Dataset | $|V|$ | $|E|$ | Source | Citation | Task |
|---|---|---|---|---|---|
| Yelp | 45954 | 3846979 | [82] | RioGNN [83], CARE-GNN [42], UNICORN [84] | Fraud Detection |
| Amazon | 11944 | 4398392 | [85] | RioGNN [83], CARE-GNN [42], SparRL [43] | Fraud Detection, Community detection |
| MIMIC-III | 28522 | 337636545 | [86] | RioGNN [83] | Diabetes Detection |
| Cora* | 2708 | 5429 | [79] | Policy-GNN [48], RL-S2V [19], GraphNAS [87], GPA [14], AGNN [88], GDPNet [89], NIPA [90], RLNet [91] | Node Classification |
| CiteSeer | 3327 | 4732 | [79] | Policy-GNN [48], RL-S2V [19], GraphNAS [87], SparRL [43], GPA [14], AGNN [88], GDPNet [89] | Node Classification |
| PubMed | 19717 | 44338 | [79] | Policy-GNN [48],NIPA [90], RL-S2V [19], GraphNAS [87], GPA [14], AGNN [88], GDPNet [89] | Node Classification |
| Twitter* | 81306 | 1768149 | [92] | SparRL [43], AdRumor-RL[24], IMGER [93] | Fraud Detection |
| Facebook | 4039 | 88234 | [92] | SparRL [43], IMGER [93], RLNet [91] | Fraud Detection |
| YouTube | 4890 | 20787 | [94] | SparRL [43], IMGER [93] | Community detection |
| Email-Eu-Core | 1005 | 16064 | [95] | SparRL [43] | Community detection |
| NELL-995* | 75492 | 237 | [3] | RF [38], AttnPath [96], RLH [97], PAAR [98], Zheng et al. [99], ADRL [100], GRL [59], MARLPaR [101], ConvE [102], MINERVA [103], RLPath [104], DAPath [105], MemoryPath [106] | KG reasoning, Link Prediction, Fact Prediction |
| WN18RR | 40493 | 11 | [80] | RF [38], RLH [97], Zheng et al. [99], ADRL [100], GRL [59], MARLPaR [101], ConvE [102], MINERVA [103], KBGAN [107] | KG reasoning |
| FB15K-237* | 14505 | 200 | [81] | RF [38], AttnPath [96], RLH [97], PAAR [98], Zheng et al. [99], ADRL [100], GRL [59], ConvE [102], MINERVA [103], KBGAN [107], RLPath [104], DAPath [105], MemoryPath [106] | KG reasoning, Link Prediction, Fact Prediction |

Note: * These datasets are used at different scales, yet they share the same data sources. We marked the maximum number of nodes and edges.

TABLE II
COMMONLY USED DATASETS FOR GRL IN GRAPH CLASSIFICATION

| Dataset | # Graphs | # Nodes(Avg.) | # Edges(Avg.) | Source | Citation |
|---|---|---|---|---|---|
| MUTAG | 188 | 17.93 | 19.79 | [108] | SUGAR [15], SubgraphX [40], XGNN [18], GraphAug [109] |
| PTC | 344 | 14.29 | 14.69 | [110] | SUGAR [15] |
| PROTEINS | 1113 | 39.06 | 72.82 | [111] | SUGAR [15], GraphAug [109] |
| D&D | 1178 | 284.32 | 715.66 | [112] | SUGAR [15] |
| NCI1 | 4110 | 29.87 | 32.30 | [113] | SUGAR [15], GraphAug [109] |
| NCI109 | 4127 | 29.68 | 32.13 | [113] | SUGAR [15], GraphAug [109] |
| BBBP | 2039 | 24.06 | 25.95 | [114] | SubgraphX [40] |
| GRAPH-SST2 | 70042 | 10.19 | 9.20 | [115] | SubgraphX [40] |

TABLE III
REPORTED EXPERIMENTAL RESULTS FOR NODE CLASSIFICATION ON THREE FREQUENTLY USED DATASETS. MISSING VALUES (-) IN THIS TABLE INDICATE THAT THIS METHOD HAS NO EXPERIMENTAL RESULTS REPORTED ON THE SPECIFIC DATASET

| Method | Cora | CiteSeer | PubMed |
|---|---|---|---|
| GCN [116]* | 81.500.0 | 70.300.0 | 79.000.0 |
| GAT [117]* | 83.00.7 | 72.50.7 | 79.00.3 |
| LGCN [118]* | 83.30.5 | 73.00.6 | 79.5 0.2 |
| DGI [119]* | 82.30.6 | 71.80.7 | 76.80.6 |
| Policy-GNN [48] | 91.91.4 | 89.72.1 | 92.12.2 |
| GraphNAS [87] | - | 73.10.9 | 79.60.4 |
| AGNN [88] | 83.6 0.3 | 73.8 0.7 | 79.7 0.4 |

Note: * These methods are non-GRL method.

### B. Graph Mining with Reinforcement Learning

*1) Network representation learning:* Network representation learning is to learn a mapping that embeds the nodes of a graph as low-dimensional vectors, while encoding a variety of structural and semantic information. These methods aim to optimize the representations so that geometric relationships in the embedding space preserve the structure of the original graph [122], [123]. The node representations obtained could effectively support extensive tasks such as node classification, node clustering, link prediction and graph classification [124],

[125]. Existing network representation learning methods commonly suffer from three challenges [15]: (1) **Low feature discrimination**. Fusing all features and relationships to obtain an overview network representation always brings the potential concern of over-smoothing, which leads to indistinguishable features of the graph. (2) **Demand for the prior knowledge**. Preservation of structural features in the form of similarity measures or motifs is always based on heuristics and requires substantial prior knowledge. (3) **Low explainability**. Many methods exploit substructures by step-by-step pooling, which loses much of the detailed structural information and leads to a lack of sufficient explainability for downstream tasks. To address the above challenges, SUGAR [15] retains structural information from the "Nodes-Subgraphs-Graphs" levels by adaptively selecting significant subgraphs with Q-learning to represent the discriminative information of the graph, and this hierarchical learning approach provides powerful representations, generalization, and explainability. The SUGAR architecture is shown in the Figure 4.

GNNs receive increasing interest due to their effectiveness in network representation learning. Generally, the message passing architecture predefined by scholars is probably not applicable to all nodes in graphs, and the multi-hop message passing mechanism is poor for passing important information
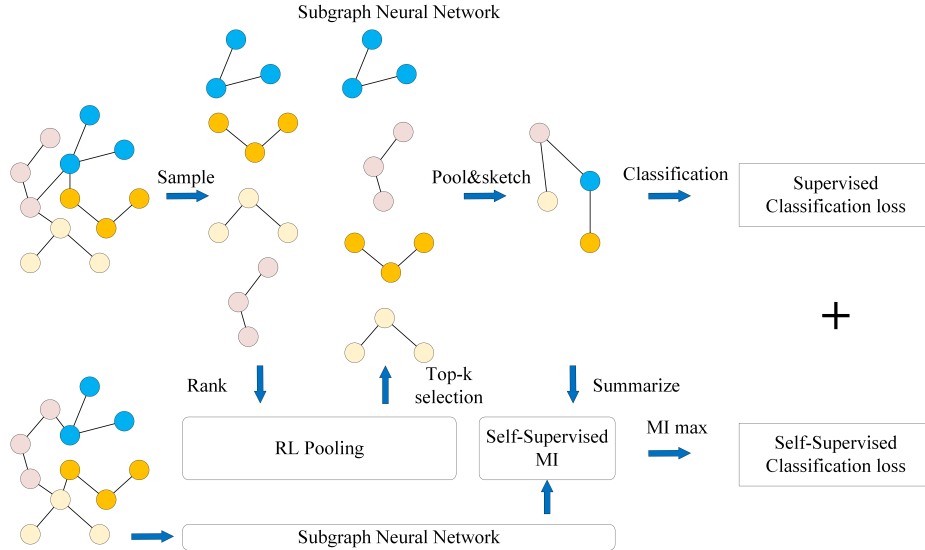
Fig. 4. An illustration of the SUGAR architecture [15]. This method consists of Subgraph Neural Network, Reinforcement Pooling Module, and Self-Supervised Mutual Information Module. (Redrawn from [15])

TABLE IV
REPORTED EXPERIMENTAL RESULTS FOR KG REASONING ON THREE FREQUENTLY USED DATASETS.

| Method | Metrics | NELL-995 | WN18RR | FB15K-237 |
|---|---|---|---|---|
| TransE [61] * | Hit@1 | 0.608 | 0.491 | 0.392 |
| | Hit@10 | 0.793 | 0.626 | 0.614 |
| | MRR | 0.715 | 0.557 | 0.494 |
| TransR [120] * | Hit@1 | 0.631 | 0.526 | 0.405 |
| | Hit@10 | 0.806 | 0.641 | 0.634 |
| | MRR | 0.727 | 0.581 | 0.532 |
| ComplEX [121] * | Hit@1 | 0.614 | 0.319 | 0.318 |
| | Hit@10 | 0.815 | 0.462 | 0.542 |
| | MRR | 0.652 | 0.428 | 0.374 |
| RLH [97] | Hit@1 | 0.692 | 0.453 | 0.342 |
| | Hit@10 | 0.873 | 0.516 | 0.648 |
| | MRR | 0.723 | 0.481 | 0.460 |
| Zheng et al. [99] | Hit@1 | 0.655 | 0.438 | 0.332 |
| | Hit@10 | 0.847 | 0.539 | 0.591 |
| | MRR | 0.731 | 0.473 | 0.422 |
| ADRL [100] | Hit@1 | 0.808 | 0.683 | 0.574 |
| | Hit@10 | 0.975 | 0.723 | 0.796 |
| | MRR | 0.916 | 0.704 | 0.712 |
| ConvE [102] | Hit@1 | 0.672 | 0.402 | 0.313 |
| | Hit@10 | 0.864 | 0.519 | 0.601 |
| | MRR | 0.747 | 0.438 | 0.410 |
| MINERVA [103] | Hit@1 | 0.663 | 0.413 | 0.217 |
| | Hit@10 | 0.831 | 0.513 | 0.456 |
| | MRR | 0.725 | 0.448 | 0.293 |

Note: * These methods are non-GRL method.

to other modules. Therefore, the novel and effective node sampling strategies and innovation of new message passing mechanisms for graph data attract more attention from scholars [16], [37], [91], [126], [127]. Batch sampling and importance sampling are commonly employed in GNNs to obtain the local structure of graph data, however, such sampling methods may cause information loss. Lai et al. [48] defined a meta-policy module and a GNN module for learning relationships and network representations between node attributes and aggregation iterations. The defined MDP samples the number of hops of the current node and neighboring nodes iteratively with meta-policy and trains the GNN by aggregating the node information within a specified within the sampled hops. Policy-GNN algorithm solves the challenge of determining the aggregation range of nodes in large-scale networks with DQN algorithm [2]. Hong et al. [37] focused on the modular RL method and designed a transformer model without morphological information to efficiently encode relationships in graphs, thus solving the challenge of the multi-hop information transfer mechanism. In addition to the parameter sharing and Buffer mechanism for network representation learning in Policy-GNN that enable the improvement of model efficiency, Yan et al. [21] proposed a novel and efficient algorithm for automatic virtual network embedding that combines the A3C [72] with GCN, which employs the GCN to automatically extract spatial features in an irregular graph topology in learning agent. Zhao et al. [128] propose that multi-view graph representation learning is implemented by reinforcing inter- and intra-graph aggregation, and propagation updates are performed with RL methods to determine the optimal filtering threshold. Furthermore, the hierarchical structure of the network is crucial for network representation learning, ACE-HGNN [129] leverages multi-agent RL method to learn the optimal curvature of the network to improve model quality and generalizability.

Existing representation learning methods based on GNNs and their variants depend on the aggregation of neighborhood information, which makes them sensitive to noise in the graph, scholars have improved the performance of network representation learning from the perspective of removing graph noisy data. GDPNet [89] employs two phases of signal neighborhood selection and representation learning to remove noisy neighborhoods as a MDP to optimize the neighborhood of each target node, and learns a policy with task-specific rewards received from the representation learning phase, allowing the model to perform graph denoising just with weak supervision from the task-specific reward signals. GAM [130], which also focuses on noisy graph data, enables to limit the learning

attention of the model to small but informative parts of the graph with attention-guided walks, The graph attention model employs the Partially Observable Markov Decision Process (POMDP) method for training to selectively process informative portions of the graph. Moreover, NetRL [131] performs network enhancement with detecting noisy links and predicting missing links to improve network analysis and modeling capabilities.

The ability of GNN models for representation learning will change significantly with slight modifications of the architecture. The design principles of the architecture require substantial domain knowledge to guide, e.g., the manual setting of hidden dimensions, aggregation functions, and parameters of the target classifier. GraphNAS [87] designs a search space covering sampling functions, aggregation functions and gated functions and searches the graph neural architectures with RL. The algorithm first uses a recurrent network to generate variable-length strings that describe the architectures of GNNs. AGNN [88], which has the same objective as GraphNAS, verifies the architecture in the predefined search space via small steps with RL-based controllers, decomposing the search space into the following six classes of actions: Hidden dimension, Attention function, Attention head, Aggregate function, Combine function, and Activation function. Furthermore, GQNAS [132] captures the structural correlations within the network layers to solve the neural architecture search with DQN and GNN.

However, the above algorithms fail to consider heterogeneous neighborhoods in the aggregation of nodes [83], causing neglect or simplification of the diverse relationship of nodes and edges in real networks, which brings a loss of heterogeneity information. Heterogeneous information networks contain multiple types of nodes and relationship, and maintain abundant and complex interdependencies between nodes, and are widely observed in real networks and practical applications [133]–[135]. Relational GNNs based on artificial meta-paths [136] or meta-graphs [137] rely on inherent entity relationships and require the support of substantial domain knowledge. Zhong et al. [138] solve the dependence on the handcrafted meta-paths via proposing a RL enhanced heterogeneous GNN model to design different meta-paths for nodes in heterogeneous information networks, and replacing the manual design of meta-paths with agent. This method discovers many meaningful meta-paths that are ignored by human knowledge. In addition to improving the design solution of meta-paths, the adoption of novel neighborhood aggregation techniques enables to effectively improve the performance of heterogeneous network representation learning. Peng et al. [83] propose a task-driven GNN framework based on multi-relational graphs that learns multi-relational node representations with the semi-supervised method, which leverages relational sampling, message passing, metric learning, and RL to guide neighbor selection within and between different relations. FinEvent [139] for social network modeling tasks allows transferring cross-lingual social event detection through modeling social messages into heterogeneous graphs. These above works further illustrate the massive role of RL methods in multi-relational networks.

The graph data augmentations provide further improvements in terms of the performance of network representation learning, and scholars believe that invariance of the learning mechanism is critical to the data augmentations. GraphAug [109] enables to avoid compromising the critical label-related information of the graph on the basis of the label-invariance of the computed graph using an automated augmentation model, thus producing label invariance at most time. They proposed a RL-based training method to maximize the estimated label-invariant probability.

GPA [14] which is employed to improve the learning performance of graph representations studies how to efficiently label the nodes in GNNs thereby reducing the annotation cost of training GNNs using an active learning method. The GPA architecture is shown in the Figure 5.

More GRL algorithms on network representation learning could be found in the Table V.

*2) Adversarial Attacks:* Following recent research showing that GNNs are trained based on node attributes and link structures in the graph, the attacker can attack the GNN by modifying the graph data used for training, thus affecting the performance of node classification [19], [140], [141]. Existing researches of adversarial attacks on GNNs focuses on modifying the connectivity between existing nodes, and the attackers inject hostile nodes with fake links into the original graph to degrade the performance of GNN in classification of existing nodes. RL-S2V [19] learns how to modify the graph structure by sequentially adding or removing links to the graph only using the feedback from downstream tasks. Q-learning algorithm is used as the implementation solution of a RL module for sequential modification of network structure. The algorithm trains the original graph structure with structure2vec to obtain the representation of each node, and then parameterizes each node with a graph neural network to obtain the Q-value, and the candidate links that need to modify are obtained with Q-learning to achieve a black-box attack. The RL-S2V architecture is shown in the Figure 6.

However, such attack policies that require frequent modifications to the network structure need a high level of complexity to avoid the attack being caught. Sun et al. [90] propose NIPA method to poison the graphs to increase the node misclassification rate of GNNs without changing the link structure between the existing nodes in the graph. Instead of manipulating the links between existing nodes, NIPA method injects fake nodes into the graph. This algorithm represents the adversarial connections and adversarial labels of the injected false nodes as MDP and designs an appropriate reward function to guide the RL agent to reduce the node classification performance of GNNs. The graph rewiring method [52] also avoids adding or removing edges resulting in significant changes to the graph structure.

Adversarial attack researches on graph neural networks also includes fraud entity detection tasks, such as opinion fraud and financial fraud. However, previous work has paid little attention to the camouflage behavior of fraudsters, which could hamper the performance of GNN-based fraud detectors during the aggregation process. These methods either fail to fit the fraud detection problems or break the end-to-end
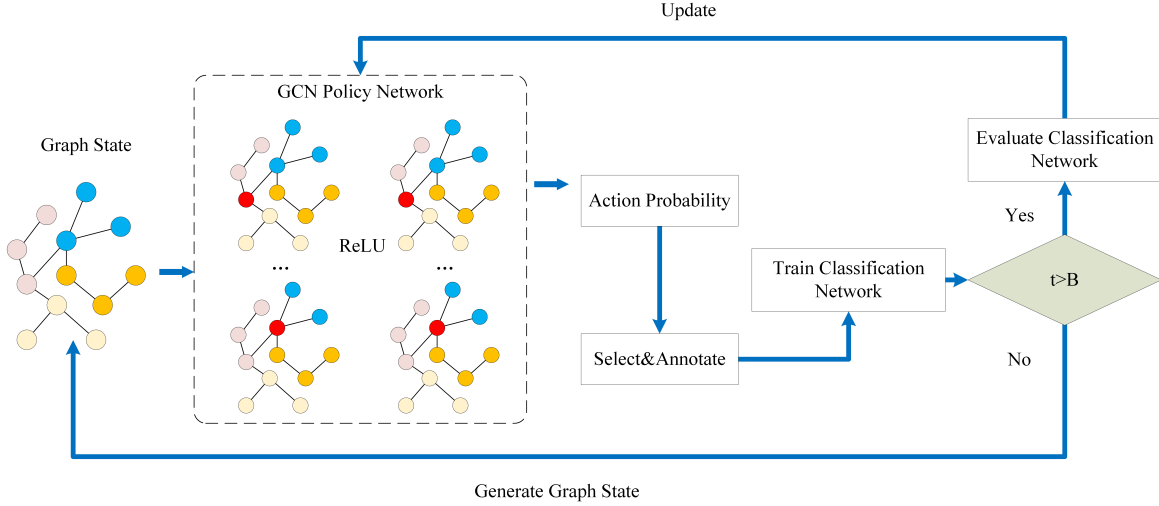
Fig. 5. An illustration of the GPA architecture [14]. In GCN Policy Network, each column denotes a layer of GNN, and the graphs in each column correspond to the feature aggregation on different nodes. (Redrawn from [14])

TABLE V
GRL ALGORITHMS ON NETWORK REPRESENTATION LEARNING TASKS. MISSING VALUES (-) IN THIS TABLE INDICATE THAT THE PERSONALIZED
TERMINATION CONDITIONS CONTRIBUTE TO THE STABILIZATION OF THE TRAINING PROCESS.

| Method | Action | State | Reward | Termination | RL | Metrics |
|---|---|---|---|---|---|---|
| SUGAR [15] | Add or minus a fixed value $\Delta k \in [0, 1]$ from $k$. | The middle graph consists of the subgraphs selected in each training round. | Define the corresponding reward value based on the classification accuracy of the previous round. | The change of $k$ among ten consecutive epochsis no more than $\Delta k$. | Q-learning | Average accuracy, Standard deviation, Training process, Testing performance, and Result visualization. |
| RL-HGNN [138] | All nodes involved in the meta-path. | The set of available relation types. | The performance improvement on the specific task comparing with the history performances. | - | DQN | Micro-F1, Macro-F1, and Time-consuming for Meta-path Design. |
| Policy-GNN [48] | The attribute of the current node. | The number of hops of the current node. | The performance improvement on the specific task com-paring with the last state. | - | DQN | Accuracy |
| RioGNN [83] | All actions when the relation $r$ is at the depth $d$ of the $l$-th layer. | The average node distance of each epoch. | The similarity as the decisive factor within the reward function. | The RL will be terminated as long as the same action appears three times in a row at the current accuracy | MDP | Accuracy |
| GraphNAS [87] | Select the specified parameters for the graph neural network architecture. | Graph neural network architecture. | Define the corresponding reward value based on the special task accuracy. | To maximize the Expected accuracy of the generated architectures. | MDP | Micro-F1 and Accuracy |
| GPA [14] | The actions are defined on the basis of the action probabilities given by the policy network. | State representation of nodes are defined on the basis of several commonly used heuristics criteria in active learning. | The classification GNNs performance score on the validation set after convergence. | Performance convergence of classifiers. | MDP | Micro-F1 and Macro-F1 |
| GDPNet [89] | Select neighbors. | Embedding of information about the current node and its neighbors. | Define the corresponding reward value based on the special task accuracy. | - | MDP | Accuracy |

learning fashion of GNNs [42]. Dou et al. [42] propose a label-aware similarity metric and a RL-based similarity-aware neighbor selector. They formulate the RL process as a Bernoulli Multi-armed Bandit (BMAB) between the neighbor selector and the GNN with similarity metric to achieve the selection threshold for adaptively finding the best neighbor when the GNN is training, and formulate the relational-aware neighborhood aggregation method to obtain the final central node representation. Lyu et al. [24] consider that the explainability of the model is as important as its effectiveness, and their proposed AdRumor-RT model enables the generation of interpretable and effective evasive attack against a GCN-based rumor detector, where the Actor-Critic algorithm is employed to implement black-box attacks.

More GRL algorithms on adversarial attacks could be found in the Table VI.

*3) Relational Reasoning:* Discovering and understanding causal mechanisms could be defined as searching for DAGs that minimize the defined score functions. RL methods have achieved excellent results in terms of causal discovery from observed data, but searching the Directed Acyclic Graphs (DAGs) space or discovering the implied conditions usually has a high complexity. The agent in RL with random policies could automatically define the search policy based on the learn the uncertain information of the policy, which can be updated rapidly with the reward signal. Therefore, Zhu et al. [142]
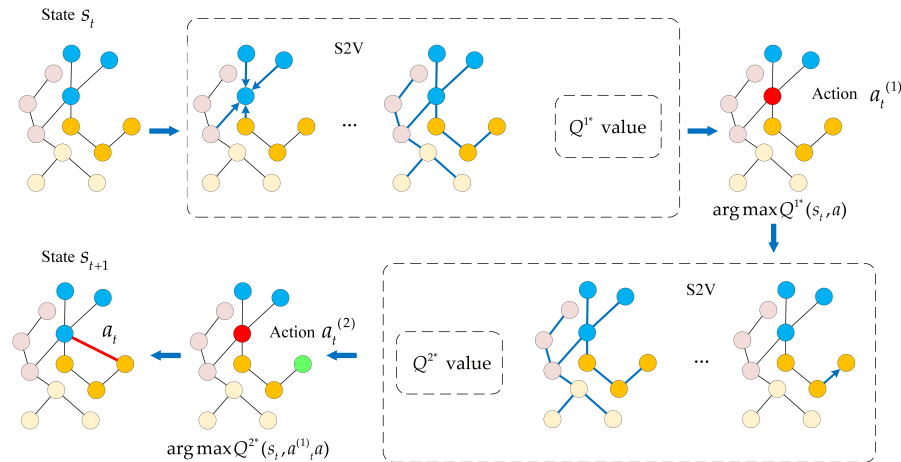
Fig. 6. An illustration of the RL-S2V architecture [19]. The algorithm performs network representation learning via S2V for calculating Q-value, adding a single edge $a_t$ is decomposed into two decision steps $a^{(1)}$ and $a^{(2)}$, with $Q^{1*}$ and $Q^{2*}$. (Redrawn from [19])

TABLE VI
GRL ALGORITHMS ON ADVERSARIAL ATTACKS TASKS

| Method | Action | State | Reward | Termination | RL | Metrics |
|---|---|---|---|---|---|---|
| NIPA [90] | Add the adversarial edges within the injected nodes between the injected nodes and the clean node and designing the adversarial labels of the injected nodes. | The intermediate poisoned graph and labels information of the injected nodes. | Guiding reward based on classifier success rate. | The agent adds budget number of edges. | DQN | Accuracy, Key statistics of the poisoned graphs, Average Degrees of Injected Nodes, and Sparsity of the Origin Graph. |
| RL-S2V [19] | Add or delete edges in the graph. | The modified graph. | The prediction confidence of the target classifier. | The agent modifies the specified number of links. | Q-learning | Accuracy |
| CARE-GNN [42] | Plus or minus a fixed small value. | The selection range of neighbor nodes. | The average distance differences between two consecutive epochs. | The RL converges in the recent ten epochs and indicates an optimal threshold until the convergence of GNN. | BMAB | ROC-AUC and Recall |
| ReWatt [52] | The action space consists of all the valid rewiring operations. | The state space of the environment consists of all the intermediate graphs generated after all the possible rewiring operations. | The reward function is developed based on the performance of the classifier. | The attack process will stop either when the number of actions reaches the budget $K$ or the attacker successfully changed the label of the slightly modified graph. | MDP | Attacking performance |

propose to leverage RL to find the underlying DAG from the graph space without the need of smooth score functions. This algorithm employs Actor-Critic as the search algorithm and outputs the graph that achieves the best reward among all the graphs generated during training. However, the problem of poor computational scores emerges in this method, and the action space composed of directed graphs is commonly difficult to be explored completely. Wang et al. [143] propose CORL method via incorporating RL methods into the ordering-based paradigm. The method describes the ordering search problem as a multi-step MDP and implements the ordering generation process with encoder-decoder structures, and finally optimizes the proposed model with RL based on the reward mechanism designed for each ordering. The generated ordering is then processed using variable selection to obtain the final causal graph. Furthermore, Sun et al. [144] combine transfer learning and RL for co-learning to leverage the prior causal knowledge for solving causal reasoning tasks.

Task-oriented Spoken Dialogue System (SDS) is a system that can continuously interact with a human to accomplish a predefined task. Chen et al. [145] design an alternative but complementary method to innovate the structure of neural networks incorporating the DQN algorithm in order to make them more adaptable to dialogue policy. The scholars have proposed some natural question generation models to provide more training data for the Q&A tasks in order to further improve the performance of the task. Chen et al. [146] focuses on the natural question generation and propose a RL-based Graph-to-Sequence (Graph2Seq) model, and the algorithm employs the self-critical sequence training (SCST) algorithm [147] to directly optimize the evaluation metric. Explicitly obtaining the preferences for recommended items and attributes of user through interactive conversations is the goal of conversational recommender systems. Deng et al. [84] leverage a graph structure to integrate recommendation and conversation components as a whole. This algorithm exploits a dynamic weighted graph to model the changing interrelationships between users, items and attributes during the conversation, and considers a graph-based MDP environment for simultaneously processing the relationships.

With the development of artificial intelligence, the knowledge graph has become the data infrastructure for many

downstream real-world applications and has attained a variety of applications in dialogue systems and knowledge reasoning [38], [100], [101], [105]. Lin et al. [102] propose a policy-based agent to extend its reasoning paths sequentially through a RL approach on the multi-hop reasoning task. MINERVA [103] leverages RL method to train an end-to-end model for the practical task of answering questions on multi-hop knowledge graph. However, these methods focusing on fixed multi-hop or single-hop reasoning consume substantial computational resources, and the incompleteness of hand-collected data affects the performance of the reasoning. Liu et al. [38] build an end-to-end dynamic knowledge graph reasoning framework with dynamic rewards, and this method integrates the embedding of actions and search histories as a policy network into a feedforward neural network for dynamic path reasoning. Sun et al. [148] propose a temporal path-based RL model to solve the temporal knowledge graph reasoning task and design a temporal-based reward function to guide the learning of the model. AttnPath [96] is a path-based framework that solves the lack of memory modules and over-reliance on pre-training of algorithms in knowledge graph reasoning, and this algorithm combines Long Short-Term Memory (LSTM) and Graph Attention Network (GAT) to design a RL mechanism capable of forcing the agent to move. RLPath [104] considers both relation choosing and entity choosing in the relational path search.

To address the problem of incomplete knowledge graphs, scholars commonly employ multi-hop reasoning to infer the missing knowledge. Motivated by the hierarchical structure commonly employed by humans when dealing with fuzzy scenarios in cognition, Wan et al. [97] suggest a hierarchical RL method to simulate human thinking patterns, where the whole reasoning process is decomposed into two steps of RL policies. In addition, the policy gradient approach [149] and REINFORCE [70] are employed to select two policies for encoding historical information and learning a structured action space. DeepPath [3] for hierarchical reasoning of knowledge graphs based on RL requires manual rules to process for obtaining more adequate hierarchical relationships. On the other hand, PAAR [98] employs a multi-hop reasoning model based on hyperbolic knowledge graph embedding and RL for the hierarchical reasoning. Hierarchical information of the knowledge graph plays an important role for downstream tasks, and another RL method [99] employing hierarchical information is to reason about the knowledge graph from a methodological perspective.

Recommendation systems are critical to various online applications such as E-commerce websites and social media platforms through providing the better item or information to users [150]. The interactive recommender system receives substantial attention as its flexible recommendation policy and optimal long-term user experience, and scholars have introduced DRL models such as DQN [151] and DDPG [75] into the IRS for decision-making and long-term planning in dynamic environments. KGQR [25] enables to incorporate graph learning and sequential decision problems in interactive recommender systems as a whole, to select candidate objects and learn user preferences from user feedback with prior
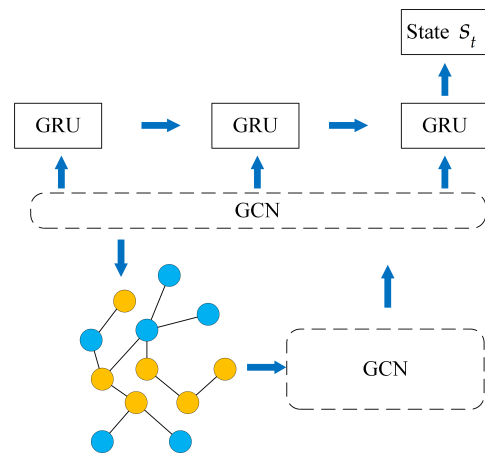


Fig. 7. An illustration of the KGQR architecture [25]. The knowledge-enhanced state representation module preserves user preferences with recurrent neural network and graph neural network. (Redrawn from [25], and we skip the candidate selection module and Q-value network in KGQR.)

knowledge, and improves the performance of RL through aggregating the semantic correlations between entities in the knowledge graph. It is capable of obtaining better recommendation performance with fewer user-item interactions. The KGQR architecture is shown in the Figure 7.

More GRL algorithms on relational reasoning could be found in the Table VII.

### C. Real World Applications with Reinforcement Learning on Graphs

The research on GRL is booming in the past few years, and GRL methods play an important role in solving various real-world applications and attract substantial attention from scholars. We summarize abundant real-world applications in transportation network optimization, recommendation systems in E-commerce networks, drug structure prediction and molecular structure generation in medical research, and network diffusion models for controlling the COVID-19 virus to further prove the hot development in the field of GRL. More GRL algorithms on real world applications could be found in the Table VIII.

*1) Explainability:* In exploring GNNs, scholars often treat them as black boxes, and this assumption lacks explanations that are readily understood by humans, which prevents their use in critical applications involving medicine, privacy, and security [115], [153]. A great number of scholars working on deciphering the explainability of such black-box model of GNNs commonly focus on explaining the explainability of the node, edge, or node feature-level in the graph data, but ignore the GNN model and the frequent subgraph. The interpretation at the subgraph-level enables more intuitive and effective description of the GNN [154]–[156]. Yuan et al. [115] categorize existing methods for explainability of GNNs into instance-level and model-level methods to advance the understanding of GNNs. In the study of GNN explainability based on subgraph-level, SubgraphX [40] propose to employ the Monte Carlo tree search method [157] to explore the critical subgraphs and

TABLE VII
GRL ALGORITHMS ON RELATIONAL REASONING TASKS. MISSING VALUES (-) IN THIS TABLE INDICATE THAT THE PERSONALIZED TERMINATION CONDITIONS CONTRIBUTE TO THE STABILIZATION OF THE TRAINING PROCESS.

| Method | Action | State | Reward | Termination | RL | Metrics |
|---|---|---|---|---|---|---|
| AttnPath [96] | The agent choosing a relation path to step forward. | The state space is consisted of the embedding part, the LSTM part and the graph attention part. | The reward is a feedback to the agent according to whether the action is valid, and whether a series of actions can lead to ground truth tail entities in a specified number of times. | - | REINFORCE | Mean Selection Rate (MSR) and Mean Replacement Rate (MRR) |
| UNICORN [84] | The actions can be selected from the candidate item set to recommend items or from the candidate attribute set to ask attributes. | All the given information for conversational recommendation, including the previous conversation history and the full graph. | The reward is defined as a weighted sum of five specific reward functions. | - | DQN | Success rate at the turn $t$ (SR$t$), Average turn (AT), Two-level hierarchical version of normalized discounted cumulative gain (hNDCG@(T,$K$)) |
| TITer [148] | The set of optional actions is sampled from the set of outgoing edges. | The state space is represented by a special quintuple. | The dynamic reward function is defined based on the search results. | - | MDP | Mean Reciprocal Rank (MRR), Hits@1/3/10 |
| RLH [97] | The set of outgoing edges of the current entity. | The tuple of entities and the relationships between them. | The reward for each step is defined based on whether the agent reaches the target entity. | - | REINFORCE | Mean average precision (MAP) and MRR |
| CORL [143] | The action space consisting of all the variables at each decision step. | The state space is defined as the embedding of the input data pre-processed with the encoder module. | The reward function consisting of episodic and dense rewards. | - | MDP | True Positive Rate (TPR), Structural Hamming Distance (SHD) |
| Zheng et al. [99] | The action space consisting of all available actions that the agent can take. | The state space consisting of the current entity, start entity, and query relation. | The reward is defined as the results of the scoring function of the knowledge representation model. | - | REINFORCE | MRR and Hits@1/10 |
| IMUP [152] | The actions are defined as the visit event. | The state is to describe the environment composed by users and the spatial knowledge graph. | The reward is defined as the weighted sum of three parts about Points of Information (POI). | - | DQN | Precision on Category, Recall on Category, Average Similarity, and Average Distance |
| ADRL [100] | The set of outgoing edges | The state space consisting of the source entity, the query relationship and the entity that is accessed. | The rewards that depend on the value function. | - | A3C | MRR and Hits@1/10 |
| GRL [59] | The agent can select its neighboring relational path to another neighboring entity at each state. | The vector representation of entities and relationships obtained with GNN. | The Adversarial rewards are defined by employing Wasserstein-GAN. | The agent reaches the target entity | DDPG | MRR, Hits@1//10, MAP, and MAE |

thus explain the prediction problem of GNNs from subgraph-level, the SubgraphX architecture is shown in the Figure 8. RioGNN [83] learns the difference in importance between different relations to obtain discriminative node embeddings, and this measure of differentiating the vector representation of nodes can enhance the quality of node embeddings while improving the explainability of the model and achieving the explainability of multi-relational GNNs from the perspective of the individual importance of different relations. Lyu et al. [24] design subgraph and node-level features to support the understanding of attack policies and detector vulnerabilities. Bacciu et al. [158] obtain perturbed policies by optimizing multi-objective scores to achieve local explanation on graph data.

In addition, the scholars believe that explaining GNNs from the model level could enhance human trust for some application domains. XGNN [18] explains GNNs by formulating graph generation as RL, and the generated graph structures enable verification, understanding, and improvement of trained GNN models.

*2) City Services:* With the rapid development of modern cities, the rapid growth of the urban population has caused many problems, such as traffic congestion and inefficient communication [159]. The analysis of problems arising in transportation networks [160]–[162] and communication networks [20], [31] with complex network methods has attracted the attention of numerous scholars. Modeling road networks and communication networks as graph data retain richer information since their structured property and assist the relevant government institutions or telecommunication service providers to better optimize traffic roads or communication links, thus alleviating the traffic congestion problems and providing better communication services to citizens.

Traffic flow prediction focus on building some prediction models to estimate the traffic flow of specific roads based on the statistical information of traffic flow within or between these regions and traffic data provided by relevant government institutions or organizations. Peng et al. [22] address the problem of incompleteness and non-temporality of most traffic flow data by modeling existing traffic flow graphs with GCPN [163] and extracting temporal features with LSTM [164]. IG-RL [165] for dynamic transportation networks focuses on the traffic signal control problem, which uses GCN to learn the entities as node embeddings and the Q-learning algorithm is employed to train the model. This method enables the representation and utilization of traffic demand and road network structure in an appropriate way regardless of the number of entities and their location in the road network.

Electronic Toll Collection (ETC) system serves an important role in alleviating the urban traffic congestion problem. Qiu et al. [10] employ a GCN to represent the value function and policy function in the Dynamic Electronic Toll Collection
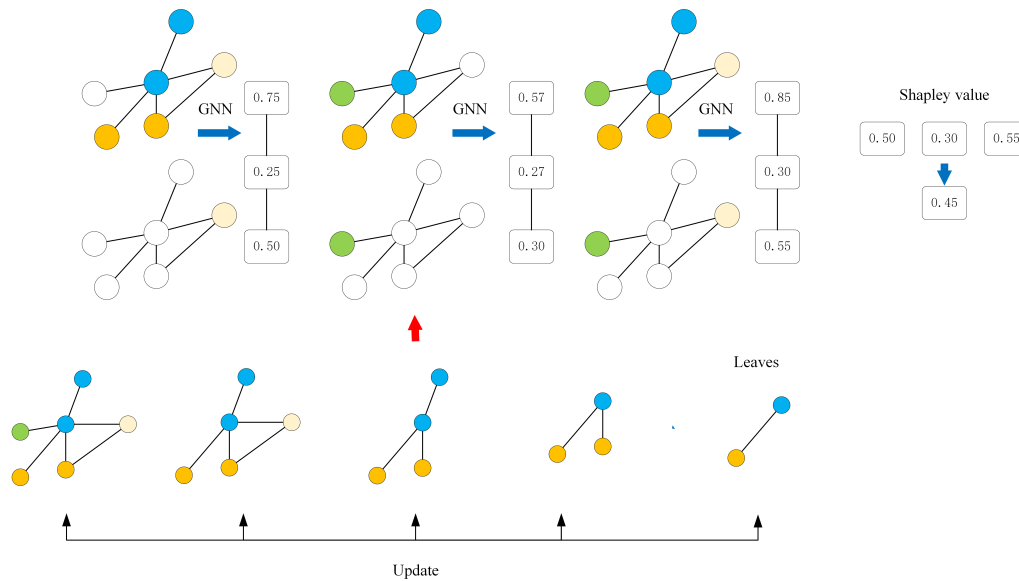
Fig. 8. An illustration of the SubgraphX architecture [40]. The algorithm selects a specified path from the search tree in root-to-leaf order corresponding to one iteration of MCTS. For each node, the importance of the subgraph is evaluated by calculating the Shapley value through Monte Carlo sampling. (Redrawn from [40])

(DETC) problem and solve the DETC with collaborative multi-agent RL algorithm. Moreover, optimization of on-demand ride-sharing services [166], lane change decision for connected autonomous vehicles [167], [168], and adaptive traffic signal control [169], [170] could also improve the operational efficiency of traffic flow in cities.

The scholars suggest a graph convolution model of multi-intelligence collaboration [41] obtain efficient collaboration policies for improving the routing in packet switching networks. Almasan et al. [171] propose to integrate GNNs with the DQN. The proposed DRL+GNN architecture enables learning and generalize over arbitrary network topologies, and is evaluated in SDN-based optical transport network scenario. For Wireless Local Area Networks (WLANs), Nakashima et al. [172] propose a DRL-based channel allocation scheme. This algorithm employs a graph convolutional layer to extract features of channel vectors with topological information and learns DDQN [76] for channel allocation policy formulation.

*3) Epidemic Control:* In the fields of epidemic containment, product marketing, and fake news detection, scholars have employed a limited number of interventions to control the observed dynamic processes partially on the graph. Yang et al. [173] propose a full-scale diffusion prediction model to integrate information from macro and micro for solving information diffusion prediction. Meirom et al. [23] propose a RL method controlling a partially-observed dynamic process on a graph by a limited number of interventions, which is successfully applied to curb the spread of an epidemic. In another epidemic control application, the RAI algorithm [174] leverage social relationships between mobile devices in the Social Internet of Things (SIoT) to assist in controlling the spread of the virus by allocating limited protection resources to influential individuals through early identification of suspected COVID-19 cases. Furthermore, some algorithms for key node finding [175], [176] and network dismantling [177] enable the

application to the epidemic control problem. To date, there has been little research on the use of GRL methods for network dynamics, and the existing methods have focused only on problems such as epidemic disease control. We believe that the network dynamics model based on GRL deserves further research and investigation.

*4) Combinatorial Optimization:* Combinatorial optimization is an interdisciplinary problem, which spans optimization, operational research, discrete mathematics, and computer science, covering abundant classical algorithms and many critical real-world applications [178]. Most real-world combinatorial optimization problems could be represented by graphs. GRL methods for combinatorial optimization allow learning automatically strategies that return feasible and high-quality outputs based on the original input data and the guidance process [179]. OpenGraphGym [180] is proposed to solve the combinatorial graph optimization problems with computed graph solutions. S2V-DQN [181] learn a policy for building solutions incrementally with Q-learning. in which the agent incrementally constructs efficient solutions based on the graph structure through adding nodes, and this greedy algorithm is a popular pattern for approximating algorithms and heuristics for combinatorial optimization, and the S2V-DQN architecture is shown in the Figure 9. Toyer et al. [182] propose an Action Schema Network for learning generalized policies for probabilistic planning problems. The AS-Net can employ a weight sharing scheme through simulating the relational structure of the planning problem, allowing the network to be applied to any problem in a given planning domain.

The dynamic version of many graphs data mining is a critical object of study for solving traffic, social, and telecommunication networks where the structure of networks in the real world evolves over time. The GTA-RL algorithm [39] is a graph-based heuristic learning algorithm for dynamic combinatorial optimization problems, which addresses the
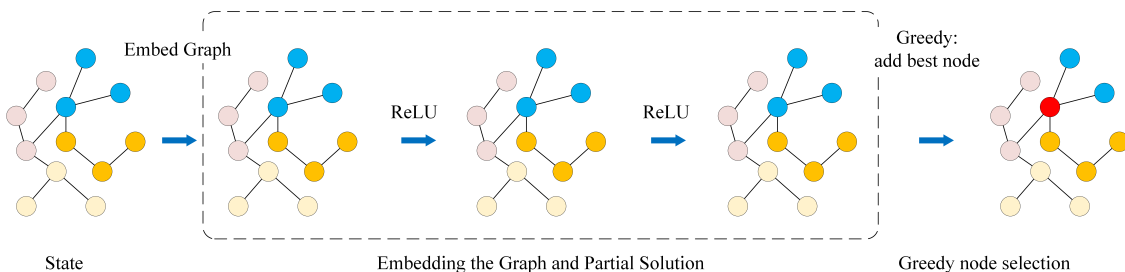
Fig. 9. An illustration of the S2V-DQN architecture [181]. This framework illustrates an instance of the S2V-DQN framework for Minimum Vertex Cove (MVC). MVC is implemented through graph embedding and iterative adding nodes. (Redrawn from [181])

dynamics of the NP-hard problem by proposing a temporal feature encoder capable of embedding instances and a decoder capable of focusing on the embedded features to find instances of a given combinatorial problem. DeepOpt [183] employs a DRL method to solve the problem of placement policies for multiple resource types in a Virtual Network Function (VNF). The algorithm proposed by Silver et al. [184] constructs the shortest path query method in spatio-temporal graphs based on static and dynamic subgraphs, and employs a Proximal Policy Optimization (PPO) algorithm [185] to train an agent to optimize the state-to-action policy to solve the path planning problem in real-world tasks.

*5) Medicine:* GRL techniques are commonly used in Clinical Decision Support (CDS) applications, Medicine Combination Prediction (MCP), chemical reaction product prediction, and brain network analysis tasks. Wang et al. [8] propose a graph convolution RL model for MCP to solve medicine correlation and sequential problems. The method represents the MCP problem as an order-free MDP and designs a DQN mechanism to learn correlative and adverse interactions between medicines. Graph Transformation Policy Network (GTPN) [9] addresses the problem of chemical reaction product prediction. They represented the entity system consisting of input reactants and reagent molecules as a labeled graph with multiple connected components by using a graph neural network, and the process of generating product molecules from reactant molecules is formulated as a series of graph transformations processes. GTPN uses A2C to find the optimal sequence of bond changes that transforms the reactants into products. Moreover, BN-GNN [186] is a brain network representation framework, which enables determine the optimal number of feature aggregations with DDQN, and the combination of abundant medical knowledge graph and DRL methods could support patients to describe disease characteristics and provide accuracy in diseases diagnosis [187].

In the pharmaceutical industry, design and discovery aids for de novo drugs serve an important research value, and many scholars have conducted studies on chemical molecule generation and optimization with RL methods [188], [189]. These methods that use RL [130], [163], [190] to pre-design the molecular structure of drugs can significantly save working time, money, and labor costs in chemical laboratories and provide an efficient aid for the design and discovery of new drugs.

## IV. OPEN RESEARCH DIRECTIONS

Although several GRL methods have been explored and proposed by scholars at present, we believe that this field is still in the initial stage and has much exploration space and research value, and we suggest some future research directions that might be of interest to scholars in this field.

**Automated RL.** Modeling of the environment, selection of RL algorithms and hyperparameter design of the model in GRL methods commonly require detailed expert exploration. The different definitions of state and action spaces, batch sizes and frequency of batch updates, and the number of timestep will result in different experimental performances. Automatic RL provides a solution for automatically making appropriate decisions about the settings of the RL process before starting to learn, enabling a non-expert way to perform, and improving the range of applications of RL. We believe that applying this automatic learning method to graph mining would significantly improve the efficiency of GRL methods and provide more efficient generation solutions for the optimal combination of graph mining algorithms and RL methods. In our survey, there are a limited number of scholars who have applied automatic RL to graph data mining tasks.

**Hierarchical RL.** Hierarchical RL methods in algorithm design allow top-level policies to focus on higher-level goals, while sub-policies are responsible for fine-grained control [193]–[195]. The scholars have proposed the implementation of the top-level policy based on manual formulation [196] and automatic policy formulation [197] to choose between sub-policies. Hierarchical policies for GRL allow models to design top-level and sub-level structures to enhance model explainability and robustness, and either using a hierarchical data processing policy or a hierarchical algorithm design policy can make the algorithms more understandable. We believe that hierarchical GRL is a worthy field for exploration.

**Multi-agent RL.** Single agents usually ignore the interests of other agents in the process of interacting with the environment. Individual agents often consider other agents in the environment as part of the environment and will fail to adapt to the instability environment during the learning process as the interaction of other agents with the environment. Multi-agent RL considers the communication between multiple agents, allowing them to collaborate learning effectively. Existing multi-agent RL methods have proposed several policies such as bidirectional channels [198], sequential transfer [199] and all-to-all channels [200] to achieve interaction between agents.

TABLE VIII
GRL ALGORITHMS ON REAL WORLD APPLICATIONS. MISSING VALUES (-) IN THIS TABLE INDICATE THAT THE PERSONALIZED TERMINATION
CONDITIONS CONTRIBUTE TO THE STABILIZATION OF THE TRAINING PROCESS.

| Method | Action | State | Reward | Termination | RL | Metrics |
|---|---|---|---|---|---|---|
| XGNN [18] | The action is defined to add an edge to the current graph by determining the starting node and the ending node of the edge. | The partially generated graph. | The reward consisting of the guidance from the trained GNNs and validated graph rules. | - | MDP | Accuracy and Metrics for explanation |
| Marl-eGCN [10] | The toll set by the transit authority on special road that has an ETC gantry. | The number of vehicles that travel to the special zone on road. | The number of vehicles which arrive at destinations. | - | MDP | Traffic throughput |
| S2V-DQN [181] | A node of original graph that is not part of the current state. | A sequence of actions (nodes) on a graph. | The change in the cost function after taking an action and transitioning to a new state. | Corresponds to tagging the node that is selected as the last action with feature equal 1. | Q-learning | Approximation ratio, Generalization ability, and Time-approximation trade-off |
| GTPN [9] | The action consisting of three consecutive sub-actions. | The intermediate graph. | The Reward is determined based on whether the model is successful in predicting | - | A2C | Coverage@k, Recall@k, and precision@K |
| IG-RL [165] | The action is defined as whether to switch to the next phase or prolong the current phase. | Current connectivity and demand in the network. | The reward for a given agent is defined as the negative sum of local queues lengths. | An episode either ends as soon as all trips are completed or after a fixed amount of time. | Q-learning | Trips Duration and Total Delay Evolution |
| GCQ [167] | The customized discrete action space ensures that does not prevent vehicle collisions during the lane change process. | The state space consisting of nodes feature, adjacency matrix, and a mask. | The reward function consisting of intention and speed rewards, collision and lane-changing penalties. | The algorithm boundary is controlled by specifying the total number of connected self-driving vehicles entering the roadway. | DQN | Episode reward and Total simulation steps per episode |
| Peng et al. [22] | The increase and decrease of the traffic flow between any two stations. | The intermediate generated graph. | The reward is defined as a sum of structural rewards and adversarial rewards. | - | PPO | Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) |
| DeepOpt [183] | The action value can be denoted with a binary list indicating the selection status of each network node. | The state information including the input of the node attributes and edge attributes. | The reward function consisting of custom penalty terms. | - | MDP | The Reject Ratio of Service Function Chain Requests, Influence of Topology Change, and Time Consumption |
| Vulcan [191] | An action selects a vertex on the graph which is not included in the partial solution and connected to the partial solution. | The state can be represented by the embedded vertex vector through graph embedding technology, including partial solution and other vertices on a graph. | The reward consisting of the cost function and the remaining path weight. | All terminals are added to the partial solution. | DQN | Gain |
| A3C+GCN [21] | An action is a valid embedding process that allocates virtual network requests onto a subset of substrate network components. | The state is the real-time representation of special network status. | Reward shaping | - | A3C | Acceptance Ratio, Long-Term Average Revenue, and Running Time |
| GCN-RL [192] | The continuous action space is defined to determine the transistor sizes. | The state consisting of transistor index, component type and the selected model feature vector for the component. | The reward is defined as a weighted sum of the normalized performance metrics. | - | DDPG | FoM and the performance of Voltage Amplifier and Voltage Amplifier |
| GMETAEXP [33] | Each action is a test input to the program, which can be text (sequences of characters) or images (2D array of characters). | The state contains the full interaction history in the episode. | Customized step-by-step rewards. | Once the agent has visited all nodes or met a custom training termination condition | MDP | Coverage performance and generalization performance |

From existing studies, scholars are trying to learn dynamic multi-agent environments and provide an abstract representation of the interactions between agents to adapt to the dynamic evolution. The collaborative policies in multi-agent RL could help scholars to solve graph mining problems with parallel and partitioned policies.

**Subgraph Patterns Mining.** In the process of network local structure analysis, many scholars have explored and discovered many novel graph data mining methods that exploit the local structures of the graphs with different strategies [15], [48], [201], [202]. In addition, they have employed batch sampling and importance sampling to obtain the local structure for network representation learning. Graph local structure selection provides natural advantages for GRL. The sequential construction of graph local structures can be represented as MDPs and guided by the performance of downstream classifi-

cation and prediction tasks. These adaptive neighbor selection methods enable the performance of GNNs to be improved and the importance of different relationships in messages passing to be analyzed in heterogeneous information networks.

**Explainability.** There is a limited number of scholars working on the explainability for graph neural networks, which is important for improving the performance of graph neural network models and assisting in understanding the intrinsic meaning and potential mechanisms of graph neural network models [203], [204]. We believe that the sequential decision-making of RL is suitable for the studies of explainability of graph structures.

**Evaluation metrics.** It is crucial to define the reward function in RL. A reasonable reward function could enable RL algorithms to converge faster and better to obtain the expected goal. However, existing GRL methods commonly employ the

performance of downstream classifiers or predictors for designing reward functions. These methods allow GRL methods to be evaluated based on the performance of downstream tasks, but the execution of GRL methods includes multiple steps, and evaluating the results of each action using downstream tasks would cause substantial resource consumption. Although many scholars have proposed novel methods to reduce this consumption or to design more efficient reward functions to guide the training process of RL, the problem of resource consumption still remains, and we believe that designing reward functions that rely on actions or providing a unified GRL evaluation framework can effectively improve the performance of the algorithm. This is an urgent problem to be solved.

## V. Conclusion

Network science has developed into an interdisciplinary field spanning physics, economics, biology, and computer science, with many critical real-world applications, and enables the modeling and analysis of the interaction of entities in complex systems [205]. In this survey, we conduct a comprehensive overview of the state-of-the-art methods in GRL, which suggests a variety of solution strategies for different graph mining tasks, in which RL methods are combined to cope with existing challenges better. Although GRL methods have been applied in many fields, the in-depth combination of more efficient RL with excellent GNNs is expected to provide better generalization and explainability and improve the ability to tackle sample complexity. We believe that GRL methods could reveal valuable information in the growing volume of graph-structured data and enable scholars and engineers to analyze and exploit graph-structured data with more clarity. We hope that this survey will help scholars to understand the key concepts of GRL and drive the field forward in the future.

## Appendix A
### Abbreviations

We summarize the full descriptions and abbreviations used in this survey to facilitate searching by scholars in Table IX.

## Appendix B
### Notations

We provide the commonly used notations and explanations in different domains in Table X.

## Appendix C
### Open-source Implementations

Here we summarize the open-source implementations of GRL in this survey in Table XI.

### References

[1] S. Levine, P. Pastor, A. Krizhevsky, J. Ibarz, and D. Quillen, "Learning hand-eye coordination for robotic grasping with deep learning and large-scale data collection," *The International journal of robotics research*, vol. 37, no. 4-5, pp. 421–436, 2018.

[2] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.

### TABLE IX
A LIST OF ABBREVIATIONS.

| Abbreviation | Full description |
| --- | --- |
| AC | Actor-Critic |
| A2C | Advantage Actor-Critic |
| A3C | Asynchronous Advantage Actor-Critic |
| AutoMl | Automated Machine Learning |
| BMAB | Bernoulli Multi-armed Bandit |
| CDS | Clinical Decision Support |
| CDQN | Cascaded Deep Q-Network |
| DRL | Deep reinforcement learning |
| DAG | Directed Acyclic Graph |
| DQN | Deep Q-Network |
| DDQN | Double Deep Q-Network |
| DPG | Deterministic Policy Gradient |
| DDPG | Deep Deterministic Policy Gradient |
| DETC | Dynamic Electronic Toll Collection |
| ETC | Electronic Toll Collection |
| GNN | Graph neural network |
| GRL | Graph Reinforcement Learning |
| GAT | Graph attention network |
| JSSP | Job Shop Scheduling Problem |
| LSTM | Long Short-Term Memory |
| MDP | Markov Decision Process |
| MCTS | Monte Carlo Tree Search |
| MI | Mutual Information |
| MVC | Minimum Vertex Cove |
| NRL | Network Representation Learning |
| NLP | Natural Language Processing |
| PPO | Proximal Policy Optimization |
| POMDP | Partially Observable Markov Decision Process |
| Q&A | Question and Answer |
| RL | Reinforcement learning |
| SCST | Self-Critical Sequence Training |
| SDS | Spoken Dialogue System |
| TSP | Traveling Salesman Problem |
| TD | Temporal Difference |

[3] W. Xiong, T. Hoang, and W. Y. Wang, "Deeppath: A reinforcement learning method for knowledge graph reasoning," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, Copenhagen, Denmark, 2017, pp. 564–573.

[4] Y. Duan, X. Chen, R. Houthooft, J. Schulman, and P. Abbeel, "Benchmarking deep reinforcement learning for continuous control," in *International conference on machine learning*, 2016, pp. 1329–1338.

[5] Y. Deng, F. Bao, Y. Kong, Z. Ren, and Q. Dai, "Deep direct reinforcement learning for financial signal representation and trading," *IEEE transactions on neural networks and learning systems*, vol. 28, no. 3, pp. 653–664, 2016.

[6] J. B. Heaton, N. G. Polson, and J. H. Witte, "Deep learning for finance: deep portfolios," *Applied Stochastic Models in Business and Industry*, vol. 33, no. 1, pp. 3–12, 2017.

[7] F. Soleymani and E. Paquet, "Deep graph convolutional reinforcement learning for financial portfolio management-deeppocket," *Expert Systems with Applications*, vol. 182, pp. 115–127, 2021.

[8] S. Wang, P. Ren, Z. Chen, Z. Ren, J. Ma, and M. de Rijke, "Order-free medicine combination prediction with graph convolutional reinforcement learning," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, Beijing, China, 2019, pp. 1623–1632.

[9] K. Do, T. Tran, and S. Venkatesh, "Graph transformation policy network for chemical reaction prediction," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, ser. KDD '19, Anchorage, AK, USA, 2019, pp. 750–760.

[10] W. Qiu, H. Chen, and B. An, "Dynamic electronic toll collection via multi-agent deep reinforcement learning with edge-based graph convolutional networks," in *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, 2019, pp. 4568–4574.

[11] M. Obara, T. Kashiyama, and Y. Sekimoto, "Deep reinforcement learning approach for train rescheduling utilizing graph theory," in *2018 IEEE International Conference on Big Data (Big Data)*, 2018, pp. 4525–4533.

[12] M. Zhang and Y. Chen, "Link prediction based on graph neural

TABLE X
COMMONLY USED NOTATIONS.

| Notation | Explanation | Domain |
|---|---|---|
| $\lvert \cdot \rvert$ | The length of a set. | GNN |
| G | A graph. | GNN |
| $A \in \{0,1\}^{m \times m}$ | The graph adjacency matrix. | GNN |
| $X_i \in R^{m \times d_i}$ | The feature matrix of a graph in layer $i$. | GNN |
| $E$ | The set of edges in a graph. $e \in E$. | GNN |
| $V$ | The set of nodes in a graph. $v \in V$. | GNN |
| $D$ | The degree matrix of A. $D_{ii} = \sum_{j=1}^{n} A_{ij}$. | GNN |
| $n$ | The number of nodes, $n = \lvert V \rvert$. | GNN |
| $m$ | The number of edges, $m = \lvert E \rvert$. | GNN |
| $d$ | The dimension of a node feature vector. | GNN |
| $W_i \in R^{d_i \times d_{i+1}}$ | Learnable model parameters. | GNN |
| $\sigma(\cdot)$ | The sigmoid activation function. | GNN |
| $\mathcal{F} : v \to e_v \in R^d$ | Mapping functions in graph embedding. | GNN |
| $(h, r, t)$ | $h$ denotes the head entity, $t$ denotes the tail entity, $r$ denotes the relationship between $h$ and $t$. | KG |
| $s_t \in \mathcal{S}$ | The set of all possible states | RL |
| $a_t \in \mathcal{A}$ | The set of actions that are available in the state. | RL |
| $r_t \in \mathcal{R}$ | The reward given to the agent by the environment after the agent performs an action. | RL |
| $\mathcal{T}$ | The state transfer function is defined by the environment. | RL |
| $\gamma$ | Discount Factor. | RL |
| $\pi : \mathcal{S} \to p(\mathcal{A})$ | Policy for agent to perform in the environment. | RL |
| $a_{i,j}$ | Non-negative learning factor | REINFORCE |
| $b_{i,j}$ | Representation function of the state for reducing the variance of the gradient estimate. | REINFORCE |
| $g_i$ | The probability density function for randomly generating unit activation-based actions. | REINFORCE |
| $L(\theta_i)$ | Loss function | DQN |

TABLE XI
A SUMMARY OF OPEN-SOURCE IMPLEMENTATIONS

| Model | Year | Framework | Link | Source |
|---|---|---|---|---|
| AGILE | 2022 | PyTorch | https://github.com/clvrai/agile | [35] |
| GTA-RL | 2022 | PyTorch | https://github.com/udeshmg/GTA-RL | [39] |
| SubgraphX | 2021 | PyTorch | https://github.com/divelab/DIG | [40] |
| SUGAR | 2021 | Tensorflow | https://github.com/RingBDStack/SUGAR | [15] |
| CORL | 2021 | PyTorch | https://github.com/huawei-noah/trustworthyAI/tree/master/gcastle | [143] |
| RioGNN | 2021 | PyTorch | https://github.com/safe-graph/RioGNN | [83] |
| IG-RL | 2021 | PyTorch | https://github.com/FXDevailly/IG-RL | [165] |
| TITer | 2021 | Python | https://github.com/JHL-HUST/TITer | [148] |
| SparRL | 2021 | PyTorch | https://github.com/rwickman/SparRL-PyTorch | [43] |
| PAAR | 2021 | PyTorch | https://github.com/seukgcode/PAAR | [98] |
| Policy-GNN | 2020 | PyTorch | https://github.com/lhenry15/Policy-GNN | [48] |
| CARE-GNN | 2020 | PyTorch | https://github.com/YingtongDou/CARE-GNN | [42] |
| RL-BIC | 2020 | Tensorflow | https://github.com/huawei-noah/trustworthyAI/tree/master/Causal_Structure_Learning/ Causal_Discovery_RL | [142] |
| KG-A2C | 2020 | PyTorch | https://github.com/rajammanabrolu/KG-A2C | [206] |
| GPA | 2020 | PyTorch | https://github.com/ShengdingHu/GraphPolicyNetworkActiveLearning | [14] |
| GAEA | 2020 | Tensorflow | https://github.com/salesforce/GAEA | [207] |
| CompNet | 2019 | PyTorch | https://github.com/WOW5678/CompNet | [8] |
| GRPI | 2019 | Python | https://github.com/LASP-UCL/Graph-RL | [122] |
| DRL+GNN | 2019 | PyTorch | https://github.com/knowledgedefinednetworking/DRL-GNN | [171] |
| GPN | 2019 | PyTorch | https://github.com/qiang-ma/graph-pointer-network | [195] |
| PGPR | 2019 | PyTorch | https://github.com/orcax/PGPR | [208] |
| DGN | 2018 | PyTorch | https://github.com/PKU-AI-Edge/DGN | [32] |
| GCPN | 2018 | Python | https://github.com/bowenliu16/rl_graph_generation | [163] |
| KG-DQN | 2018 | PyTorch | https://github.com/rajammanabrolu/KG-DQN | [209] |
| ASNets | 2018 | Tensorflow | https://github.com/qxcv/asnets | [182] |
| S2V-DQN | 2017 | C+Python | https://github.com/Hanjun-Dai/graph_comb_opt | [181] |
| DeepPath | 2017 | Tensorflow | https://github.com/xwhan/DeepPath | [3] |
| MINERVA | 2017 | Tensorflow | https://github.com/shehzaadzd/MINERVA | [103] |
| KBGAN | 2017 | PyTorch | https://github.com/cai-lw/KBGAN | [107] |

networks," *Advances in neural information processing systems*, vol. 31, 2018.

[13] T. Trouillon, J. Welbl, S. Riedel, E. Gaussier, and G. Bouchard, "Complex embeddings for simple link prediction," in *Proceedings of The 33rd International Conference on Machine Learning*, M. F. Balcan and K. Q. Weinberger, Eds., vol. 48, 2016, pp. 2071–2080.

[14] S. Hu, Z. Xiong, M. Qu, X. Yuan, M.-A. Côté, Z. Liu, and J. Tang, "Graph policy network for transferable active learning on graphs," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 10 174–10 185.

[15] Q. Sun, J. Li, H. Peng, J. Wu, Y. Ning, P. S. Yu, and L. He, "Sugar: Subgraph neural network with reinforcement pooling and self-

supervised mutual information mechanism," in *Proceedings of the Web Conference 2021*, Ljubljana, Slovenia, 2021, pp. 2081–2091.

[16] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, Long Beach, California, USA, 2017, pp. 1025–1035.

[17] H. Yuan and S. Ji, "Structpool: Structured graph pooling via conditional random fields," in *Proceedings of the 8th International Conference on Learning Representations*, 2020.

[18] H. Yuan, J. Tang, X. Hu, and S. Ji, "Xgnn: Towards model-level explanations of graph neural networks," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 430–438.

[19] H. Dai, H. Li, T. Tian, X. Huang, L. Wang, J. Zhu, and L. Song, "Adversarial attack on graph structured data," in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80, 2018, pp. 1115–1124.

[20] P. T. A. Quang, Y. Hadjadj-Aoul, and A. Outtagarts, "A deep reinforcement learning approach for vnf forwarding graph embedding," *IEEE Transactions on Network and Service Management*, vol. 16, no. 4, pp. 1318–1331, 2019.

[21] Z. Yan, J. Ge, Y. Wu, L. Li, and T. Li, "Automatic virtual network embedding: A deep reinforcement learning approach with graph convolutional networks," *IEEE Journal on Selected Areas in Communications*, vol. 38, no. 6, pp. 1040–1057, 2020.

[22] H. Peng, B. Du, M. Liu, M. Liu, S. Ji, S. Wang, X. Zhang, and L. He, "Dynamic graph convolutional network for long-term traffic flow prediction with reinforcement learning," *Information Sciences*, vol. 578, pp. 401–416, 2021.

[23] E. Meirom, H. Maron, S. Mannor, and G. Chechik, "Controlling graph dynamics with reinforcement learning and graph neural networks," in *International Conference on Machine Learning*, vol. 139, 2021, pp. 7565–7577.

[24] Y. Lyu, X. Yang, J. Liu, S. Xie, and X. Zhang, "Interpretable and effective reinforcement learning for attacking against graph-based rumor detection," *arXiv preprint arXiv:2201.05819*, 2022.

[25] S. Zhou, X. Dai, H. Chen, W. Zhang, K. Ren, R. Tang, X. He, and Y. Yu, "Interactive recommender system via knowledge graph-enhanced reinforcement learning," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 179–188.

[26] Y. Lei, H. Pei, H. Yan, and W. Li, "Reinforcement learning based recommendation with graph convolutional q-network," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 1757–1760.

[27] Y. Wu, X. Xi, and J. He, "Afgsl: Automatic feature generation based on graph structure learning," *Knowledge-Based Systems*, vol. 238, p. 107835, 2022.

[28] M. Zhang, X. Yu, J. Rong, and L. Ou, "Graph pruning for model compression," *Applied Intelligence*, pp. 1–13, 2022.

[29] S. Yu, A. Mazaheri, and A. Jannesari, "Gnn-rl compression: Topology-aware network pruning using multi-stage graph embedding and reinforcement learning," *arXiv preprint arXiv:2102.03214*, 2021.

[30] S. Yu, A. Mazaheri, and A. Jannesari, "Auto graph encoder-decoder for neural network pruning," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 6362–6372.

[31] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, "Fully decentralized multi-agent reinforcement learning with networked agents," in *International Conference on Machine Learning*, vol. 80, 2018, pp. 5872–5881.

[32] J. Jiang, C. Dun, T. Huang, and Z. Lu, "Graph convolutional reinforcement learning," *arXiv preprint arXiv:1810.09202*, 2018.

[33] H. Dai, Y. Li, C. Wang, R. Singh, P.-S. Huang, and P. Kohli, "Learning transferable graph exploration," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[34] J. Jin, S. Zhou, W. Zhang, T. He, Y. Yu, and R. Fakoor, "Graph-enhanced exploration for goal-oriented reinforcement learning," in *Proceedings of the Tenth International Conference on Learning Representations*, 2022.

[35] A. Jain, N. Kosaka, K.-M. Kim, and J. J. Lim, "Know your action set: Learning action relations for reinforcement learning," in *Proceedings of the Tenth International Conference on Learning Representations*, 2022.

[36] Z. Zhang, Z. Yang, H. Liu, P. Tokekar, and F. Huang, "Reinforcement learning under a multi-agent predictive state representation model: Method and theory," in *Proceedings of the Tenth International Conference on Learning Representations*, 2022.

[37] S. Hong, D. Yoon, and K.-E. Kim, "Structure-aware transformer policy for inhomogeneous multi-task reinforcement learning," in *International Conference on Learning Representations*, 2022.

[38] H. Liu, S. Zhou, C. Chen, T. Gao, J. Xu, and M. Shu, "Dynamic knowledge graph reasoning based on deep reinforcement learning," *Knowledge-Based Systems*, vol. 241, p. 108235, 2022.

[39] U. Gunarathna, R. Borovica-Gajic, S. Karunasekara, and E. Tanin, "Solving dynamic graph problems with multi-attention deep reinforcement learning," *arXiv preprint arXiv:2201.04895*, 2022.

[40] H. Yuan, H. Yu, J. Wang, K. Li, and S. Ji, "On explainability of graph neural networks via subgraph explorations," in *Proceedings of the 38th International Conference on Machine Learning*, vol. 139, 2021, pp. 12 241–12 252.

[41] J. Jiang, C. Dun, T. Huang, and Z. Lu, "Graph convolutional reinforcement learning," *arXiv preprint arXiv:1810.09202*, 2018.

[42] Y. Dou, Z. Liu, L. Sun, Y. Deng, H. Peng, and P. S. Yu, "Enhancing graph neural network-based fraud detectors against camouflaged fraudsters," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, Virtual Event, Ireland, 2020, pp. 315–324.

[43] R. Wickman, X. Zhang, and W. Li, "Sparrl: Graph sparsification via deep reinforcement learning," *arXiv preprint arXiv:2112.01565*, 2021.

[44] K. Rusek, J. Surez-Varela, A. Mestres, P. Barlet-Ros, and A. Cabellos-Aparicio, "Unveiling the potential of graph neural networks for network modeling and optimization in sdn," in *Proceedings of the 2019 ACM Symposium on SDN Researc*, San Jose, CA, USA, 2019.

[45] J. Surez-Varela, S. Carol-Bosch, K. Rusek, P. Almasan, M. Arias, P. Barlet-Ros, and A. Cabellos-Aparicio, "Challenging the generalization capabilities of graph neural networks for network modeling," in *Proceedings of the ACM SIGCOMM 2019 Conference Posters and Demos*, Beijing, China, 2019, pp. 114–115.

[46] R. Wang, F. Fang, J. Cui, and W. Zheng, "Learning self-driven collective dynamics with graph networks," *Scientific reports*, vol. 12, no. 1, pp. 1–11, 2022.

[47] N. Shervashidze, P. Schweitzer, E. J. van Leeuwen, K. Mehlhorn, and K. M. Borgwardt, "Weisfeiler-lehman graph kernels," *The Journal of Machine Learning Research*, vol. 12, pp. 2539–2561, 2011.

[48] K.-H. Lai, D. Zha, K. Zhou, and X. Hu, "Policy-gnn: Aggregation optimization for graph neural networks," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Virtual Event, CA, USA, 2020, pp. 461–471.

[49] D. Zhang, J. Yin, X. Zhu, and C. Zhang, "Network representation learning: A survey," *IEEE Transactions on Big Data*, vol. 6, no. 1, pp. 3–28, 2020.

[50] M. Sun, J. Tang, H. Li, B. Li, C. Xiao, Y. Chen, and D. Song, "Data poisoning attack against unsupervised node embedding methods," *arXiv preprint arXiv:1810.12881*, 2018.

[51] L. Sun, Y. Dou, C. Yang, J. Wang, P. S. Yu, L. He, and B. Li, "Adversarial attack and defense on graph data: A survey," *arXiv preprint arXiv:1812.10528*, 2018.

[52] Y. Ma, S. Wang, T. Derr, L. Wu, and J. Tang, "Attacking graph convolutional networks via rewiring," *arXiv preprint arXiv:1906.03750*, 2019.

[53] Z. Xi, R. Pang, S. Ji, and T. Wang, "Graph backdoor," in *30th USENIX Security Symposium (USENIX Security 21)*, 2021, pp. 1523–1540.

[54] Z. Zhang, J. Jia, B. Wang, and N. Z. Gong, "Backdoor attacks to graph neural networks," in *Proceedings of the 26th ACM Symposium on Access Control Models and Technologies*, New York, NY, USA, 2021, pp. 15–26.

[55] S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives, "Dbpedia: A nucleus for a web of open data," in *The semantic web*, 2007, pp. 722–735.

[56] K. Bollacker, C. Evans, P. Paritosh, T. Sturge, and J. Taylor, "Freebase: A collaboratively created graph database for structuring human knowledge," in *Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data*, New York, NY, USA, 2008, pp. 1247–1250.

[57] F. Mahdisoltani, J. Biega, and F. Suchanek, "Yago3: A knowledge base from multilingual wikipedias," in *7th biennial conference on innovative data systems research*, Asilomar, California, USA, 2015.

[58] X. Zhao, L. Chen, and H. Chen, "A weighted heterogeneous graph-based dialog system," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–6, 2021.

[59] Q. Wang, Y. Ji, Y. Hao, and J. Cao, "Grl: Knowledge graph completion with gan-based reinforcement learning," *Knowledge-Based Systems*, vol. 209, p. 106421, 2020.

[60] M. Gardner, P. Talukdar, B. Kisiel, and T. Mitchell, "Improving learning and inference in a large knowledge-base using latent syntactic cues," in *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, Seattle, Washington, USA, 2013, pp. 833–838.

[61] A. Bordes, N. Usunier, A. Garcia-Durán, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Proceedings of the 26th International Conference on Neural Information Processing Systems*, vol. 2, Lake Tahoe, Nevada, 2013, pp. 2787–2795.

[62] Z. Wang, J. Zhang, J. Feng, and Z. Chen, "Knowledge graph embedding by translating on hyperplanes," in *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence*, Québec City, Québec, Canada, 2014, pp. 1112–1119.

[63] T. Wang, R. Liao, J. Ba, and S. Fidler, "Nervenet: Learning structured policy with graph neural networks," in *International conference on learning representations*, 2018.

[64] J. Schrittwieser, I. Antonoglou, T. Hubert, K. Simonyan, L. Sifre, S. Schmitt, A. Guez, E. Lockhart, D. Hassabis, T. Graepel, T. Lillicrap, and D. Silver, "Mastering atari, go, chess and shogi by planning with a learned model," *Nature*, vol. 588, no. 7839, pp. 604–609, 2020.

[65] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, and M. Lanctot, "Mastering the game of go with deep neural networks and tree search," *nature*, vol. 529, no. 7587, pp. 484–489, 2016.

[66] T. M. Moerland, J. Broekens, and C. M. Jonker, "Model-based reinforcement learning: A survey," *arXiv preprint arXiv:2006.16712*, 2020.

[67] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[68] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.

[69] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, no. 3, pp. 279–292, 1992.

[70] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3, pp. 229–256, 1992.

[71] V. R. Konda and J. N. Tsitsiklis, "Actor-critic algorithms," in *Advances in neural information processing systems*, vol. 12, Denver, Colorado, USA, 1999, pp. 1008–1014.

[72] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Harley, T. P. Lillicrap, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *Proceedings of the 33rd International Conference on International Conference on Machine Learning*, vol. 48, 2016, pp. 1928–1937.

[73] J. X. Wang, Z. Kurth-Nelson, D. Tirumala, H. Soyer, J. Z. Leibo, R. Munos, C. Blundell, D. Kumaran, and M. Botvinick, "Learning to reinforcement learn," *arXiv preprint arXiv:1611.05763*, 2016.

[74] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *Proceedings of the 31st International Conference on International Conference on Machine Learning*, vol. 32, 2014, pp. I–387–I–395.

[75] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.

[76] H. v. Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*, Phoenix, Arizona, 2016, pp. 2094–2100.

[77] M. Hausknecht and P. Stone, "Deep recurrent q-learning for partially observable mdps," *arXiv preprint arXiv:1507.06527*, 2015.

[78] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *Proceedings of The 33rd International Conference on Machine Learning*, vol. 48, New York, New York, USA, 20–22 Jun 2016, pp. 1995–2003.

[79] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, "Collective classification in network data," *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.

[80] T. Dettmers, P. Minervini, P. Stenetorp, and S. Riedel, "Convolutional 2d knowledge graph embeddings," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, New Orleans, Louisiana, USA, 2018.

[81] K. Toutanova, D. Chen, P. Pantel, H. Poon, P. Choudhury, and M. Gamon, "Representing text for joint embedding of text and knowledge bases," in *Proceedings of the 2015 conference on empirical methods in natural language processing*, 2015, pp. 1499–1509.

[82] A. Mukherjee, V. Venkataraman, B. Liu, and N. Glance, "What yelp fake review filter might be doing?" in *Proceedings of the International AAAI Conference on Web and Social Media*, vol. 7, 2013, pp. 409–418.

[83] H. Peng, R. Zhang, Y. Dou, R. Yang, J. Zhang, and P. S. Yu, "Reinforced neighborhood selection guided multi-relational graph neural networks," *ACM Transactions on Information Systems (TOIS)*, vol. 40, no. 4, pp. 1–46, 2021.

[84] Y. Deng, Y. Li, F. Sun, B. Ding, and W. Lam, "Unified conversational recommendation policy learning via graph-based reinforcement learning," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1431–1441.

[85] J. J. McAuley and J. Leskovec, "From amateurs to connoisseurs: modeling the evolution of user expertise through online reviews," in *Proceedings of the 22nd international conference on World Wide Web*, 2013, pp. 897–908.

[86] A. L. Goldberger, L. A. Amaral, L. Glass, J. M. Hausdorff, P. C. Ivanov, R. G. Mark, J. E. Mietus, G. B. Moody, C.-K. Peng, and H. E. Stanley, "Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals," *circulation*, vol. 101, no. 23, pp. e215–e220, 2000.

[87] Y. Gao, H. Yang, P. Zhang, C. Zhou, and Y. Hu, "Graphnas: Graph neural architecture search with reinforcement learning," *arXiv preprint arXiv:1904.09981*, 2019.

[88] K. Zhou, Q. Song, X. Huang, and X. Hu, "Auto-gnn: Neural architecture search of graph neural networks," *arXiv preprint arXiv:1909.03184*, 2019.

[89] L. Wang, W. Yu, W. Wang, W. Cheng, W. Zhang, H. Zha, X. He, and H. Chen, "Learning robust representations with graph denoising policy network," in *2019 IEEE International Conference on Data Mining (ICDM)*, 2019, pp. 1378–1383.

[90] Y. Sun, S. Wang, X. Tang, T.-Y. Hsieh, and V. Honavar, "Adversarial attacks on graph neural networks via node injections: A hierarchical reinforcement learning approach," in *Proceedings of The Web Conference 2020*, 2020, pp. 673–683.

[91] S. Shen, Y. Fu, A. L. Jia, H. Su, Q. Wang, C. Wang, and Y. Dou, "Learning network representation through reinforcement learning," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 3537–3541.

[92] J. Leskovec and J. Mcauley, "Learning to discover social circles in ego networks," in *Advances in Neural Information Processing Systems*, F. Pereira, C. Burges, L. Bottou, and K. Weinberger, Eds., vol. 25, 2012.

[93] C. Wang, Y. Liu, X. Gao, and G. Chen, "A reinforcement learning model for influence maximization in social networks," in *International Conference on Database Systems for Advanced Applications*, Cham, 2021, pp. 701–709.

[94] A. Mislove, M. Marcon, K. P. Gummadi, P. Druschel, and B. Bhattacharjee, "Measurement and analysis of online social networks," in *Proceedings of the 7th ACM SIGCOMM Conference on Internet Measurement*, San Diego, California, USA, 2007, pp. 29–42.

[95] J. Leskovec, J. Kleinberg, and C. Faloutsos, "Graph evolution: Densification and shrinking diameters," *ACM transactions on Knowledge Discovery from Data (TKDD)*, vol. 1, no. 1, pp. 2–es, 2007.

[96] H. Wang, S. Li, R. Pan, and M. Mao, "Incorporating graph attention mechanism into knowledge graph reasoning based on deep reinforcement learning," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Hong Kong, China, 2019, pp. 2623–2631.

[97] G. Wan, S. Pan, C. Gong, C. Zhou, and G. Haffari, "Reasoning like human: Hierarchical reinforcement learning for knowledge graph reasoning," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, Yokohama, Yokohama, Japan, 2021, pp. 1926–1932.

[98] X. Zhou, P. Wang, Q. Luo, and Z. Pan, "Multi-hop knowledge graph reasoning based on hyperbolic knowledge graph embedding and reinforcement learning," in *The 10th International Joint Conference on Knowledge Graphs*, Virtual Event, Thailand, 2021, pp. 1–9.

[99] M. Zheng, Y. Zhou, and Q. Cui, "Hierarchical policy network with multi-agent for knowledge graph reasoning based on reinforcement learning," in *International Conference on Knowledge Science, Engineering and Management*, Cham, 2021, pp. 445–457.

[100] Q. Wang, Y. Hao, and J. Cao, "Adrl: An attention-based deep reinforcement learning framework for knowledge graph reasoning," *Knowledge-Based Systems*, vol. 197, p. 105910, 2020.

[101] Z. Li, X. Jin, S. Guan, Y. Wang, and X. Cheng, "Path reasoning over knowledge graph: A multi-agent and reinforcement learning based

method," in *2018 IEEE International Conference on Data Mining Workshops (ICDMW)*, 2018, pp. 929–936.

[102] X. V. Lin, R. Socher, and C. Xiong, "Multi-hop knowledge graph reasoning with reward shaping," in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, Brussels, Belgium, 2018, pp. 3243–3253.

[103] R. Das, S. Dhuliawala, M. Zaheer, L. Vilnis, I. Durugkar, A. Krishnamurthy, A. Smola, and A. McCallum, "Go for a walk and arrive at the answer: Reasoning over paths in knowledge bases using reinforcement learning," *arXiv preprint arXiv:1711.05851*, 2017.

[104] L. Chen, J. Cui, X. Tang, Y. Qian, Y. Li, and Y. Zhang, "Rlpath: a knowledge graph link prediction method using reinforcement learning based attentive relation path searching and representation learning," *Applied Intelligence*, vol. 52, no. 4, pp. 4715–4726, 2022.

[105] P. Tiwari, H. Zhu, and H. M. Pandey, "Dapath: Distance-aware knowledge graph reasoning based on deep reinforcement learning," *Neural Networks*, vol. 135, pp. 1–12, 2021.

[106] S. Li, H. Wang, R. Pan, and M. Mao, "Memorypath: A deep reinforcement learning framework for incorporating memory component into knowledge graph reasoning," *Neurocomputing*, vol. 419, pp. 273–286, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0925231220312959

[107] L. Cai and W. Y. Wang, "Kbgan: Adversarial learning for knowledge graph embeddings," in *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 2018, pp. 1470–1480.

[108] A. K. Debnath, R. L. Lopez de Compadre, G. Debnath, A. J. Shusterman, and C. Hansch, "Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. correlation with molecular orbital energies and hydrophobicity," *Journal of medicinal chemistry*, vol. 34, no. 2, pp. 786–797, 1991.

[109] Y. Luo, M. McThrow, W. Y. Au, T. Komikado, K. Uchino, K. Maruhash, and S. Ji, "Automated data augmentations for graph classification," *arXiv preprint arXiv:2202.13248*, 2022.

[110] H. Toivonen, A. Srinivasan, R. D. King, S. Kramer, and C. Helma, "Statistical evaluation of the predictive toxicology challenge 2000-2001," *Bioinformatics*, vol. 19, no. 10, pp. 1183–1193, 2003.

[111] K. M. Borgwardt, C. S. Ong, S. Schnauer, S. Vishwanathan, A. J. Smola, and H.-P. Kriegel, "Protein function prediction via graph kernels," *Bioinformatics*, vol. 21, pp. i47–i56, 2005.

[112] P. D. Dobson and A. J. Doig, "Distinguishing enzyme structures from non-enzymes without alignments," *Journal of molecular biology*, vol. 330, no. 4, pp. 771–783, 2003.

[113] N. Wale, I. A. Watson, and G. Karypis, "Comparison of descriptor spaces for chemical compound retrieval and classification," *Knowledge and Information Systems*, vol. 14, no. 3, pp. 347–375, 2008.

[114] Z. Wu, B. Ramsundar, E. N. Feinberg, J. Gomes, C. Geniesse, A. S. Pappu, K. Leswing, and V. Pande, "Moleculenet: a benchmark for molecular machine learning," *Chemical science*, vol. 9, no. 2, pp. 513–530, 2018.

[115] H. Yuan, H. Yu, S. Gui, and S. Ji, "Explainability in graph neural networks: A taxonomic survey," *arXiv preprint arXiv:2012.15445*, 2020.

[116] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

[117] P. Velikovi, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[118] H. Gao, Z. Wang, and S. Ji, "Large-scale learnable graph convolutional networks," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, London, United Kingdom, 2018, pp. 1416–1424.

[119] P. Veličković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax," *arXiv preprint arXiv:1809.10341*, 2018.

[120] Y. Lin, Z. Liu, M. Sun, Y. Liu, and X. Zhu, "Learning entity and relation embeddings for knowledge graph completion," in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, Austin, Texas, 2015, pp. 2181–2187.

[121] W.-t. Yih, K. Toutanova, J. C. Platt, and C. Meek, "Learning discriminative projections for text similarity measures," in *Proceedings of the Fifteenth Conference on Computational Natural Language Learning*, Portland, Oregon, 2011, pp. 247–256.

[122] S. Madjiheurem and L. Toni, "Representation learning on graphs: A reinforcement learning application," in *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, vol. 89, 2019, pp. 3391–3399.

[123] D. Chen, M. Nie, H. Zhang, Z. Wang, and D. Wang, "Network embedding algorithm taking in variational graph autoencoder," *Mathematics*, vol. 10, no. 3, p. 485, 2022.

[124] B. Li and D. Pi, "Network representation learning: a systematic literature review," *Neural Computing and Applications*, vol. 32, no. 21, pp. 16 647–16 679, 2020.

[125] X. Huang, D. Chen, T. Ren, and D. Wang, "A survey of community detection methods in multilayer networks," *Data Mining and Knowledge Discovery*, vol. 35, no. 1, pp. 1–45, 2021.

[126] J. Chen, T. Ma, and C. Xiao, "Fastgcn: fast learning with graph convolutional networks via importance sampling," *arXiv preprint arXiv:1801.10247*, 2018.

[127] Z. Li, Y. Sun, S. Tang, C. Zhang, and H. Ma, "Reinforcement learning with dual attention guided graph convolution for relation extraction," in *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 2021, pp. 946–953.

[128] X. Zhao, Q. Dai, J. Wu, H. Peng, M. Liu, X. Bai, J. Tan, S. Wang, and P. Yu, "Multi-view tensor graph neural networks through reinforced aggregation," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2022.

[129] X. Fu, J. Li, J. Wu, Q. Sun, C. Ji, S. Wang, J. Tan, H. Peng, and P. S. Yu, "Ace-hgnn: Adaptive curvature exploration hyperbolic graph neural network," in *2021 IEEE International Conference on Data Mining (ICDM)*, 2021, pp. 111–120.

[130] J. B. Lee, R. Rossi, and X. Kong, "Graph classification using structural attention," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, London, United Kingdom, 2018, pp. 1666–1674.

[131] J. Xu, Y. Yang, S. Pu, Y. Fu, J. Feng, W. Jiang, J. Lu, and C. Wang, "Netrl: Task-aware network denoising via deep reinforcement learning," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2021.

[132] Y. Qin, X. Wang, P. Cui, and W. Zhu, "Gqnas: Graph q network for neural architecture search," in *Proceedings of the IEEE International Conference on Data Mining (ICDM)*, 2021, pp. 1288–1293.

[133] J. Jin, J. Qin, Y. Fang, K. Du, W. Zhang, Y. Yu, Z. Zhang, and A. J. Smola, "An efficient neighborhood-based interaction model for recommendation on heterogeneous graph," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Virtual Event, CA, USA, 2020, pp. 75–84.

[134] D. Seyler, P. Chandar, and M. Davis, "An information retrieval framework for contextual suggestion based on heterogeneous information network embeddings," in *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, Ann Arbor, MI, USA, 2018, pp. 953–956.

[135] H. Peng, J. Li, Z. Wang, R. Yang, M. Liu, M. Zhang, P. Yu, and L. He, "Lifelong property price prediction: A case study for the toronto real estate market," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–1, 2021.

[136] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, "Heterogeneous graph attention network," in *The World Wide Web Conference*, San Francisco, CA, USA, 2019, pp. 2022–2032.

[137] C. Yang, Y. Feng, P. Li, Y. Shi, and J. Han, "Meta-graph based hin spectral embedding: Methods, analyses, and insights," in *2018 IEEE International Conference on Data Mining (ICDM)*, 2018, pp. 657–666.

[138] Z. Zhong, C.-T. Li, and J. Pang, "Reinforcement learning enhanced heterogeneous graph neural network," *arXiv preprint arXiv:2010.13735*, 2020.

[139] H. Peng, R. Zhang, S. Li, Y. Cao, S. Pan, and P. Yu, "Reinforced, incremental and cross-lingual event detection from social messages," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2022.

[140] D. Zügner, A. Akbarnejad, and S. Günnemann, "Adversarial attacks on neural networks for graph data," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, London, United Kingdom, 2018, pp. 2847–2856.

[141] J. Dineen, A. Haque, and M. Bielskas, "Reinforcement learning for data poisoning on graph neural networks," in *International Conference on Social Computing, Behavioral-Cultural Modeling and Prediction and Behavior Representation in Modeling and Simulation*, 2021, pp. 141–150.

[142] S. Zhu, I. Ng, and Z. Chen, "Causal discovery with reinforcement learning," *arXiv preprint arXiv:1906.04477*, 2019.

[143] X. Wang, Y. Du, S. Zhu, L. Ke, Z. Chen, J. Hao, and J. Wang, "Ordering-based causal discovery with reinforcement learning," *arXiv preprint arXiv:2105.06631*, 2021.

[144] Y. Sun, K. Zhang, and C. Sun, "Model-based transfer reinforcement learning based on graphical model representations," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–14, 2021.

[145] L. Chen, B. Tan, S. Long, and K. Yu, "Structured dialogue policy with graph neural networks," in *Proceedings of the 27th International Conference on Computational Linguistics*, Santa Fe, New Mexico, USA, 2018, pp. 1257–1268.

[146] Y. Chen, L. Wu, and M. J. Zaki, "Reinforcement learning based graph-to-sequence model for natural question generation," *arXiv preprint arXiv:1908.04942*, 2019.

[147] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1179–1195.

[148] H. Sun, J. Zhong, Y. Ma, Z. Han, and K. He, "Timetraveler: Reinforcement learning for temporal knowledge graph forecasting," *arXiv preprint arXiv:2109.04101*, 2021.

[149] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Proceedings of the 12th International Conference on Neural Information Processing Systems*, Denver, CO, 1999, pp. 1057–1063.

[150] W. Song, Z. Duan, Z. Yang, H. Zhu, M. Zhang, and J. Tang, "Ekar: An explainable method for knowledge aware recommendation," *arXiv preprint arXiv: 1906.09506*, 2019.

[151] X. Zhao, L. Zhang, Z. Ding, L. Xia, J. Tang, and D. Yin, "Recommendations with negative feedback via pairwise deep reinforcement learning," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, London, United Kingdom, 2018, pp. 1040–1048.

[152] P. Wang, K. Liu, L. Jiang, X. Li, and Y. Fu, "Incremental mobile user profiling: Reinforcement learning with spatial knowledge graph for modeling event streams," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, Virtual Event, CA, USA, 2020, pp. 853–861.

[153] M. Yousefi, N. Mtetwa, Y. Zhang, and H. Tianfield, "A reinforcement learning approach for attack graph analysis," in *2018 17th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/ 12th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, 2018, pp. 212–217.

[154] R. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec, "Gnnexplainer: Generating explanations for graph neural networks," in *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, 2019.

[155] D. Luo, W. Cheng, D. Xu, W. Yu, B. Zong, H. Chen, and X. Zhang, "Parameterized explainer for graph neural network," *Advances in neural information processing systems*, vol. 33, pp. 19 620–19 631, 2020.

[156] M. Vu and M. T. Thai, "Pgm-explainer: Probabilistic graphical model explanations for graph neural networks," *Advances in neural information processing systems*, vol. 33, pp. 12 225–12 235, 2020.

[157] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and D. Hassabis, "Mastering the game of go without human knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.

[158] D. Bacciu and D. Numeroso, "Explaining deep graph networks via input perturbation," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–12, 2022.

[159] P. Shang, X. Liu, C. Yu, G. Yan, Q. Xiang, and X. Mi, "A new ensemble deep graph reinforcement learning network for spatio-temporal traffic volume forecasting in a freeway network," *Digital Signal Processing*, vol. 123, p. 103419, 2022.

[160] Y. Wang, Y. Tong, C. Long, P. Xu, K. Xu, and W. Lv, "Adaptive dynamic bipartite graph matching: A reinforcement learning approach," in *Proceedings of The IEEE 35th International Conference on Data Engineering (ICDE)*, 2019, pp. 1478–1489.

[161] T. Nishi, K. Otaki, K. Hayakawa, and T. Yoshimura, "Traffic signal control based on reinforcement learning with graph convolutional neural nets," in *Proceedings of the 2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, 2018, pp. 877–883.

[162] H. Zhu, T. Shou, R. Guo, Z. Jiang, Z. Wang, Z. Wang, Z. Yu, W. Zhang, C. Wang, and L. Chen, "Redpacketbike: A graph-based demand modeling and crowd-driven station rebalancing framework for bike sharing systems," *IEEE Transactions on Mobile Computing*, pp. 1–1, 2022.

[163] J. You, B. Liu, R. Ying, V. Pande, and J. Leskovec, "Graph convolutional policy network for goal-directed molecular graph generation," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, Montréal, Canada, 2018, pp. 6412–6422.

[164] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[165] F.-X. Devailly, D. Larocque, and L. Charlin, "Ig-rl: Inductive graph reinforcement learning for massive-scale traffic signal control," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–12, 2021.

[166] Z. Yu and M. Hu, "Deep reinforcement learning with graph representation for vehicle repositioning," *IEEE Transactions on Intelligent Transportation Systems*, pp. 1–14, 2021.

[167] S. Chen, J. Dong, P. Ha, Y. Li, and S. Labi, "Graph neural network and reinforcement learning for multi-agent cooperative control of connected autonomous vehicles," *ComputerAided Civil and Infrastructure Engineering*, vol. 36, no. 7, pp. 838–857, 2021.

[168] H. Zuo, H. Si, N. Ding, X. Wang, and G. Tan, "Coordinated learning for lane changing based on coordination graph and reinforcement learning," in *Artificial Intelligence in China*, 2020, pp. 599–607.

[169] Z. Zeng, "Graphlight: Graph-based reinforcement learning for traffic signal control," in *2021 IEEE 6th International Conference on Computer and Communication Systems (ICCCS)*, 2021, pp. 645–650.

[170] Y. Wang, T. Xu, X. Niu, C. Tan, E. Chen, and H. Xiong, "Stmarl: A spatio-temporal multi-agent reinforcement learning approach for cooperative traffic light control," *IEEE Transactions on Mobile Computing*, pp. 1–1, 2020.

[171] P. Almasan, J. Surez-Varela, A. Badia-Sampera, K. Rusek, P. Barlet-Ros, and A. Cabellos-Aparicio, "Deep reinforcement learning meets graph neural networks: Exploring a routing optimization use case," *arXiv preprint arXiv:1910.07421*, 2019.

[172] K. Nakashima, S. Kamiya, K. Ohtsu, K. Yamamoto, T. Nishio, and M. Morikura, "Deep reinforcement learning-based channel allocation for wireless lans with graph convolutional networks," *IEEE Access*, vol. 8, pp. 31 823–31 834, 2020.

[173] C. Yang, H. Wang, J. Tang, C. Shi, M. Sun, G. Cui, and Z. Liu, "Full-scale information diffusion prediction with reinforced recurrent networks," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–13, 2021.

[174] B. Wang, Y. Sun, T. Q. Duong, L. D. Nguyen, and L. Hanzo, "Risk-aware identification of highly suspected covid-19 cases in social iot: A joint graph theory and reinforcement learning approach," *IEEE Access*, vol. 8, pp. 115 655–115 661, 2020.

[175] C. Fan, L. Zeng, Y. Sun, and Y.-Y. Liu, "Finding key players in complex networks through deep reinforcement learning," *Nature machine intelligence*, vol. 2, no. 6, pp. 317–324, 2020.

[176] L. Ma, Z. Shao, X. Li, Q. Lin, J. Li, V. C. M. Leung, and A. K. Nandi, "Influence maximization in complex networks by using evolutionary deep reinforcement learning," *IEEE Transactions on Emerging Topics in Computational Intelligence*, pp. 1–15, 2022.

[177] D. Yan, W. Xie, and Y. Zhang, "Hypernetwork dismantling via deep reinforcement learning," *arXiv preprint arXiv:2104.14332*, 2021.

[178] J. Yan, S. Yang, and E. Hancock, "Learning for graph matching and related combinatorial optimization problems," in *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, Yokohama, Yokohama, Japan, 2021.

[179] C. K. Joshi, Q. Cappart, L.-M. Rousseau, and T. Laurent, "Learning TSP Requires Rethinking Generalization," in *27th International Conference on Principles and Practice of Constraint Programming (CP 2021)*, vol. 210, Dagstuhl, Germany, 2021, pp. 33:1–33:21.

[180] W. Zheng, D. Wang, and F. Song, "Opengraphgym: a parallel reinforcement learning framework for graph optimization problems," in *International Conference on Computational Science*, 2020, pp. 439–452.

[181] H. Dai, E. B. Khalil, Y. Zhang, B. Dilkina, and L. Song, "Learning combinatorial optimization algorithms over graphs," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2018, pp. 6351–6361.

[182] S. Toyer, F. Trevizan, S. Thiébaux, and L. Xie, "Action schema networks: Generalised policies with deep learning," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in Artificial Intelligence*, New Orleans, Louisiana, USA, 2018.

[183] P. Sun, J. Lan, J. Li, Z. Guo, and Y. Hu, "Combining deep reinforcement learning with graph neural networks for optimal vnf placement," *IEEE Communications Letters*, vol. 25, no. 1, pp. 176–180, 2020.

[184] S. H. Silva, A. Alaeddini, and P. Najafirad, "Temporal graph traversals using reinforcement learning with proximal policy optimization," *IEEE Access*, vol. 8, pp. 63 910–63 922, 2020.

[185] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[186] X. Zhao, J. Wu, H. Peng, A. Beheshti, J. Monaghan, D. McAlpine, H. Hernandez-Perez, M. Dras, Q. Dai, Y. Li *et al.*, "Deep reinforcement learning guided graph neural networks for brain network analysis," *arXiv preprint arXiv:2203.10093*, 2022.

[187] Y. Jia, Z. Tan, and J. Zhang, "Dkdr: An approach of knowledge graph and deep reinforcement learning for disease diagnosis," in *2019 IEEE Intl Conf on Parallel & Distributed Processing with Applications, Big Data & Cloud Computing, Sustainable Computing & Communications, Social Computing & Networking (ISPA/BDCloud/SocialCom/SustainCom)*, 2019, pp. 1303–1308.

[188] Y. Khemchandani, S. OHagan, S. Samanta, N. Swainston, T. J. Roberts, D. Bollegala, and D. B. Kell, "Deepgraphmolgen, a multi-objective, computational strategy for generating molecules with desirable properties: a graph convolution and reinforcement learning approach," *Journal of cheminformatics*, vol. 12, no. 1, pp. 1–17, 2020.

[189] W. Jin, R. Barzilay, and T. Jaakkola, "Multi-objective molecule generation using interpretable substructures," in *International conference on machine learning*, vol. 119, 2020, pp. 4849–4859.

[190] N. De Cao and T. Kipf, "Molgan: An implicit generative model for small molecular graphs," *arXiv preprint arXiv:1805.11973*, 2018.

[191] H. Du, Z. Yan, Q. Xiang, and Q. Zhan, "Vulcan: Solving the steiner tree problem with graph neural networks and deep reinforcement learning," *arXiv preprint arXiv:2111.10810*, 2021.

[192] H. Wang, K. Wang, J. Yang, L. Shen, N. Sun, H. S. Lee, and S. Han, "Gcn-rl circuit designer: Transferable transistor sizing with graph neural networks and reinforcement learning," in *Proceedings of the 57th ACM/IEEE Design Automation Conference (DAC)*, 2020, pp. 1–6.

[193] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine, "Diversity is all you need: Learning skills without a reward function," *arXiv preprint arXiv:1802.06070*, 2018.

[194] O. Nachum, S. S. Gu, H. Lee, and S. Levine, "Data-efficient hierarchical reinforcement learning," *Advances in neural information processing systems*, vol. 31, 2018.

[195] Q. Ma, S. Ge, D. He, D. Thaker, and I. Drori, "Combinatorial optimization by graph pointer networks and hierarchical reinforcement learning," *arXiv preprint arXiv:1911.04936*, 2019.

[196] C. Tessler, S. Givony, T. Zahavy, D. Mankowitz, and S. Mannor, "A deep hierarchical approach to lifelong learning in minecraft," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 31, 2017.

[197] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu, "Feudal networks for hierarchical reinforcement learning," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, 2017, pp. 3540–3549.

[198] P. Peng, Y. Wen, Y. Yang, Q. Yuan, Z. Tang, H. Long, and J. Wang, "Multiagent bidirectionally-coordinated nets: Emergence of human-level coordination in learning to play starcraft combat games," *arXiv preprint arXiv:1703.10069*, 2017.

[199] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," in *Advances in neural information processing systems*, vol. 29, 2016.

[200] S. Sukhbaatar and R. Fergus, "Learning multiagent communication with backpropagation," *Advances in neural information processing systems*, vol. 29, 2016.

[201] D. Chen, M. Nie, J. Yan, D. Wang, and Q. Gan, "Network representation learning algorithm based on complete subgraph folding," *Mathematics*, vol. 10, no. 4, p. 581, 2022.

[202] W. Zhao, Y. Li, T. Fan, and F. Wu, "A novel embedding learning framework for relation completion and recommendation based on graph neural network and multi-task learning," *Soft Computing*, pp. 1–13, 2022.

[203] H. Van Seijen, M. Fatemi, J. Romoff, R. Laroche, T. Barnes, and J. Tsang, "Hybrid reward architecture for reinforcement learning," in *Proceedings of the Advances in Neural Information Processing Systems*, vol. 30, 2017.

[204] H. Kawano, "Hierarchical sub-task decomposition for reinforcement learning of multi-robot delivery mission," in *Proceedings of the IEEE International Conference on Robotics and Automation*, 2013, pp. 828–835.

[205] C. Seshadhri, A. Sharma, A. Stolman, and A. Goel, "The impossibility of low-rank representations for triangle-rich complex networks," *Proceedings of the National Academy of Sciences*, vol. 117, no. 11, pp. 5631–5637, 2020.

[206] P. Ammanabrolu and M. Hausknecht, "Graph constrained reinforcement learning for natural language action spaces," *arXiv preprint arXiv:2001.08837*, 2020.

[207] G. S. Ramachandran, I. Brugere, L. R. Varshney, and C. Xiong, "Gaea: Graph augmentation for equitable access via reinforcement learning," in *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, 2021, pp. 884–894.

[208] Y. Xian, Z. Fu, S. Muthukrishnan, G. de Melo, and Y. Zhang, "Reinforcement knowledge graph reasoning for explainable recommendation," in *Proceedings of the 42nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, Paris, France, 2019, pp. 285–294.

[209] P. Ammanabrolu and M. Riedl, "Playing text-adventure games with graph-based deep reinforcement learning," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 2019, pp. 3557–3565.

**Mingshuo Nie** received the master's degree in 2021 in software engineering from the Software College, Northeastern University, Shenyang, China, where he is currently working toward the Ph.D. degree in software engineering. His current research focuses on graph reinforcement learning, graph mining, link prediction, and graph neural networks.

**Dongming Chen** received the Ph.D. degree in 2006 in computer system architecture from Northeastern University, Shenyang, China, where he is currently a professor at Software College, Northeastern University, China. He is a member of IEEE Computer Society, senior member of China Computer Federation (CCF), and Senior member of China Institute of Communications (CIC). His current research focuses on deep reinforcement Learning, social network and Big Data analysis.

**Dongqi Wang** received the Ph.D. degree in 2011 in computer system architecture from Northeastern University, Shenyang, China, where he is a lecturer in Software College, Northeastern University, China. His current research focuses on network security and social network analysis. He is also a member of CCF.