# Continual Learning on Graphs: Challenges, Solutions, and Opportunities

Xikun Zhang, Dongjin Song*, *Member, IEEE*, and Dacheng Tao*, *Fellow, IEEE*

**Abstract**—Continual learning on graph data has recently attracted paramount attention for its aim to resolve the catastrophic forgetting problem on existing tasks while adapting the sequentially updated model to newly emerged graph tasks. While there have been efforts to summarize progress on continual learning research over Euclidean data, *e.g.*, images and texts, a systematic review of progress in continual learning on graphs, *a.k.a*, continual graph learning (CGL) or lifelong graph learning, is still demanding. Graph data are far more complex in terms of data structures and application scenarios, making CGL task settings, model designs, and applications extremely challenging. To bridge the gap, we provide a comprehensive review of existing continual graph learning (CGL) algorithms by elucidating the different task settings and categorizing the existing methods based on their characteristics. We compare the CGL methods with traditional continual learning techniques and analyze the applicability of the traditional continual learning techniques to CGL tasks. Additionally, we review the benchmark works that are crucial to CGL research. Finally, we discuss the remaining challenges and propose several future directions. We will maintain an up-to-date GitHub repository featuring a comprehensive list of CGL algorithms, accessible at https://github.com/UConn-DSIS/Survey-of-Continual-Learning-on-Graphs.

**Index Terms**—Graph representation learning, continual learning, graph neural networks, continual graph learning, CGL, lifelong learning, lifelong graph learning, continual graph representation learning.

✦

## 1 INTRODUCTION

IN the field of graph representation learning, traditional methods typically assume static graphs, *i.e.*, the structure of graphs as well as their attributed nodes and edges remain constant. However, in many real-world applications, including both node-level and graph-level scenarios, graphs may evolve constantly. This means that either new types of nodes and their associated edges may appear or new types of graphs may be collected. In this case, the graph data may exhibit distribution shift, and a desired model is expected to continually accommodate the new distributions without forgetting the previously learnt knowledge. For example, new categories of research papers (graph nodes) and the accompanying citations (graph edges) will constantly emerge in a citation network. To automatically classify the papers, a document classifier is expected to continually adapt to the distribution of new categories while maintaining the learnt knowledge of the previously observed categories at the same time [1], [2], [3]. For drug discovery, new molecule properties and new molecule categories may be encountered intermittently, and a molecule property predictor has to fit its parameters for new patterns without compromising their prediction performance over existing molecule categories or properties [1], [2]. In such continual learning scenarios, naively adapting the model to new tasks will encounter the catastrophic forgetting problem, *i.e.*, drastic performance drop on existing tasks after the model parameters are adapted for new tasks.

A naive solution is to retrain the model over all previously observed data whenever a new task emerges. However, this could be infeasible either due to the intractable retraining cost or potential privacy/legal regulation issues. In light of this, Continual Graph Learning (CGL), which aims to continually learn new tasks without forgetting previously learnt knowledge, has recently received increasingly more attention from various areas [1], [2], [4], [5], [6], [7], [8], [9], [10], [11], [12], [13], [14], [15], [16], [17], [18], [19], [20], [21], [22], [23]. Due to the complexity of graph data, existing CGL methods are highly heterogeneous in terms of the targeted graph types, learning settings, basic techniques, and evaluation metrics. First, existing methods may propose general techniques or specialized techniques for certain application scenarios (*e.g.* knowledge graph, recommender system, etc.). Second, according to the availability of the task identity during testing, CGL methods may adopt different settings including task-incremental learning (task-IL), domain-incremental learning (domain-IL), and class-incremental learning (class-IL). Besides the incremental setting, the graph learning tasks may also focus on different granularity, including node-level tasks and graph-level tasks. Third, the methodologies employed by various existing approaches align with distinct prevailing streams, such as regularization, memory replay, and parameter isolation. Finally, different from standard learning settings, the performance of a CGL model is concerned with different perspectives including the overall performance, the performance decrease (forgetting), the inter-task interfere, *etc*. Therefore, depending on the research target, different works may adopt different metrics for evaluating the models. These orthogonal dimensions of CGL research present complex challenges and significant obstacles for researchers attempting to enter the field.

- *Mr X. Zhang and Prof D. Tao are with the School of Computer Science, in the Faculty of Engineering, at The University of Sydney, 6 Cleveland St, Darlington, NSW 2008, Australia. Email: xzha0505@uni.sydney.edu.au, dacheng.tao@sydney.edu.au.*
- *Dr. D. Song is with the Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, the United States. Email: dongjin.song@uconn.edu.*
  *\* indicates corresponding authors.*

To this end, in this paper, we provide a systematic review of the existing methods from the aforementioned four different perspectives. In Section 2, we briefly introduce the context information of continual learning, CGL, as well as the difference between CGL and other related research areas on dynamic evolving graphs. In Section 3, we explain the problem setup of CGL from the perspectives of basic concepts, task sequence construction, task granularity, and different incremental learning scenarios. In Section 4, we first systematically review the related CGL works following the three mainstreams, including the regularization based, memory-replay based, and parameter-isolation based methods, and then analyze the applicability of traditional continual learning techniques. In Section 6, we introduce the recently proposed benchmark works that aim to provide a consistent experimental setting and fair platform for comparing different techniques. Finally, in Section 7, we discuss existing challenges and highlight the promising future directions for CGL research.

## 2 BACKGROUND

### 2.1 Continual Learning

Continual learning (CL) [1], [24], [25], [26], [27] targets the catastrophic forgetting problem, *i.e.* the phenomenon that a model's performance on previous tasks decreases drastically after learning the subsequent new tasks. During training, the learning can be formulated as training a model consecutively on a sequence of tasks. When learning each task, the model can only access the data of the current task, while access to previous tasks is usually not allowed. In continual learning, the data of different tasks typically exhibit different distributions, *e.g.* images of different classes. Accordingly, after learning a new task, if the model parameters are adapted solely for the new task, the learnt knowledge of the previous tasks may be overwritten, leading to the catastrophic forgetting problem over previous tasks. During testing, the model will be tested on each learnt task. A desired model should be able to maintain its performance on every learnt task. The evaluation of continual learning models is detailed in Section 5.

### 2.2 Graph Representation Learning

Graph representation learning [28], [29], [30], [31], [32], [33], [34], [35], [36], [37] aims to generate qualified representations for nodes, edges, or the entire graphs, which can further be used for downstream tasks like node classification [28], [38], link prediction [39], [40], [41], graph classification [42], [43], [44], *etc..* Graph neural networks (GNNs) are currently the prevailing approaches. To obtain either node/edge level or graph level representations, GNNs would first generate node/edge level representations. The most popular GNNs follow the message passing neural network (MPNN) framework [42], and can be formulated as,

$$\mathbf{m}_v^{l+1} = \sum_{u \in \mathcal{N}^1(v)} \mathrm{M}_l(\mathbf{h}_v^l, \mathbf{h}_u^l, \mathbf{x}_{v,u}^e; \boldsymbol{\theta}_l^{\mathrm{M}}), \qquad (1)$$

$$\mathbf{h}_v^{l+1} = \mathrm{U}_l(\mathbf{h}_v^l, \mathbf{m}_v^{l+1}; \boldsymbol{\theta}_l^{\mathrm{U}}), \qquad (2)$$

As shown in the formulations above, different from learning on independent data (*e.g.* images), generating representations on graph data requires properly capturing the valuable topological information (*e.g.* through the message passing). Similarly, compared to classic CL on independent data, CGL also has to properly preserve the highly valuable topological information. Besides, the topological connections may cause some continual learning techniques to be inapplicable. For example, to generate the representation of a single node, the message passing based GNNs would require aggregating information from multi-hop neighbors, and the methods based on memory replay, which rely on storing individual data points, become untenable.

### 2.3 Differentiation from Other Related Works

In addition to CGL efforts aimed at mitigating the issue of forgetting, there exists research at the intersection of graph representation learning and evolving graphs that primarily addresses alternative challenges [45], [46], [47], [48], [49], [50], [51], [52], [53], [54], [55], which could be confusing. In this subsection, we clarify the crucial difference between CGL and some easily confusing topics such as dynamic graph learning [56], [57], [58], [59], [60], [61], [62], [63], [64], [65] and few-shot graph learning [66], [67], [68], [68], [69], [70]. Dynamic graph learning primarily aims to capture the evolving graph structure and maintain up-to-date graph representations, with access to all prior information, rather than tackling the forgetting issue. Conversely, CGL focuses on resolving the forgetting problem, and the data of previous tasks are typically unavailable. An exception is CGL with inter-task edges, which permits aggregation of past task information via inter-task edges during the GNNs' neighborhood aggregation process. Nevertheless, the labels of prior task data remain inaccessible. Few-shot graph learning is designed for rapid model adaptation to new tasks. During training, few-shot learning models can access all tasks simultaneously, which is not the case for CGL (CGL with inter-task edges deviates slightly). During evaluation, few-shot learning models are tested on new tasks after initial fine-tuning, whereas CGL models are evaluated on existing tasks without any fine-tuning.

Finally, this paper differs from an earlier survey [71] on CGL in several key aspects. Febrinanto *et al.* [71] emerged during the nascent stages of Continual Graph Learning (CGL) research, incorporating a limited scope of studies due to the field's infancy. In contrast, our survey, benefiting from over a year of rapid advance in CGL [72], [73], [74], [75] and the deep understanding of the CGL problem based on our previous technical contributions in this field [1], [76], [77], [78], [79], encompasses a comprehensive review of state-of-the-art techniques. Our survey not only elucidates technical approaches but also offers a systematic exploration of CGL's problem setup. In addition, we also examine different methods' applicability across diverse application scenarios, as summarized in Table 1. Finally, our survey adopts a more analytical perspective. We integrate the discussed CGL research within a unified mathematical framework, providing a detailed exposition of CGL settings and methodologies with an emphasis on technical precision.
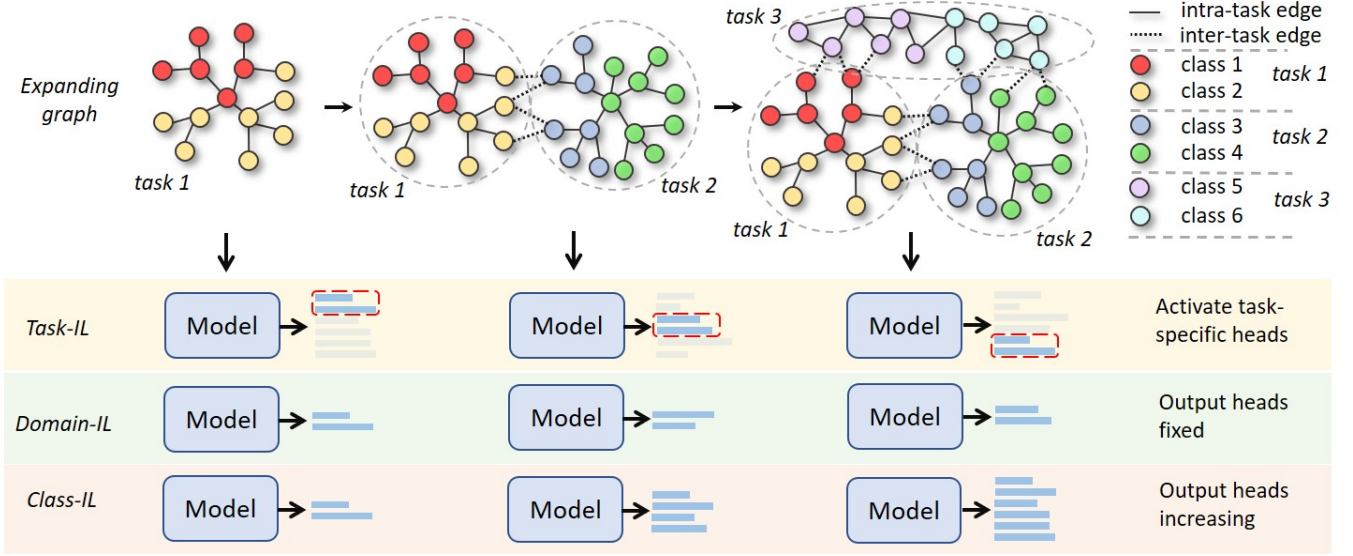
Fig. 1: Illustration of the different incremental settings.

## 3 PROBLEM SETUPS

To cover different graph evolving scenarios, the continual learning process can be formulated on a sequence of graph tasks: $\mathcal{S} = \{\mathcal{G}_1, \mathcal{G}_2, ..., \mathcal{G}_T\}$. Each graph task $\mathcal{G}_\tau$ represents the entire graph (or a set of independent graphs) that has grown (or increased) from task 1 to $\tau$. Each $\mathcal{G}_\tau$ consists of a node set $\mathbb{V}_\tau$ containing all the nodes, and an edge set $\mathbb{E}_\tau$ containing all the edges. $\mathbb{E}_\tau$ can also be represented as an adjacency matrix $\mathbf{A}_\tau \in \mathbb{R}^{|\mathbb{V}_\tau| \times |\mathbb{V}_\tau|}$ (connected graphs for node-level CGL task and disconnected independent graphs for graph-level CGL tasks). Each entry $\mathbf{A}_\tau^{u,v}$ denotes the edge between node $u$ and $v$ ($\mathbf{A}_\tau^{u,v} = 0$ if node $u$ and $v$ are not connected by an edge). The number of edges connected to a node is referred to as its degree, and the degree of all nodes within a graph can be stored as the diagonal entries of a degree matrix $\mathbf{D}_\tau \in \mathbb{R}^{|\mathbb{V}_\tau| \times |\mathbb{V}_\tau|}$, $i.e.$ $\mathbf{D}_\tau^{u,u}$ is the degree of node $u$. In practice, $\mathbf{D}_\tau$ is often used to normalize the adjacency matrix, $i.e.$ $\hat{\mathbf{A}}_\tau = \mathbf{D}_\tau^{-\frac{1}{2}} \mathbf{A}_\tau \mathbf{D}_\tau^{-\frac{1}{2}}$.

To conveniently refer to the nodes coming in different tasks, we denote the new nodes in a task $\tau$ as $\mathbb{V}_\tau^{new}$. Accordingly, we denote the induced subgraph based on $\mathbb{V}_\tau^{new}$, $i.e.$ all nodes in $\mathbb{V}_\tau^{new}$ together with all edges that connect nodes in $\mathbb{V}_\tau^{new}$, as $\mathcal{G}_\tau^{new}$. $\mathcal{G}_\tau^{new}$ is the task-specific subgraph of the new task $\tau$. Then, all the edges within each task specific subgraph are referred to as intra-task edges [1], while the other edges are inter-task edges connecting nodes across different tasks. We denote the intra-task edges of task $\tau$ as $\mathbb{E}_\tau^{intra}$, and the inter-task edges after the arrival of task $\tau$ as $\mathbb{E}_\tau^{inter} = \mathbb{E}_\tau \setminus \bigcup_{i=1}^\tau \mathbb{E}_\tau^{intra}$.

The formulation in this subsection applies to both node/edge-level CGL and graph-level CGL, which will be detailed in Section 3.1. For node/edge-level CGL, which typically involves expanding networks ($e.g.$ citation networks), each graph $\mathcal{G}_\tau$ is the entire network that has grown from task 1 to $\tau$. While for graph-level CGL, which typically involves independent graphs ($e.g.$ molecule graphs), each graph $\mathcal{G}_\tau$ is a set of disconnected independent graphs.

### 3.1 Different Prediction Granularity

Different from learning on independent data, graph related tasks have different granularity that corresponds to different real-world scenarios and requires different task sequence constructions.

Node/edge-level learning focuses on generating representations and predictions for individual nodes/edges. For example, given a citation network, classifying the papers is modeled as a node classification task (node-level). With a protein interaction network, in which the nodes are proteins and the edges are interactions between proteins, predicting the missing interactions is modeled as a link prediction (edge-level) task. Since nodes and edges belong to the same granularity and edge-related tasks are often based on the learnt node representation, we will use node-level to denote both node- and edge-level tasks in the following for simplicity.

Formally, depending on whether inter-task edges are allowed to be preserved, node-level tasks can have two different formulations. First, without the inter-task edges, the representation of a task $\tau$ is:

$$\{\mathbf{h}_v | v \in \mathbb{V}_\tau^{new}\} = \mathrm{f}(\mathcal{G}_\tau^{new}; \boldsymbol{\theta}_\tau), \qquad (3)$$

where $\mathbf{h}_v$ is the representation of node $v$ generated by the model $\mathrm{f}(\cdot; \boldsymbol{\theta}_\tau)$ parameterized by $\boldsymbol{\theta}_\tau$ (the model parameters during learning task $\tau$). When the inter-task edges are preserved, then the input for learning task $\tau$ is not limited to $\mathcal{G}_\tau^{new}$. Since GNNs could aggregate information from previously observed nodes via the inter-tasks, the input becomes the entire graph $\mathcal{G}_\tau$, $i.e.$,

$$\{\mathbf{h}_v | v \in \mathbb{V}_\tau^{new}\} = \mathrm{f}(\mathcal{G}_\tau; \boldsymbol{\theta}_\tau). \qquad (4)$$

When dealing with node-level tasks, the incoming subgraph $\mathcal{G}_\tau^{new}$ is typically a connected graph. While for graph-level tasks, $\mathcal{G}_\tau^{new}$ is a set of disconnected graphs, $e.g.$, a set of molecule graphs. In this case, the graph-level representations can be formulated as:

$$\{\mathbf{h}_g | g \in \mathcal{G}_\tau^{new}\} = \mathrm{f}(\mathcal{G}_\tau^{new}; \boldsymbol{\theta}_\tau), \qquad (5)$$
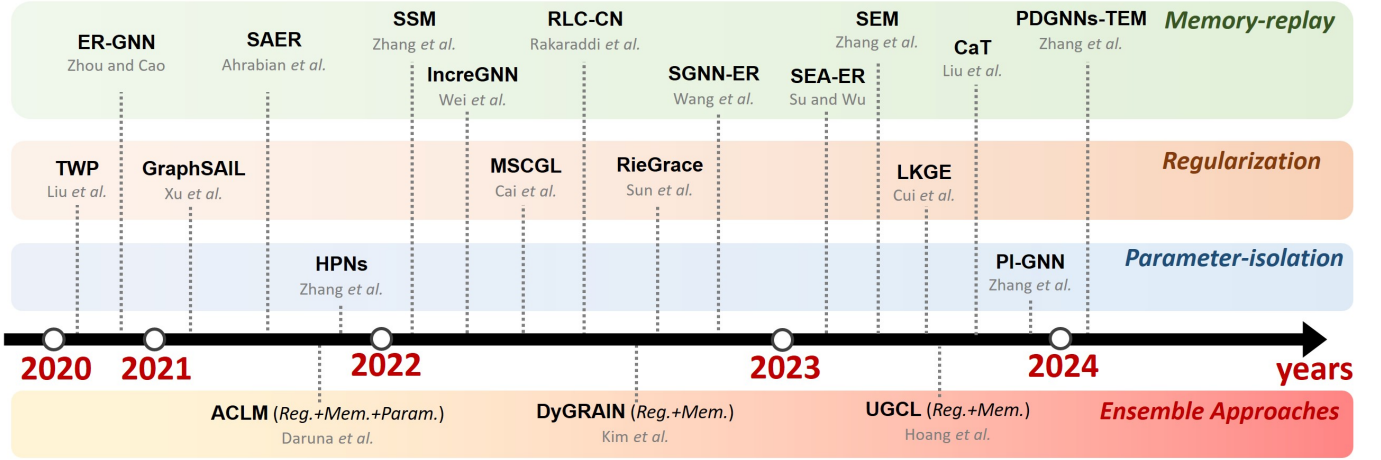
Fig. 2: Timeline of the development of different CGL techniques.

where $g$ denotes the disconnected components in $\mathcal{G}_\tau^{new}$. Since graph-level CGL deals with individual graphs (disconnected components), there is no inter-task edge across different tasks.

### 3.2 Different Incremental Scenarios

According to whether the task identities are provided during testing and whether a model is required to figure out the task identities, CGL, similar to classic continual learning, can be categorized into task-incremental learning (task-IL), domain-incremental learning (domain-IL), and class-incremental learning (class-IL), as illustrated in Figure 1.

#### 3.2.1 Task-IL

In task-IL, the task identities are revealed to the model during testing, therefore the model is not required to identify the given tasks. For classification tasks, existing models typically increase their output dimensions to accommodate new tasks (*e.g.*, new classes), and only activate the corresponding dimensions for each given task during testing. For example, in molecular property prediction tasks, each new task could be predicting whether a given molecule exhibits a certain new property.

#### 3.2.2 Domain-IL

Domain-IL refers to a scenario where the domain of the data changes across tasks, while the task remains static. Therefore, the semantic meaning of a model's output dimensions is fixed. For example, one possible scenario for continual learning on knowledge graphs is to sequentially learn on graphs with different entities and relations, while the prediction task is always the completion of the triplets. Another example is the data split according to the temporal orders. The task could be the same across different temporal periods, while different temporal periods may contain data following significantly diverse regimes (domains).

#### 3.2.3 Class-IL

Class-IL is the most challenging among the three scenarios. During testing, task identities are inaccessible and the model has to identify the given tasks. For classification tasks, a
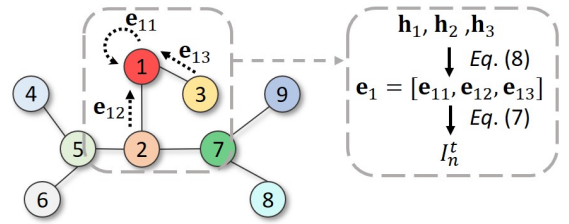


Fig. 3: Topology preservation design in TWP [2].

model typically increases the output dimensions to allocate new classes when they arrive and has to pick up the correct class among all learnt classes, unlike task-IL which only requires distinguishing between classes within a known task.

## 4 METHODS

Similar to traditional continual learning, CGL also approaches the problem from the perspectives of restricting the change in the model parameters, isolating and protecting the parameters that are important for previously learnt tasks, and replaying representative data from previous tasks to remind the model of the previously learnt patterns (Figure 2). However, a key challenge of CGL is the necessity to properly preserve the topological structure of the data, which is crucial information contained in the graph data. In this section, we will introduce the CGL techniques following this categorization, which is also summarized in Table 1.

### 4.1 Regularization based methods

Since the reason of the forgetting is that the model parameters trained for previous tasks are modified after being adapted to new tasks, traditional regularization based methods [91], [92] add penalty terms to prevent the parameters from being drastically changed. However, these methods do not explicitly preserve the topology of the graph data. Targeting this insufficiency, topology-aware weight preserving (TWP) [2] proposes to explicitly preserve the topology learnt on previous tasks via regularization on the model weights.

TABLE 1: Summarization of different CGL techniques.

| Method | Applications | Task Granularity | Technique | Characteristics |
|---|---|---|---|---|
| TWP [2] | General | Node/Graph | Reg. | Preserve the topology learnt from previous tasks |
| RieGrace [75] | General | Node | Reg. | Maintain previous knowledge via knowledge distillation |
| GraphSAIL [80] | Recommender Systems | Node | Reg. | Local and global structure preservation, node information preservation |
| MSCGL [81] | General | Node | Reg. | Parameter changes orthogonal to previous parameters |
| LKGE [82] | Knowledge Graph | Node | Reg. | Alleviating forgetting issue with l2 regularization |
| ER-GNN [83] | General | Node | Mem. | Replay representative nodes |
| SSM [77] | General | Node | Mem. | Replay representative sparsified computation subgraphs |
| SEM [78] | General | Node | Mem. | Sparsify computation subgraphs based on information bottleneck |
| PDGNNs-TEM [79] | General | Node | Mem. | Replay representative topology-aware embeddings |
| IncreGNN [84] | General | Node | Mem. | Replay nodes according to their influence |
| RLC-CN [85] | General | Node | Mem. | Model structure adaption and dark experience replay |
| SGNN-ER [86] | General | Node | Mem. | Model retraining with generated fake historical data |
| SAER [87] | Recommender System | Node | Mem. | Buffer the representative user-item pairs based on reservoir sampling |
| SEA-ER [88] | General | Node | Mem. | Minimize the structural difference between the memory buffer and the original graph |
| CaT [73] | General | Node | Mem. | Train the model solely on balanced condensed graphs from all tasks |
| HPNs [3] | General | Node | Para. | Extracting and storing basic features to encourage knowledge sharing across tasks, model expanding to accommodate new patterns |
| PI-GNN [74] | General | Node | Para. | Separate parameters for encoding stable and changed graph parts |
| DyGRAIN [89] | General | Node | Mem.+Reg. | Alleviate catastrophic forgetting and concept shift of previous task nodes via memory replay and knowledge distillation |
| ACLM [90] | Knowledge Graph | Node | Mem.+Reg.+Para. | Adapting general CL techniques to CGL tasks |
| UGCL [72] | General | Node/Graph | Mem.+Reg. | Memory replay & local/global structure preservation |

Denoting the model parameters after learning the $t$-th task as $\boldsymbol{\theta}_t$, the learning of regularization based methods on the $(t+1)$-th can be generally formulated as,

$$\boldsymbol{\theta}_{t+1} = \arg\min_{\boldsymbol{\theta}} \mathcal{L}_{t+1} + \sum_{n=1}^{t} I_n \otimes (\boldsymbol{\theta} - \boldsymbol{\theta}_n^*)^2. \quad (6)$$

In Equation 6, $\boldsymbol{\theta}_n^*$ is the optimal parameters for the $n$-th task, and the importance scores $I_n$ indicate the importance of each parameter for the performance in task $n$. In other words, the change of the parameters that are important to previous tasks are strongly penalized, while the less important parameters are more free to adapt to the new tasks. In different methods, $I_n$ are calculated with different strategies. In TWP, the importance scores consist of two parts. The first part is calculated for each weight of the model based on the sensitivity of the loss on the weight, which is measured with the gradient of the loss with respect to the weight. This part resembles the strategies adopted by EWC [91] and MAS [92], and serves to preserve the weights that are important to previous tasks. We denote the corresponding importance scores as $I_n^p$. The second part, which is the key contribution of TWP, is the sensitivity of the learnt graph structure on the weight. For backbones like graph attention network (GAT) that explicitly learns the strength of the edges via the attention scores, the learnt

structures are deemed as the attention score on each weight, and the sensitivity is measured as the gradient of the squared $l_2$ norm of the attention scores with respect to the weights. Denoting the attention score between two nodes $i$ and $j$ as $e_{ij}$, the attention score between a node $i$ and its all neighbors can be represented as $\mathbf{e}_i = [e_{i1}, ..., e_{i|\mathcal{N}_i|}]$. After that, the sensitivity (importance) of the topology learnt in task $n$ with respect to the model parameters is defined as the gradient of the squared $l2$ norm with respect to the parameters (Figure 3).

$$I_n^t = \frac{\partial \left(\left|[\mathbf{e}_1, ..., \mathbf{e}_{|V^n|}]\right|\right)}{\partial \boldsymbol{\theta}} \quad (7)$$

For backbones without attention mechanism, the attention score between two nodes is calculated with a non-parametric attention mechanism,

$$e_{ij} = \mathbf{h}_i^T \tanh(\mathbf{h}_j) \quad (8)$$

With the importance scores for both the model performance and the topology, the final importance is obtained by their weighted summation,

$$I_n = \lambda_p I_n^p + \lambda_t I_n^t. \quad (9)$$

Besides the learnt attention scores, the graph curvature is another important topological property of the graph data, and is crucial for graph learning [93], [94]. When learning

on a sequence of graphs, considering that the curvature of the incoming graph may constantly change, Riemannian Graph Continual Learner (RieGrace) [75] further proposes Adaptive Riemannian GCN (AdaRGCN) for accommodating the change in the graph curvature and Label-free Lorentz Distillation for alleviating the forgetting problem. Specifically, AdaRGCN can constantly adapt and capture the topology of the given graphs with different curvatures. And the Label-free Lorentz Distillation serves to maintaining the learnt knowledge when learning on the graph sequence. Considering the possible scarcity of labels in real worlds applications, Label-free Lorentz Distillation trains a model based on contrastive learning without any label. It consists of two parts, including a cross-layer intra-distillation and a cross-model inter-distillation. The intra-distillation serves to learn the pattern of the current graph, and the contrastive learning is conducted by maximizing the agreement between the representations of each node at different layers of the model. Formally, the objective function is

$$\mathcal{J}(\mathbf{x}_i^{s,L}, \mathbf{x}_i^{s,H}) = \log \frac{\exp \mathrm{Sim}^{\mathcal{L}}(\mathbf{x}_i^{s,L}, \mathbf{x}_i^{s,H})}{\sum_{j=1}^{|\mathcal{V}|} \mathbb{I}\{i \neq j\} \exp \mathrm{Sim}(\mathbf{x}_i^{s,L}, \mathbf{x}_i^{s,H})}. \tag{10}$$

In Equation 10, $\mathbf{x}_i^{s,L}$ and $\mathbf{x}_i^{s,H}$ denote the low- and high-level representations of the $i$-th node, extracted by shallow and deep layers of the model. Since these two representations correspond to the same node, the objective is designed to maximize their similarity, which is measured by the adopted function $\mathrm{Sim}^{\mathcal{L}}(\cdot, \cdot)$. And $\mathcal{V}$ is the node set of the entire graph. Besides learning the pattern of the current graph, the knowledge learnt from the previous tasks should also be distilled into the current model, which is achieved by the inter-distillation loss,

$$\mathcal{J}(\mathbf{x}_i^{t,H}, \mathbf{x}_i^{s,H}) = \log \frac{\exp \mathrm{Sim}^{\mathcal{L}}(\mathbf{x}_i^{t,H}, \mathbf{x}_i^{s,H})}{\sum_{j=1}^{|\mathcal{V}|} \mathbb{I}\{i \neq j\} \exp \mathrm{Sim}(\mathbf{x}_i^{t,H}, \mathbf{x}_i^{s,H})}. \tag{11}$$

Equation 11 has the same formulation as Equation 10, while the difference is that the contrastive learning is between the high-level representation of the teacher model and the student model. These two loss on each node are then summed up and balanced with a parameter $\lambda$, $\mathcal{J}_{overall} = \sum_{i=1}^{|\mathcal{V}|} \mathcal{J}(\mathbf{x}_i^{s,L}, \mathbf{x}_i^{s,H}) + \lambda \sum_{i=1}^{|\mathcal{V}|} \mathcal{J}(\mathbf{x}_i^{t,H}, \mathbf{x}_i^{s,H})$. By optimizing these two losses simultaneously, the current student model contains knowledge from both the current graph and the previously learnt model, and will serve as the teacher model when learning the next task.

Besides RieGrace, distillation is also adopted by Universal Graph Continual Learning (UGCL) [72]. UGCL is a technique based on both memory replay and knowledge distillation. The memory replay module is constructed by storing representative nodes as well as the associated neighborhood in a buffer, which are then replayed to the model when learning on new tasks to avoid forgetting. After learning each task, the memory buffer will be updated with two steps. First, all of the observed data of the current task will be added into the buffer. If the resulting buffer size exceeds the budget, a replacement process will be triggered. The replacement process will iteratively remove one random example from

the largest class in the buffer until the buffer size is not larger than the budget. The knowledge distillation module of UGCL is designed to preserve the learned topological information from both local and global perspectives. For local structure, UGCL adopts the strategy to preserve the difference between a node and its neighborhood. Denoting the representation of a node $v$ generated by the current/old models as $\mathbf{z}_v^c/\mathbf{z}_v^o$, the local structure between a node $v$ and its neighbourhood $\mathcal{N}(v)$ in a buffered graph $g$ is formulated as,

$$S_{(g,v)}^c = \mathbf{z}_v^c - \frac{1}{|\mathcal{N}(v)|} \sum_{w \in \mathcal{N}(v)} \mathbf{z}_w^c \tag{12}$$

$$S_{(g,v)}^o = \mathbf{z}_v^o - \frac{1}{|\mathcal{N}(v)|} \sum_{w \in \mathcal{N}(v)} \mathbf{z}_w^o. \tag{13}$$

Based on the local structure difference scores of both the current and old models ($S_{(g,v)}^c$, $S_{(g,v)}^c$), the objective for preserving the local structure difference is,

$$\mathcal{L}^{LS} = \frac{1}{K_g \cdot K_n} \sum_{j=1}^{K_g} \sum_{i=1}^{K_v} [1 - \cos(S_{(j,i)}^c, S_{(j,i)}^o)], \tag{14}$$

where $K_g$ and $K_v$ denote the number of buffered graphs and the number of nodes selected from the buffered graphs. Finally, to preserve the global structure, the objective is formulated as minimizing the cosine distance between the graph-level representations output from the current and old models,

$$\mathcal{L}^{GS} = \frac{1}{K_g} \sum_{j=1}^{K_g} [1 - \cos(\mathbf{g}_j^o, \mathbf{g}_j^c)], \tag{15}$$

where $\mathbf{g}_j^o/\mathbf{g}_j^c$ denotes the representations of graph $j$ generated by the current and old model.

The methods introduced above mainly focus on the forgetting problem. However, different from traditional continual learning, in which the new and old data are independent, the newly incoming data (graph nodes) in CGL could form edges connecting to the existing nodes. Since GNNs generate the prediction for a node based on its multi-hop neighbors, the newly formed connections between new and old nodes will cause concept drift for the old nodes and alter their representation [1]. Therefore, targeting both this challenge and the catastrophic forgetting problem, DyGRAIN [89] is designed to first identify the nodes that are most influenced by new nodes and the nodes that are most vulnerable to forgetting. After that, model retraining and knowledge distillation are applied to alleviate the influence of the new nodes and the forgetting problem. For the detection of the vulnerable nodes, the changing neighborhood problem is termed as the time-varying receptive field in the DyGRAIN paper. Based on the neighborhood aggregation of GNNs, the nodes most influenced by the changed receptive field are detected based on the adjacency matrix to conduct the neighborhood aggregation. Specifically, denoting the previously observed nodes as $\mathcal{M}$, the nodes in the $t$-th task as $\mathcal{V}^{(t)}$, an indicator matrix $\mathbf{V}$ can be derived to represent the influence of each node after certain rounds of message passing.

$$\mathbf{V}_{ii}^{(0)} = \begin{cases} 1 & if \quad i \in \mathcal{V}^{(t)} \\ 0 & if \quad i \in \mathcal{M}, \end{cases} \quad \mathbf{V}^{(l)} = \mathbf{P}^{(t)} \mathbf{V}^{(l-1)}, \tag{16}$$

where $\mathbf{P}^{(t)} = \mathbf{A}^{(t)}\mathbf{D}^{-1}$ is the adjacency matrix $\mathbf{A}^{(t)}$ normalized by the degree matrix $\mathbf{D}$. By iteratively applying $\mathbf{P}^{(t)}$ for $L$ times, the entries in $\mathbf{V}^{(L)}$ denote how much information (*i.e.* influence) would be propagated from each node to the $L$-hop neighbors. Accordingly, a Structural Influence (SI) score of a node $i$ is defined as:

$$\pi_{SI}(i) = \sum_{j \in \mathcal{M} \cup \mathcal{V}^{(t)}} \sum_{l=1}^{L} \mathbf{V}_{ij}^{(l)}. \tag{17}$$

Based on the structural influence score defined above, a subset of $k$ nodes with the highest scores (*i.e.* $\mathcal{S}_{struct}$) are selected as the ones that are most vulnerable to the influence of the new nodes. Besides the purely structural influence, the authors also consider the influence based on the change in the node embeddings. Denoting the GNN model as $f_{\boldsymbol{\theta}}(\cdot)$, the representation of a node $i$ after the arrival of task $t$ can be denoted as, $\mathbf{h}_i^{(t)} = f_{\boldsymbol{\theta}}(\mathbf{x}_i, \mathbf{A}^{(t)}; \boldsymbol{\theta}^{(t-1)})$. And the influence of the $t$-th task is derived based on the change in the node embeddings and defined as the Feature Influence (FI) score,

$$\pi_{FI}(v_i) = 1 - \frac{\mathbf{h}_i^{(t-1)} \cdot \mathbf{h}_i^{(t)}}{||\mathbf{h}_i^{(t-1)}|| \cdot ||\mathbf{h}_i^{(t)}||}. \tag{18}$$

Similarly, with the obtained feature influence score, the top $k$ nodes, *i.e.* $\mathcal{S}_{feat}$, are also selected as the vulnerable ones. Above all, to mitigate the influence from the time-varying receptive field, the union of the two sets of nodes, *i.e.* $\mathcal{S}_{IP} = \mathcal{S}_{struct} \cup \mathcal{S}_{feat}$, will be included in the memory for retraining the model.

Similar to the procedure to mitigate the influence of time-varying receptive field, the catastrophic forgetting is also dealt with by first identifying the most vulnerable nodes. Since the nodes within the $L$-hop neighbors are already sampled, they will be excluded during the identification of the nodes that are vulnerable to catastrophic forgetting, and the remaining node is denoted as $\mathcal{V}_C$. With $\mathcal{V}_C$, the loss on the previous task is,

$$\mathcal{L}^{(t-1)} = l(\mathcal{V}_C, \tilde{A}^{(t-1)}; f_{\boldsymbol{\theta}^{(t-1)}}). \tag{19}$$

Since lower loss values correspond to the more representative nodes (the model is well adapted for them), from $\mathcal{V}_C$, a subset of nodes is selected by lower $q$ percentile according to the loss value, $\mathcal{V}' = \{i \in \mathcal{V}_C | \mathcal{L}_i^{(t-1)} \leq P_q(\mathcal{L}^{(t-1)})\}$. To find the most vulnerable nodes, the first step is to estimate the new loss in task $t$ of the nodes $\mathcal{V}'$ as,

$$\hat{\mathcal{L}^{(t)}} \approx l(\mathcal{V}', \tilde{\mathbf{A}}^t; f_{\boldsymbol{\theta}^{(t-1)}}). \tag{20}$$

Then the loss-based importance score of a node $i$ is defined as the increase in the loss,

$$\pi_{CF}(i) = \hat{\mathcal{L}}_i^{(t)} - \mathcal{L}_i^{(t-1)}, \tag{21}$$

and the top $k$ nodes are the ones that with the largest loss and are selected as the most vulnerable ones to the new nodes, *i.e.* $\mathcal{S}_{KD}$. The role of $\mathcal{S}_{KD}$ is two fold. First, knowledge distillation is conducted to protect the performance of these vulnerable nodes by minimizing the discrepancy between the node representations generated by the new and old models. Second, $\mathcal{S}_{KD}$ is also used for retraining the model together with the nodes in $\mathcal{S}_{IP}$. Above all, by retraining the model with the nodes in $\mathcal{S}_{IP}$, DyGRAIN is capable of ensuring that

the old nodes whose neighborhood (receptive field) is altered by new nodes with up-to-date representations. Through the knowledge distillation and model retraining with nodes in $\mathcal{S}_{KD}$, the catastrophic forgetting can also be alleviated.

Unlike the grid data, *e.g.* images, videos, and text, graph data has tremendous different branches with significantly different properties. Accordingly, besides the above works for general CGL on any graph data, there are also methods specialized for a certain kind of graphs. For example, specially designed for incremental learning for recommender systems, Graph Structure Aware Incremental Learning (GraphSAIL) designs three knowledge distillation techniques to avoid the forgetting problem, including local structure distillation, global structure distillation, and self-embedding distillation. Since the nuclear operation of GCNs is the neighborhood aggregation over the local neighbors, preserving the local contextual structure is crucial. In GraphSAIL, this local structure is reflected through the affinity scores between the center nodes and their neighbors. And the affinity is maintained by knowledge distillation between the new and old tasks. Mathematically, it can be formulated as,

$$\mathcal{L}_{local} = \Big( \frac{1}{|\mathcal{U}| \sum_{u \in \mathcal{U}} (\mathbf{e}_u^{(t-1)} \cdot \mathbf{c}_{u,N_u^{t-1}} - \mathbf{e}_u^t \cdot \mathbf{c}_{u,N_u^{t-1}}^t)^2} + \frac{1}{|I|} \sum_{i \in \mathcal{I}} (\mathbf{e}_i^{t-1} \cdot \mathbf{c}_{i,N_i^{t-1}}^t)^2 \Big), \tag{22}$$

where $\mathcal{U}$ and $\mathcal{I}$ are the user set and item set, $\mathbf{e}_u^t$ is the embedding of the user $u$ during learning task $t$, and $c$ denotes the average of the neighboring user nodes or item nodes,

$$\mathbf{c}_{u,N_u^{t-1}} = \frac{1}{|\mathcal{N}_u^{(t-1)}|} \sum_{i' \in N_u^{t-1}} \mathbf{e}_{i'}^t, \tag{23}$$

$$\mathbf{c}_{i,N_u^{t-1}} = \frac{1}{|\mathcal{N}_i^{(t-1)}|} \sum_{u' \in N_i^{t-1}} \mathbf{e}_{u'}^t. \tag{24}$$

Besides the local structure denoting the context of a node, the global position of a node within the graph is also crucial. For example, the distance among the users can reflect certain preference groups. Therefore, a global structure distillation is also designed. Specifically, to capture the global structure, a set of anchor nodes is first calculated for both the users and the items based on K-means clustering algorithm. Then each node has a similarity distribution over all the anchors denoting the relative position of the node within the graph,

$$GS_{u,\mathcal{A}_u^{t,k}}^t = \frac{e^{\text{SIM}(\mathbf{e}_u^t, \mathcal{A}_u^{t,k})/\tau}}{\sum_{k'=1}^{K} e^{\text{SIM}(\mathbf{e}_u^t, \mathcal{A}_u^{t,k'})/\tau}}, \text{SIM}(\mathbf{a}, \mathbf{b}) = \mathbf{a}^{\mathrm{T}}\mathbf{b}, \tag{25}$$

where $\mathcal{A}_u^t$ denotes the anchors calculated during task $t$. The goal of global structure preserving is to maintain the similarity distribution across different tasks,

$$\mathcal{S}_{u,\mathcal{A}_u} = D_{KL}(GS_{u,\mathcal{A}_u^t}^t || GS_{u,\mathcal{A}_u^{t-1}}^{t-1}). \tag{26}$$

Finally, besides the topological information including both the local and global structures, the information contained in

each user/item node is also highly valuable. To protect this part of information, a self-distillation is proposed,

$$\mathcal{L}_{self} = \big( \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{\eta_u}{||\eta_U||_2} ||\mathbf{e}_u^{t-1} - \mathbf{e}_u^t||_2$$

$$+ \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} \frac{\eta_i}{||\eta_I||_2} ||\mathbf{e}_i^{t-1} - \mathbf{e}_u^t||_2 \big), \tag{27}$$

where $\eta_u = \frac{|\mathcal{N}_u^{t-1}|}{|\mathcal{N}_u^t|}$ and $\eta_i = \frac{|\mathcal{N}_i^{t-1}|}{|\mathcal{N}_i^t|}$ are the normalizing factors.

Multi-modal Structure-evolving Continual Graph Learning (MSCGL) [81] is proposed to tackle the challenges of CGL with multi-modal data. Specifically, MSCGL consists of two parts, including a Neural Architecture Search (NAS) module serving to adaptively optimize the model architecture for accommodating new tasks, and a Group Sparse Regularization (GRS) module for preserving crucial information of the previously learnt tasks. The search space of the model architecture is designed as multiple choices of different GNNs, in which the different modalities are processed with separate GNNs. Denoting the distribution of the model architecture a as $P(\mathrm{a}; \boldsymbol{\theta})$ parameterized by $\boldsymbol{\theta}$, the objective can be formulated as a bi-level optimization including the maximization of the expected accuracy $\mathcal{E}[\mathcal{R}(\mathrm{a}(w*, G))]$ and minimization of the training loss $\mathcal{L}_{train}(\mathrm{a}(w, G))$,

$$\max \mathcal{E}[\mathcal{R}(a(w*, G))], \tag{28}$$
$$s.t. \quad w* = \arg\min_w \mathcal{L}_{train}(\mathrm{a}(w, G)).$$

With the model obtained from NAS, a Group Sparse Regularization (GSR) is proposed to sparsify the network. During the sparsification, two restrictions proposed by [95] are adopted to ensure that the parameter changes are block-sparse and orthogonal to the previous parameters, so that the obtained new structure has less severe forgetting problem on the previous tasks.

In the area of knowledge graph research, Lifelong Knowledge Graph Embedding (LKGE) [82] is constructed as a specialized technique for knowledge graph based on embedding transfer and regularization. Similar to the general CGL, the sequential learning process can still be formulated as training the model on a growing graph sequence $\mathcal{S} = \{\mathcal{G}_1, \mathcal{G}_2, ..., \mathcal{G}_T\}$. Within the context of the knowledge graph, the graph is composed of facts, each of which is a tuple $(s, r, o)$ denoting the existence of a relation $r$ between the subject entity $s$ and object entity $o$. Besides, the new entities and relations emerging at the task $\tau$ are denoted as $\mathcal{E}_\tau$ and $\mathcal{R}_\tau$. LKGE mainly consists of three parts, including a masked knowledge graph autoencoder as the backbone model for knowledge graph embedding, an embedding transfer technique for transferring the learnt knowledge into unseen new entities and relations, and an embedding regularization module for alleviating the catastrophic forgetting problem. First of all, the masked knowledge graph autoencoder aims to capture information from both new and old tasks. Therefore, it is designed to minimize the distance between the embeddings and representations reconstructed by a masked neighborhood subgraph containing both old and potential new entities.

$$\mathcal{L}_{MAE} = \sum_{e \in \mathcal{E}_\tau} ||\mathbf{e}_\tau - \bar{\mathbf{e}}_\tau|| + \sum_{r \in \mathcal{R}_\tau} ||\mathbf{r}_\tau - \bar{\mathbf{r}}_\tau||. \tag{29}$$

In Equation 29, $\mathbf{e}_\tau$ and $\mathbf{r}_\tau$ are the embeddings of an entity $e$ and a relation $r$ in task $\tau$. The masked autoencoder $\mathrm{MAE}(\cdot)$ takes all facts involving an entity $e$ / relation $r$ from the observed tasks, and then generate a representation $\bar{\mathbf{e}}_\tau = \mathrm{MAE}(\cup_{j=1}^\tau \mathcal{N}_j(e))$ / $\bar{\mathbf{r}}_\tau = \mathrm{MAE}(\cup_{j=1}^\tau \mathcal{N}_j(r))$ to encourage the embeddings to contain information from both new and old facts. Then, inspired by the idea of TransE [96] ($\mathbf{s} + \mathbf{r} \approx \mathbf{o}$), the encoder is designed to represent the subject entity as the subtraction between object and relation ($\mathrm{f}_s(\mathbf{r}, \mathbf{o}) = \mathbf{o} - \mathbf{r}$), the relation entity as the subtraction between object and subject ($\mathrm{f}_s(\mathbf{s}, \mathbf{o}) = \mathbf{o} - \mathbf{s}$). Accordingly, the encoders can be derived as:

$$\bar{\mathbf{e}}_\tau = \frac{\sum_{j=1}^\tau \sum_{(s,r,o) \in \mathcal{N}_j(e)} \mathrm{f}_s(\mathbf{r}_\tau, \mathbf{o}_\tau)}{\sum_{j=1}^\tau |\mathcal{N}_j(e)|} \tag{30}$$

$$\bar{\mathbf{r}}_\tau = \frac{\sum_{j=1}^\tau \sum_{(s,r,o) \in \mathcal{N}_j(r)} \mathrm{f}_r(\mathbf{s}_\tau, \mathbf{o}_\tau)}{\sum_{j=1}^\tau |\mathcal{N}_j(r)|}. \tag{31}$$

However, in continual learning, the access to all previous data is now allowed, and the aim is approximated as,

$$\bar{\mathbf{e}}_\tau \approx \frac{\sum_{j=1}^{\tau-1} |\mathcal{N}_j(e)| \mathbf{e}_{\tau-1} + \sum_{(s,r,o) \in \mathcal{N}_j(e)} \mathrm{f}_s(\mathbf{r}_\tau, \mathbf{o}_\tau)}{\sum_{j=1}^\tau |\mathcal{N}_j(e)|} \tag{32}$$

$$\bar{\mathbf{r}}_\tau \approx \frac{\sum_{j=1}^{\tau-1} |\mathcal{N}_j(r)| \mathbf{r}_{\tau-1} + \sum_{(s,r,o) \in \mathcal{N}_j(r)} \mathrm{f}_s(\mathbf{s}_\tau, \mathbf{o}_\tau)}{\sum_{j=1}^\tau |\mathcal{N}_j(r)|}. \tag{33}$$

The learning on the new task $\tau$ is also based on TransE, and is formulated as,

$$\mathcal{L}_{new} = \sum_{(s,r,o) \in D_\tau} \max(0, \gamma + \mathrm{f}(\mathbf{s}, \mathbf{r}, \mathbf{o}) - \mathrm{f}(\mathbf{s}', \mathbf{r}, \mathbf{o}')), \tag{34}$$

where $D_\tau$ contains all facts in task $\tau$, $\gamma$ is the margin hyperparameter, and $(\mathbf{s}', \mathbf{r}', \mathbf{o}')$ denotes the embedding of a negative fact, which is generated by replacing the entities in a real fact with random entities. Second, the embedding transfer technique transfers the learned knowledge to unseen new entities by initializing the unseen entities as a aggefation of the learned entities,

$$\mathbf{e}_\tau = \frac{1}{|\mathcal{N}_\tau(e)|} \sum_{(e,r,o) \in \mathcal{N}_\tau(e)} \mathrm{f}_s(\mathbf{r}_{\tau-1}, \mathbf{o}_{\tau-1}). \tag{35}$$

Third, a regularization term is designed for both entity and relation embeddings for alleviating the forgetting issue,

$$\mathcal{L}_{old} = \sum_{e \in \mathcal{E}_{\tau-1}} \omega(e) ||\mathbf{e}_\tau - \mathbf{e}_{\tau-1}||_2^2 + \sum_{r \in \mathcal{R}_{\tau-1}} \omega(r) ||\mathbf{r}_\tau - \mathbf{r}_{\tau-1}||_2^2, \tag{36}$$

in which the contribution of the regularization terms are controlled by $\omega(\cdot)$. $\omega(\cdot)$ is calculated based on the number of observed data, i.e. $\omega(e) = 1 - \frac{|\mathcal{N}_\tau(e)|}{\sum_{j=1}^\tau |\mathcal{N}_j(e)|}$, $\omega(r) = 1 - \frac{|\mathcal{N}_\tau(r)|}{\sum_{j=1}^\tau |\mathcal{N}_j(r)|}$. Finally, the total loss function consists of all of the terms introduced so far,

$$\mathcal{L} = \mathcal{L}_{new} + \alpha \mathcal{L}_{old} + \beta \mathcal{L}_{MAE} \tag{37}$$

Besides the works on developing novel CGL techniques, researchers are also interested in exploring the effectiveness of existing continual learning techniques for specific CGL problems. For example, in fake news detection tasks, to continually accommodating new data without retraining the model on all previous data , which is prohibitively expensive,
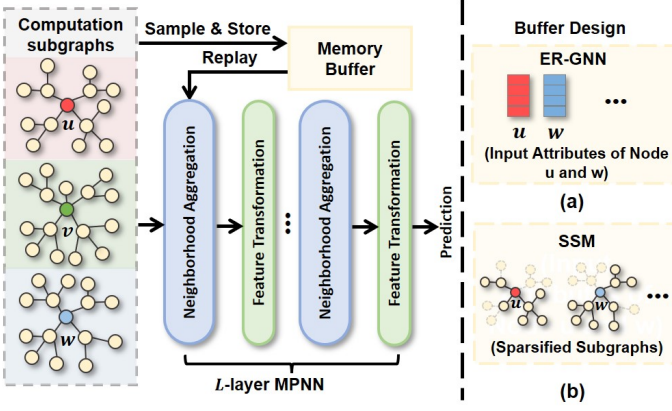
Fig. 4: Pipelines of two representative memory replay based techniques. (a) ER-GNN [83] stores single nodes. (b) SSM [77] stores sparsified computation subgraphs.

[58] adopts the Elastic Weight Consolidation (EWC) [91] and Gradient Episodic Memory (GEM) [97], which successfully protect the performance on both the previous data and new data empirically. Both EWC and GEM are model-agnostic continual learning methods, and do not explicitly consider the topological information within the graph data.

### 4.2 Memory Replay based Methods

Memory replay based methods prevent forgetting by retraining the model with representative data from previous tasks. Traditional continual learning methods process individual data without interactions, which can be simply sampled and stored in a buffer. However, for learning on graphs, memory replay based methods would meet the memory explosion challenge [1]. When generating the representation of a node (a datum), Graph Neural Networks (GNNs) typically aggregate information from multi-hop neighbors. Therefore, to regenerate the representation of a single node, the information from an exponentially expanding neighborhood have to be stored. On dense graphs, the memory consumption would easily become intractable [1]. Due to this challenge, Experience Replay Graph Neural Network (ER-GNN) [83] (Figure 4 (a)) chooses to ignore the graph topology and only store the attributes of one single node for regenerating the representation. Denoting the node buffer as $\mathcal{B}$, the loss function of ER-GNN when learning the $n$-th task can be formulated as

$$\mathcal{L}_n(\boldsymbol{\theta}, \mathcal{D}_i^{tr}, \mathcal{B}) = \beta \mathcal{L}_n(\boldsymbol{\theta}, \mathcal{D}_i^{tr}) + (1 - \beta)\mathcal{L}_n(\boldsymbol{\theta}, \mathcal{B}), \quad (38)$$

where

$$\beta = \frac{|\mathcal{B}|}{|\mathcal{D}_i^{tr}| + |\mathcal{B}|}. \quad (39)$$

In other words, with the memory buffer, the total training loss is a weighted summation of the loss calculated with the data from the current task and the buffered data. And the contribution of these two parts are balanced based on their relative sizes. To properly populate the memory buffer, three different strategies are proposed in ER-GNN to select the representative nodes, including Mean of Feature (MF), Coverage Maximization (CM), and Influence Maximization

(IM). The MF strategy is developed based on the intuition that the average feature vector of each class is representative, therefore the selected nodes should be the ones that are the closest to the average feature vector. The average feature vector (prototype) for a class $l$ can be obtained in two ways,

$$\mathbf{c}_l = \frac{1}{|V_l|}\sum_{(v \in V)}\mathbf{x}_v, \mathbf{c}_l = \frac{1}{|V_l|}\sum_{(v \in V)}\mathbf{h}_v, \quad (40)$$

where $V_l$ denotes the node set of the class $l$, $\mathbf{x}_v$ is the input attributes of the node $v$, and $\mathbf{h}_v$ is the representation of $v$ generated by the GNN. The second strategy, CM, aims to maximize the diversity of the selected nodes. Specifically, for each node in each task, the number of nodes from different classes within a certain distance $d$ is first counted as,

$$\mathcal{N}(v) = \{u | \mathrm{dist}(v, u) < d, \mathbf{y}_v \neq \mathbf{y}_u\}. \quad (41)$$

In equation 41, $y_v$ denotes the label of node $v$, $\mathrm{dist}(\cdot, \cdot)$ measures the distance between two nodes, and $d$ is a specified threshold. The lower $|\mathcal{N}(v)|$ a node $v$ has, the more distant node $v$ is from the nodes from different classes. Therefore, to maximize the coverage of the selected nodes in the representation space, CM is designed to select the several nodes with the lowest $|\mathcal{N}(v)|$. The third strategy, IF, aims to find the nodes with the most influence on the model parameters. Accordingly, a naive approach to probe the influence of a node is removing it and retraining the model to check the parameter change, which is very inefficient. Therefore, based on an alternative approach [98] to replace the node removal with loss upweighting, the authors propose to approximate the importance of a training node by its influence on the loss of testing nodes. For a training node $v$ and a testing node $u$, the influence of node $v$ during task $n$ is derived as,

$$I(v, u) = -\nabla_{\boldsymbol{\theta}}\mathcal{L}_n(u)^T \mathbf{H}_{\boldsymbol{\theta}}^{-1}\nabla_{\boldsymbol{\theta}}\mathcal{L}_n(v), \quad (42)$$

where $\mathbf{H}_{\boldsymbol{\theta}}$ is the Hessian matrix of the loss function. Finally, the influence on all testing nodes are summed up as the approximation of the node $v$.

Although being easy to implement, the valuable topological information is not preserved in ER-GNN. To preserve the graph topology and maintain a tractable memory consumption at the same time, Sparsified Subgraph Memory (SSM) [77] is proposed to store sparsified computation subgraphs. Given a computation subgraph to store, SSM would sample a fixed number of nodes iteratively from the 1-hop to a specified $K$-hop neighbors. This hop-by-hop sampling manner can ensure the connectivity of the sparsified subgraph. After storing the sparsified subgraphs in the memory buffer $\mathcal{SSM}$, the total loss for learning on a task $\tau$ is derived as,

$$\mathcal{L} = \underbrace{\sum_{u \in \mathbb{V}_\tau} l(\mathrm{f}(\mathcal{G}_u^{sub}; \boldsymbol{\theta}), \mathbf{y}_u)}_{\text{loss of the current task } \mathcal{L}_\tau}$$
$$+ \lambda \underbrace{\sum_{\bar{\mathcal{G}}_v^{sub} \in \mathcal{SSM}} l(\mathrm{f}(\bar{\mathcal{G}}_v^{sub}; \boldsymbol{\theta}), \mathbf{y}_v)}_{\text{auxiliary loss } \mathcal{L}_{aux}}. \quad (43)$$

Based on SSM, a Ricci-curvature based Subgraph Episodic Memory (SEM) [78] is further constructed to sparsify the

computation subgraphs based on the information bottleneck of each edge. The information bottleneck of is measured with its Ricci-curvature, which is a theoretically justified metric for evaluating the smoothness of information flow via each edge [93]. The Ricci-curvature of an edge $e_{uv}$ is defined as:

$$Ric(u,v) := \frac{1}{d_u} + \frac{1}{d_v} - 2 + 2\frac{|\#_{\triangle}(u,v)|}{\max\{d_u,d_v\}} + \frac{|\#_{\triangle}(u,v)|}{\min\{d_u,d_v\}}$$
$$+ \frac{\gamma_{\max}^{-1}}{\max\{d_u,d_v\}}(|\#_{\square}^u(u,v)| + |\#_{\square}^v(u,v)|),$$
(44)

where $d_u$ is the degree of node $u$, $\#_{\triangle}(u,v) := \mathcal{N}^1(u) \cap \mathcal{N}^1(v)$ denotes the number of triangles formed with $u$, $v$ and one of their common 1-hop neighbors. $\#_{\square}^v(u,v) := \{w \in \mathcal{N}^1(u) | w \neq v, \exists q \in (\mathcal{N}^1(w) \cap \mathcal{N}^1(v))\backslash\mathcal{N}^1(u)\}$ denotes the number of 4-cycles. Both of them count the higher order paths for information flow. $\gamma_{\max} = \max\{\max_{p\in\#_{\square}^u(u,v)}\{|\mathcal{N}^1(p) \cap \#_{\square}^v(u,v)\backslash\mathcal{N}^1(u)| - 1\}, \max_{q\in\#_{\square}^v(u,v)}\{|\mathcal{N}^1(q) \cap \#_{\square}^u(u,v)\backslash\mathcal{N}^1(v)| - 1\}\}$ is the maximal number of 4-cycles at edge $e_{u,v}$ via a common node, and serves as a normalization term. However, since Ricci-curvature has a high computational complexity, a random walk based surrogate is introduced to indicate the relative magnitude of the Ricci-curvature of each edge. Finally, the edges are selected by the sampling from the random walk distribution.

SSM is a model agnostic method that can be implemented with any GNN. Another recently proposed memory based CGL approach, PDGNNs-TEM [79], formulates a framework to store the complete topological information of each computation subgraph with a single embedding vector. Specifically, PDGNNs-TEM consists of a general GNN framework named as Parameter Decoupled Graph Neural Networks (PDGNNs) and a memory buffer named as Topology-aware Embedding Memory (TEM). As analyzed by [77], the key challenge in storing complete topological information of graph data is the memory explosion problem, *i.e.* the size of the computation subgraph grows exponentially with the number of hops. Therefore, PDGNNs is proposed to decouple the trainable parameters from the input computation subgraph, so that the nodes in the computation subgraph do not interact with the trainable parameters individually, and retraining the model only needs an overall information of each computation subgraph. Specifically, instead of iteratively aggregating the neighboring nodes and conducting node feature transformation, PDGNNs first encode the entire computation subgraph into a topology-aware embedding with a non-parametric function to capture the topological information of the computation subgraph,

$$\mathbf{e}_v = \mathrm{f}_{topo}(\mathcal{G}_v^{sub}).$$
(45)

In Equation 45, $\mathbf{e}_v$ denotes the topology-aware embedding (TE) for node $v$, $\mathcal{G}_v^{sub}$ is the computation subgraph of node $v$, and $\mathrm{f}_{topo}(\cdot)$ is the non-parametric function that can be instantiated with different forms. In the paper, both linear formulation and non-linear formulation of $\mathrm{f}_{topo}(\cdot)$ are explored. After that, the obtained TEs are then fed into a trainable function to generate the final output prediction,

$$\hat{\mathbf{y}}_v = \mathrm{f}_{out}(\mathbf{e}_v; \boldsymbol{\theta}).$$
(46)

With the PDGNNs framework, for retraining the model with a certain node $v$, the entire computation subgraph $\mathcal{G}_v^{sub}$ is no longer needed. Instead, only the TE of $v$ needs to be stored for re-training the trainable function $\mathrm{f}_{out}(\cdot)$. Accordingly, the Topology-aware Embedding Memory ($\mathcal{TEM}$) is developed to store the selected TEs for each task,

$$\mathcal{TEM} = \mathcal{TEM} \bigcup \mathrm{sampler}(\{\mathbf{e}_v \mid v \in \mathbb{V}_\tau\}, n),$$
(47)

where $\mathrm{sampler}(\cdot, \cdot)$ is the chosen sampling strategy for populating the memory buffer, $n$ is the budget for storage, and $\mathbb{V}_\tau\}$ denotes the node set of the current task. Compared to ER-GNN and SSM, TEM can maintain complete topological information for each node with only one emebdding vector, therefore is highly efficient in terms of both space complexity and computation complexity during memory replay. But TEM has to be accompanied with a GNN following the PDGNNs framework, while ER-GNN and SSM can be implemented with any GNN.

As mentioned before in Section 4.1, the inter-task edge connections would alter the neighborhood of the nodes in the previous tasks and cause concept shift, while the methods mentioned above do not explicitly consider this problem. Targeting this challenge, Structure-Evolution-Aware Experience Replay (SEA-ER) [88] is proposed to explicitly consider the evolution of the graph structure when populating the memory buffer. Denoting the memory buffer for the $j$-th task as $P_j$, the objective for populating the buffer is to maximize the structural similarity between the selected samples and the rest nodes,

$$\min_{P_j \subset \mathcal{V}_j} \max_{u\in\mathcal{V}_j\backslash P_j} \min_{v\in P_j} \mathrm{d}_{\mathrm{spd}(u,v)}, \quad s.t.|P_j| = b,$$
(48)

where $\mathrm{d}_{\mathrm{spd}(u,v)}$ denotes the shortest path distance between node $u$ and $v$, $b$ is the memory budget for task $j$, and $\mathcal{V}_j$ is the node set of the task $j$. Since the evolution of the graph structure may also influence the distribution of the data in previous tasks, besides the sampling strategy, an importance reweighting technique is also designed in SEA-ER to rescale the contribution of different nodes. Specifically, the objective function with the node reweighting for learning on task $i$ is formulated as,

$$\mathcal{L} = \frac{1}{|\mathcal{V}_i^{train}|} \sum_{v\in\mathcal{V}_i^{train}} \mathcal{L} + \sum_{P_j\in\mathcal{P}_i} \frac{1}{|P_j|} \sum_{v\in P_j} \beta_v\mathcal{L}(v),$$
(49)

where $\mathcal{P}_i = \{P_1, ..., P_{i-1}\}$ denotes the memory buffer during learning task $i$, which contains representative data selected from all previous tasks. $\beta_v$ is the reweighting coefficient on each node $v$, and is obtained by kernel mean matching (KMM) [99], which can be mathematically formulated as the solution of the following problem,

$$\min_\beta || \sum_{P_j\in\mathcal{P}_i} \sum_{v\in P_j} \beta_v\phi(h_v) - \phi(h_v')||^2, s.t.B_l \leq \beta_v \leq B_u,$$
(50)

where $h_v$ and $h_v'$ are node representations generated by the GNN model based on graph from task $i$ and $i-1$, and $\phi(\cdot)$ is the function that maps the node representations into the reproducing Hilbert space associated by the adopted kernel. Besides, $B_l$ and $B_u$ are the lower bound and upper bound on $\beta$ to make sure that the weight is suitable for the majority

of the nodes. With the formulation above, the reweighting procedure is actually minimizing the distance between the node representations on new and old tasks by finding the most appropriate weight values.

Similar to SEA-ER, Incremental Graph Neural Networks (IncreGNN) [84] also considers the influence of inter-task edges on nodes in the previous nodes. To tackle this challenge, IncreGNN not only proposes an approach to evaluate the node importance based on both their influence to the new nodes and their influence in the original graph, but also adopts a regularization technique to further boost the performance. Specifically, the newly incoming graph data is denoted as the base change group. Based on the base change group, the 1-hop neighbors of the base change group in the previous task graph is defined as the first order change group, the 1-hop neighbors of the first order change group is defined as the second order change group, and so forth. Since the influence of the new nodes on the previous nodes decreases with the increase of the order of the changing group, the ratio of the budget for node sampling at different order of change group is designed as,

$$\frac{1}{i}/(1 + \frac{1}{2}, ..., +\frac{1}{K}),\qquad(51)$$

where $i$ is the number of order that is being sampled, $K$ is the total number of change groups. After defining the importance of the groups of nodes with different distance to the new nodes, the importance of an individual node $v$ is evaluated by the personalized PageRank algorithm run on the previous task graph $\mathcal{G}^{t-1}$, and is denoted as $\pi_v$. Then, denoting the number of nodes to sample at $k$-order of change group as $n_k$, the nodes at each order with higher personalized PageRank values are selected as the ones with more correlation with the new nodes,

$$I(\mathcal{G}^t) = \cup_{k=1}^K \{v_i | \pi_{v_i} > \pi_{v_j}, i \in [1, n_k], j \notin [1, n_k]\}. \quad (52)$$

The above process only select the old nodes affected by the new nodes. To also maintain the knowledge learnt from the old nodes without inter-task edges to new nodes (unaffected nodes), the next step is to evaluate the importance of each unaffected node for sampling and storage. Specifically, K-means clustering is first conducted to divide the unaffected nodes into $K$ clusters. Then, the node importance within each cluster is determined by the node degree, and the nodes with the highest degrees are selected and stored in the memory buffer. Finally, IncreGNN also adopts the Memory Aware Synapse (MAS) [92] as the regularization on the model parameters to further enhance the capability to retain learnt knowledge.

Different from the approaches above, which retrain the same model structure with buffered data, [85] proposes to also dynamically adjust the model to accommodate new tasks. The proposed model consists of two parts including a Reinforcement Learning based Controller (RLC) that decides the addition and pruning of the model structure, and a Child Network (CN) which is a evolvable GNN backbone controlled by the RLC. The RLC is based on reinforcement learning, which aims to learn an optimal policy $\pi(a_t|s_t))$ to choose the most appropriate action $a_t$ to modify the model based on the model state at the $t$-the task, so that the expected summation of the future rewards $R_t = \sum_{k=0}^{\infty} \gamma_{t+k+1}^k$ is

maximized. The decaying factor $\gamma \in (0, 1]$ serves to gradually shrink the contribution of the expected rewards at distant time steps. Accordingly, the Q-function of RLC for taking an action $a_t$ with a model state $s_t$ is,

$$Q_\pi(s_t, a_t) = \mathrm{E}_\pi[R_t|S = s_t, A = a_t]. \qquad (53)$$

The RLC in [85] is designed as an LSTM, which outputs a set of $m$ actions $a_{1:m}$ for adjusting the number of features in each of the $m$ layers of the CN. The search space for the action is defined as the addition and deletion operation of the CN hidden layers. When RLC is properly trained, the optimal actions for adjusting the CN is formulated as,

$$\hat{a}_{1:m} = \arg\max_{a \in \mathcal{S}} \mathbb{E}[R(a_{1:m}, s_t)]. \qquad (54)$$

Besides the RLC to optimize the structure of CN, to further enhance the continual learning capability of the model, the Dark Experience Replay [100] is also adopted after being extended to graph data. Specifically, besides storing the data and the labels of the selected examples, their logits output by the CN is also buffered. During replay, the model retraining loss would accordingly consists of two parts including the classification loss calculated with the buffered labels and data, as well as the knowledge distillation loss calculated based on the stored logits. The total loss including the loss for learning the new task is then formulated as,

$$\mathcal{L} = \mathcal{L}_t + \alpha||\mathrm{f}_t(\mathcal{X}^i) - l_t||_2 + \beta\mathcal{L}_{cls}(\mathrm{f}_t(\mathcal{X}^i, \mathcal{Y}^i)), \qquad (55)$$

where $\mathrm{f}_t(\cdot)$ denotes the current model, $l_i$ is the logits stored from the $i$-the task, and $\mathcal{X}^i$ and $\mathcal{Y}^i$ are the data and labels stored from task $i$.

Besides the models that store real data from learnt tasks, in Streaming Graph Neural Networks via Generative Replay (SGNN-ER), an auxiliary generative model is adopted to retrain the GNN with generated fake historical data. Overall, the generative model would generate the neighborhood (in the form of a random walk with restart sequence) based on the GAN framework. Denoting the generator and discriminator parameterized by $\phi^t$ and $\varphi^t$ at task $t$ as $\mathrm{G}_{\phi^t}$ and $\mathrm{D}_{\varphi^t}$, the learning objective can be formulated as a confrontation between the generator and the discriminator,

$$\min_{\phi^t} \max_{\varphi^t} \mathrm{V}(\mathrm{G}_{\phi^t}, \mathrm{D}_{\varphi^t}) = \mathrm{E}_{v \sim p_{data}(v)}[\log \mathrm{D}_{\phi^t}(v)] \qquad (56)$$
$$+ \mathrm{E}_{z \sim p_z(z)}[\log(1 - \mathrm{D}_{\varphi^t}(\mathrm{G}_{\phi^t}(z)))].$$

In SGNN-ER, the generated data is used for knowledge distillation between the old and new models to avoid forgetting, and the total loss can be formulated as

$$\mathcal{L} = r\mathrm{E}_{v \sim \mathcal{G}_A^t}[l(\mathrm{F}_{\boldsymbol{\theta}^t}(v)), y_v]$$
$$+ (1 - r)\mathrm{E}_{v' \sim G_\phi^{t-1}}[l(\mathrm{F}_{\boldsymbol{\theta}^t}(v'), \mathrm{F}_{\boldsymbol{\theta}^{t-1}}(v'))], \qquad (57)$$

where $\mathcal{G}_A^t$ is the affected part at task $t$, including the newly added subgraph and the previous nodes connected to the newly added subgraph.

Originally, the Generative Adversarial Network (GAN) framework is used for generating independent data. However, to generate the representation of a single node, GNNs would also take its multi-hop neighborhood with complex topological dependencies as input. Therefore, SGNN-ER proposes to generate random walk with restart (RWR) sequences, which are then converted into neighborhood subgraphs

based on the generated edge connectivity. Considering the neighborhood of some previous nodes is changed by the newly added nodes, to avoid enforcing the patterns that already disappear because of the neighborhood change, SGNN-ER proposes a forgetting mechanism to filter out the nodes that are significantly affected by the new nodes during memory replay. The mechanism consists of two steps. First, a set of significantly affected nodes $\mathcal{V}_C^t$ is first detected based on the change in their representations when the graph grows from $\mathcal{G}^{t-1}$ to $\mathcal{G}^t$,

$$\mathcal{V}_C^t = \{v| ||F_{\boldsymbol{\theta}^{t-1}}(v, \mathcal{G}^t) - F_{\boldsymbol{\theta}^{t-1}}(v, \mathcal{G}^{t-1})|| > \delta\}, \quad (58)$$

where $\delta$ is the threshold on the change of representations. Since $\mathcal{V}_C^t$ denotes the severely affected nodes, among the generated fake nodes for replay, the ones with high similarity to the nodes in $\mathcal{V}_C^t$ are rejected with high probability, which is formulated as,

$$p_{reject}(v) = \max(p_{sim}(v, u), u \in \mathcal{V}_C^t) \times p_r, \quad (59)$$

where $p_r$ is a hyperparameter controlling the total number of nodes to delete, and $p_{sim}(\cdot, \cdot)$ is the similarity function for measuring the similarity between two nodes,

$$p_{sim}(v, u) = \sigma(-||F_{\boldsymbol{\theta}^{t-1}}(v, \mathcal{G}^{t-1}) - F_{\boldsymbol{\theta}}(u, \mathcal{G}^{t-1})||). \quad (60)$$

Another related work in this branch is the Condense and Train (CaT) framework [73]. Different from the works that generate fake data for previous tasks, to balance the contribution of different tasks, CaT is designed to train the model only on the generated balanced graphs. Specifically, when a new task ($\mathcal{G}_\tau$) comes, CaT will first condense it based on the graph node embeddings generated by the same GNN encoder,

$$\tilde{\mathcal{G}}_\tau = \arg\min_{\tilde{\mathcal{G}}_\tau^*} \text{Dist}(\text{GNN}(\mathcal{G}_\tau; \theta), \text{GNN}(\tilde{\mathcal{G}}_\tau^*; \theta)). \quad (61)$$

Denoting the node embeddings $\text{GNN}(\mathcal{G}_\tau; \theta)$ and $\text{GNN}(\tilde{\mathcal{G}}_\tau^*; \theta)$ as $\mathbf{E}_\tau$ and $\mathbf{E}_\tau^*$, the distant function $\text{Dist}(\cdot, \cdot)$ is formulated as a weighted summation over the mean discrepancy between the node embeddings of each class in the two graphs,

$$\text{Dist}(\mathbf{E}_\tau, \mathbf{E}_\tau^*) = \sum_{c \in \mathcal{C}_\tau} \frac{|\mathbf{E}_{\tau,c}|}{|\mathbf{E}_\tau|} \cdot ||\text{Mean}(\mathbf{E}_{\tau,c}) - \text{Mean}(\mathbf{E}_{\tau,c}^*)||. \quad (62)$$

In equation 62 above, $\frac{|\mathbf{E}_{\tau,c}|}{|\mathbf{E}_\tau|}$ denotes the ratio of data belonging to class $c$ in the entire set of data, and $\text{Mean}(\cdot)$ refer to the mean of the embedding vectors. Next, to faithfully reflect the distance between the graphs in different embedding spaces, $\text{GNN}(\cdot; \theta)$ in CaT is initialized with different random parameters $\theta$, and the distance is averaged over all random parameter initialization. The obtained condensed graph of the current task is updated into a memory,

$$\mathcal{M}_\tau = \mathcal{M}_{\tau-1} \cup \{\tilde{\mathcal{G}}_\tau\}, \quad (63)$$

where $\mathcal{M}_\tau$ denotes the memory buffer containing all condensed graphs from the observed $\tau$ tasks. Finally, the model is trained solely on $\mathcal{M}_\tau$ for a balanced training.

The methods mentioned above are general CGL methods that are applicable to any graph data. For certain graph related tasks, due to the special properties of the graph data, some specialized techniques may be developed. For example,
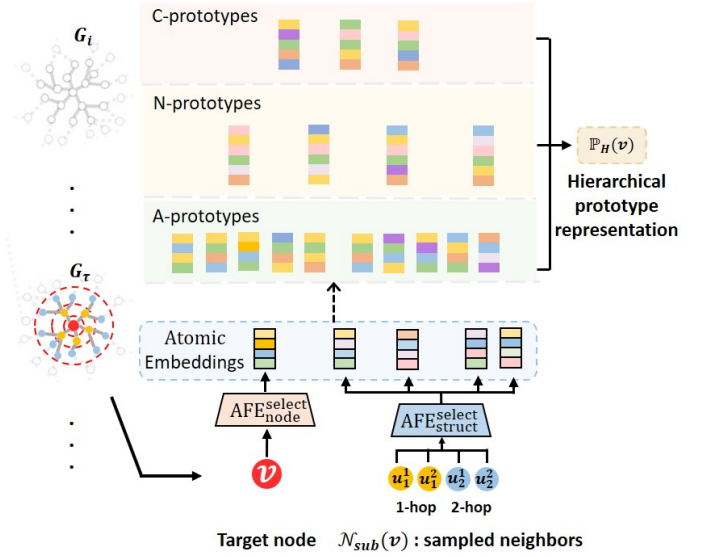


Fig. 5: Pipelines of HPNs [3] stores sparsified computation subgraphs.

specially designed for incremental learning in graph based recommender system, Structure Aware Experience Replay (SAER) [87] samples and stores representative user-item pairs based on reservoir sampling, so that the previously learnt user behaviour patterns can be maintained when learning new tasks. Specifically, the reservoir is a set of previously observed interactions denoted as $\mathcal{H}$. When a new set of interactions, *i.e.* $\mathcal{H}'$, is given, a subset $\hat{\mathcal{H}}$ will firstly be sampled from $\mathcal{H}$ according to an adopted sampling strategy. Then the model will be trained with $\mathcal{H}' \cup \hat{\mathcal{H}}$. In [87], two strategies are adopted. The first one is the uniform sampling,

$$\hat{\mathcal{H}} = \{(u, i) \sim U(\mathcal{H})\} \quad s.t \quad |\hat{\mathcal{H}}| = \gamma|\mathcal{H}|, \quad (64)$$

where $\gamma$ is a hyperparameter controlling the size of the sampled set. The uniform sampling serves as a naive baseline, and a more complex strategy is also adopted to balance the sampled interactions according to their degree, the probability for sampling an interaction $(u, i)$ is,

$$P(u, i) = \frac{1/d_u^T}{\sum_{(\hat{u},\hat{i}) \in \mathcal{H}} 1/d_{\hat{u}}^T}, \quad (65)$$

in which $d_u$ is the degree of the user $u$, and the temperature $T$ is a hyperparameter for regulating the smoothness of the probability distribution. The pairwise data structure is commonly adopted by recommender systems, but is not available for general graph data. Therefore this technique is not generally applicable.

Currently, memory replay based methods are among the most effective approaches and are easy to implement. But it requires extra space to store representative data.

## 4.3 Parameter isolation based methods

The final category, parameter-isolation based methods, protects the model's performance on previous tasks by entirely or partially isolating the model parameters for different tasks. Existing parameter isolation based CGL approaches are scarce, and one representative method is Hierarchical

Prototype Networks (HPNs) that propose to dynamically increment feature extractors and prototypes for accommodating new patterns. Specifically, HPNs consist of a set of Atomic Feature Extractors (AFEs) for extracting basic features from the given data, and three levels of hierarchical prototypes for storing learnt patterns. Given an input node, the AFEs will first extract basic features based on both the node attributes and its neighborhood relationship. Accordingly, the AFEs consist of two parts. One set of AFEs, denoted as $\text{AFE}_{\text{node}}$, serve to generating atomic node embeddings based on node attributes. The other set, denoted as $\text{AFE}_{\text{struct}}$, serve to capture the topological structure based on the neighborhood of the given node, and generate atomic structure embeddings accordingly. The obtained atomic embeddings correspond to the most basic features, and will then be matched to the existing atom-level prototypes (A-prototypes) according to their cosine similarity.

$$\text{Sim}_{E \to A}(v) = \{\frac{\mathbf{e}_i^T \mathbf{p}}{||\mathbf{e}_i||_2 ||\mathbf{p}||_2} | \mathbf{e}_i \in \mathbb{E}_A^{\text{select}}(v), \mathbf{p} \in \mathbb{P}_A\}, \quad (66)$$

where $\mathbb{E}_A^{\text{select}}(v)$ contains the atomic embeddings of the node $v$, and $\mathbb{P}_A$ is the set of A-prototypes. When the cosine distance between an A-prototype and an embedding is smaller than a threshold $t_A$, the A-prototype and the embedding is regarded as matched. If an embedding cannot be matched to any existing prototype, then it is regarded as new knowledge and will be used to initialize a new prototype. After A-prototype matching, the matched A-prototype will be further embedded and matched to higher level prototypes including the node-level prototypes (N-prototypes) and the class-level prototypes (C-prototypes). Each node will be matched to one N-prototype denoting the overall information of the entire node, and one C-prototype representing the properties of the class it belongs to. The union of the selected prototypes at different level are the final representation of the node $v$ (Figure 5). Besides a novel model, The authors of [3] also propose a theoretical framework to analyze the continual learning capability of a model and justify that the proposed HPNs can theoretically eliminate the forgetting problem when the hyperparameters are properly set.

Different from the majority of CGL works that solely focus on the growth of the graph, Parameter Isolation GNN (PI-GNN) [74] studies the CGL problem in the dynamic graph setting, in which both node addition and deletion are considered. Specifically, when the training of a $k$-layer GNN proceeds from task $\tau - 1$ to $\tau$, the graph evolves from $\mathcal{G}_{\tau-1}$ to $\mathcal{G}_\tau$. In this process, $\mathcal{G}_{\tau-1}$ can be divided into two parts: $\mathcal{G}_{\tau-1}^{unstable}$ and $\mathcal{G}_{\tau-1}^{stable}$. $\mathcal{G}_{\tau-1}^{unstable}$ consists of the nodes whose $k$-hop neighborhood is altered by the graph evolution (e.g., being connected to new nodes), while $\mathcal{G}_{\tau-1}^{stable}$ refers to the remaining part. Then, the changed part of the graph from $\mathcal{G}_{\tau-1}$ to $\mathcal{G}_\tau$ is denoted as $\Delta\mathcal{G}_\tau$. When learning on a new task $\tau$, the main idea of PI-GNN is to first preserve the stable knowledge in $\mathcal{G}_{\tau-1}^{stable}$, and then allocate new separated parameters to accommodate the new knowledge. The stable knowledge preservation is formulated as only optimizing the parameters on $\mathcal{G}_{\tau-1}^{stable}$,

$$\min_{\theta_{\tau-1}} \mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}^{stable}; \theta_{\tau-1})\big) \quad (67)$$

$$i.e. \min_{\theta_{\tau-1}} \mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}; \theta_{\tau-1})\big) - \mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}^{unstable}; \theta_{\tau-1})\big). \quad (68)$$

Since computing $\mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}; \theta_{\tau-1})\big)$ is time consuming, and the access to previous data may not always be available, $\mathcal{G}_{\tau-1}$ is approximated with a set of representative data $\mathcal{G}_{\tau-1}^{memory}$ sampled from $\mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}; \theta_{\tau-1})\big)$, i.e.,

$$\theta_{\tau-1}^{stable} = \arg\min_{\theta_{\tau-1}} \mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}^{memory}; \theta_{\tau-1})\big) - \quad (69)$$
$$\beta\mathcal{L}\big(\text{f}(\mathcal{G}_{\tau-1}^{unstable}; \theta_{\tau-1})\big).$$

Since the memory data $\mathcal{G}_{\tau-1}^{memory}$ has a different size compared to the original graph $\mathcal{G}_{\tau-1}$, the hyperparameter $\beta$ is introduced for the size imbalance issue. $\theta_{\tau-1}^{stable}$ is fixed after encoding the previous knowledge, and new parameters $\theta_\tau^{new}$ are allocated for encoding the new knowledge in $\Delta\mathcal{G}_\tau$ and $\mathcal{G}_{\tau-1}^{memory}$. Finally, the parameters obtained after learning task $\tau$ is $\theta_\tau = \theta_{\tau-1}^{stable} \oplus \theta_\tau^{new}$.

## 4.4 Application of Traditional Continual Learning Techniques

Although not specialized for graph data, some of the traditional CL techniques for Euclidean data are both data- and model-agnostic, therefore are applicable to CGL. For example, regularization based methods like Elastic Weight Consolidation (EWC) [91] and Memory Aware Synapses (MAS) [92] that only add regularization terms to the model weights, memory replay based models like Gradient Episodic Memory (GEM) [97] that stores the gradients of the losses on previous tasks and clip the current gradients to avoid increasing the previous task losses, and parameter-isolation based methods like Progressive Neural Networks [101] and Supermasks in Superposition (SupSup) [102] that entirely or partially separate parameters for different tasks. The key characteristic of these works that makes them applicable to CGL is that they are agnostic of the model structure and the data structure. While the other methods not satisfying this criterion are not directly applicable to CGL. For example, most memory replay based methods would directly store representative data from previous tasks. However, on a growing graph, since GNNs generate the representation of one datum (node) not only based on itself but also on its multi-hop neighbors (computation subgraph), and different GNNs may adopt different computation subgraph constructions, it is unclear which part of the graph should be stored.

Targeting the inapplicability of some of the traditional CL methods, Feature Graph Network (FCN) [103] is proposed to accommodate the data by transforming each node into an independent feature graph, so that the data structure is no longer an obstacle for applying traditional continual learning techniques. When generating the feature graph, the topological information is also implicitly encoded. FCN provides a general approach to enable traditional CL techniques in any graph data. But for some special graphs, like knowledge graphs, since the data storage format is more tidy, traditional CL techniques are applicable with little modification. For example, considering that the knowledge graph in robotics applications needs to be frequently updated with new information, [90] extends several continual learning techniques including Progressive Neural Networks [101], Copy weight Re-init Deep Generative Replay [104] to knowledge graph embedding for avoiding learning all the concepts afresh. We
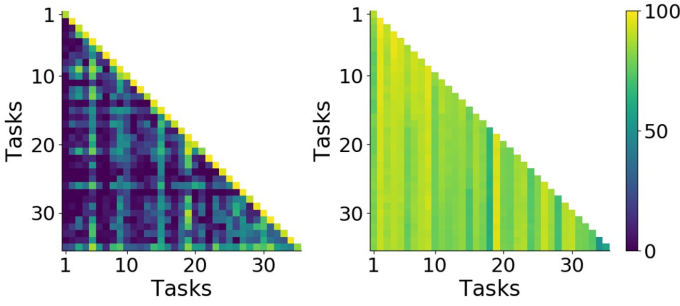
Fig. 6: An example of two visualized performance matrices provided in CGLB [1]

denote the ensemble of methods proposed by [90] as Adapted Continual Learning Methods (ACLM) for convenience.

## 5 EVALUATION OF CGL MODELS

Different from standard learning setting that is only concerned with one task, the evaluation of continual learning models have to consider the model's performance on all learnt tasks. Therefore, the most through approach to evaluate a continual learning model is to show its performance on each previous task after learning each new task. Formally, this result could be represented as a performance matrix $M^p \in \mathbb{R}^{T \times T}$ [1]. Each entry $M^p_{i,j}$ is the model's performance on the $j$-th task after the model has learnt the first $i$ tasks. The performance matrix contains the raw performance during the entire learning process, which can be visualized as a heatmap as shown in Figure 6. Based on the performance matrix, multiple different metrics are adopted by the researchers. For example, considering that the performance matrix is inconvenient for performance comparisons across different methods, the average performance (AP) and average forgetting (AF) can be used [1]. These two metrics are widely adopted in CGL as well as traditional continual learning works, while the names may be different in different works. For example, they are named as Average Accuracy (ACC) and Backward Transfer (BWT) in [97], [105], Average Performance (AP) and Average Forgetting (AF) in [1], [83] and Performance Mean (PM) and Forgetting Mean (FM) in [2]. Broadly speaking, the entries of the performance matrix may not only be accuracy and can also be other metrics like F1 scores. Therefore, we follow the names (AP and AF) in CGLB [1]. Given a task sequence of length $T$, AP is calculated as the average of the performance on all previous tasks after learning the entire task sequence $\frac{\sum_{j=1}^{T} M^p_{T,j}}{T}$. The AF is calculated as the average of the forgetting (performance decrease). The forgetting has two slightly different versions in different works. First, the forgetting of a task $j$ after learning the final task $T$ could be defined as $f_j = M^p_{T,j} - M^p_{j,j}$, which indicates how much the performance decreases from the time when the model has just learnt task $j$ to the time when the model has learnt all tasks [1], [97]. Another definition compares the final performance with the best performance ever obtained when the model is trained from task 1 to task $T$, i.e. $f_j = \max_{i=1,...,T-1} M^p_{i,j} - M^p_{T,j}$. Different from the first definition, this definition compare the final performance on

task $j$ with the best performance ever obtained on task $j$ before learning on the final task $T$. Another difference between these two definitions is that the sign of the computation. With the definition of the forgetting, the average forgetting (AF) can be naturally defined as $\frac{\sum_{j=1}^{T-1} f_j}{T-1}$. Further, to reflect the learning dynamics when learning sequentially on the tasks, the AP and AF after learning each task can be calculated, i.e. $\left\{ \frac{\sum_{j=1}^{i} M^p_{i,j}}{i} | i = 1, ..., T \right\}$ and $\left\{ \frac{\sum_{j=1}^{i-1} M^p_{i,j} - M^p_{j,j}}{i-1} | i = 2, ..., T \right\}$. Besides AP and AF, researchers are also interested in the forward transfer during learning on the task sequence, and a forward transfer is also defined [97] as, $\frac{\sum_{j=2}^{T} M^p_{j-1,j} - \bar{b_j}}{T-1}$, where $\bar{b_j}$ denotes the performance on task $j$ with random model initialization. Accordingly, $M^p_{j-1,j} - \bar{b_j}$ reflects how does the learning from task 1 to task $j-1$ benefit the performance on task $j$, i.e. forward transfer.

## 6 CGL BENCHEMARKS & DATASETS

Since CGL is a newly emerging area, different works adopt different datasets and conduct experiments under different settings, a fair comparison across different methods becomes difficult, which greatly impedes the development of the area. Accordingly, several benchmark works are proposed to provide a consistent setting and dataset protocol for developing and comparing CGL techniques. In 2021, a benchmark for CGL with graph-level tasks is introduced [8]. The benchmark is constructed based on three datasets: MNIST, CIFAR10, and OGBG-PPA. MNIST and CIFAR10 are image datasets that are originally used for computer vision. However, following the process in [106], the images are converted to graphs. Given an image, with the SLIC algorithm [107], the pixels within a small local area with homogeneous intensity are defined as a super-pixel, corresponding to a node in the converted graph. The topology (graph edges) is constructed as the k-nearest neighbor adjacency matrix. For the task construction, both MNIST and CIFAR10 contains 10 classes in total, and the classes are divided into five 2-class tasks. The OGBG-PPA [108] dataset contains a set of protein association neighborhoods (subgraphs) extracted from the large protein protein interaction (PPI) network. OGBG-PPA consists of 37 classes, each of which corresponds to a taxonomy group. The 37 groups are also split into 5 tasks. The first task contains 17 classes while the other four classes contain 5 classes. Constructing the tasks by splitting the classes of the data is a common practice for CGL, because the data in different classes generally follow different distributions, which easily cause catastrophic forgetting and is ideal for CGL tasks. Besides the benchmark task construction, experiments of 4 baselines including Naive, EWC, Replay, and LwF are conducted under the class-IL scenario. The accompanied code of this benchmark is provided through https://github.com/diningphil/continual_learning_for_graphs.

In 2022, the Continual Graph Learning Benchmark (CGLB) [1] is proposed, which is a more comprehensive benchmark covering 4 node-level datasets and 3 graph-level datasets under both task-IL and class-IL scenarios. The node-level datasets include CoraFull-CL, Arxiv-CL, Reddit-CL, and Products-CL. Both CoraFull-CL and Arxiv-CL are constructed based on citation networks provided in the public datasets CoraFull [109] and OGB-Arxiv [108].

Reddit-CL is constructed from social network consisting of Reddit posts [110]. Products-CL is constructed from the OGB-Products dataset [108], which is based on the Amazon product co-purchasing network. In CGLB, to maximize the number of tasks and increase the continual learning difficulty, in CGLB, all of these node-level tasks are split into 2-class tasks. For node-level CGL tasks, the growth of the graph would bring in inter-task edges connecting the nodes in previous and new tasks, and cause concept drift to the previous nodes connected by the inter-task edges. Accordingly, whether to preserve the inter-task edges when increasing the tasks is studied separately in CGLB. All of these node-level benchmark datasets support both task-IL and class-IL learning scenarios. On the node-level benchmark datasets, preliminary experimental results of multiple baseline methods are also provided under different settings. For graph-level CGL tasks, three benchmark datasets are constructed in CGLB including SIDER-tIL, Tox21-tIL, and Aromaticity-CL. SIDER-tIL is constructed from the SIDER dataset [111], which contains 1,427 drugs falling into 27 different classes. Since SIDER is a multi-label dataset, it is naturally constructed into a task-IL dataset with 27 tasks. Tox21-tIL is constructed from the Tox21 dataset [111], which is a multi-label dataset containing 8,014 molecules with 12 labels. Similar to SIDER, Tox21 is also constructed into a task-IL dataset with 12 tasks. Aromaticity-CL is constructed from the PubChem BioAssay Dataset [112], which contains 3,945 molecules falling into 40 classes. Since PubChem BioAssay Dataset is a multi-class dataset, the constructed Aromaticity-CL contains fifteen 2-class tasks, and can be used for both task-IL and class-IL. In Aromaticity-CL, some classes with very few data are removed. The accompanied code of CGLB is provided via https://github.com/QueuQ/CGLB.

Finally, we list all commonly used datasets in Table 2. Among these datasets, some are constructed as standard CGL tasks (CoraFull, OGB-Arxiv, Reddit and OGB-Products are included in CGLB [1]), while the others may be configured differently in different works.

# 7 OPPORTUNITIES, FUTURE DIRECTIONS AND DISCUSSIONS

Research on CGL has achieved significant advancements, manifesting both in enhanced performance and the exploration of more complex and practical applications. However, as a newly emerging and fast-growing area, there are still multiple challenges to tackle and promising directions to explore. Moreover, CGL is highly relevant to many research areas, and could potentially benefit a vast majority of the graph learning community. In this section, we first introduce several key challenges to be tackled, and then introduce the topics that transcend the boundary of current CGL research and intersect with other areas.

## 7.1 Trade-off between Effectiveness and Space Complexity

Among the three strategies, regularization based methods are the most memory efficient since they only require little extra storage space, while memory-replay based methods and parameter-isolation based methods require significant extra space for storing the representative data or expanded network structure. However, as shown in recent works including benchmark work [1] and technical works [3], [77], [79], [117], the regularization based method is not as effective as the other two approaches. The reason is that although the regularization term can effectively restrict the change of the model parameters to prevent the forgetting of previous tasks, it also reduces the model's capacity for adapting to the new tasks. On the contrary, neither memory-replay based methods [77], [118], [119] nor the parameter-isolation based methods put restrictions on the model parameters. Moreover, the parameter-isolation based methods can even expand the model and increase the capacity so that the model can better adapt to the new tasks [101], [102], [120]. For the memory-replay based methods, an increase in the buffer size directly correlates with enhanced capacity for preserving information from preceding tasks, resulting in improved performance outcomes. For the parameter-isolation based methods, more budget for model expansion is also preferred for better performance. In a word, currently, the CGL techniques have to maintain a trade-off between the space complexity and the effectiveness in preventing forgetting. When the number of tasks becomes relatively large, how to achieve better performance with a limited memory space budget is still challenging.

## 7.2 Task-free CGL

Existing CGL techniques mostly assume that the task boundaries are available during training. For regularization based methods, task boundaries are used for matching the stored parameter importance to the corresponding tasks. For memory replay based methods, task boundaries help allocate memory budget for different tasks. For parameter-isolation based methods, the boundaries are also crucial to determine when to expand the network. However, in real-world applications, data from different tasks may come in as a mixture without a clear task boundary, raising multiple challenges. First, without the given task identities, it is critical to investigate how to properly detect the distribution shift or the emergence of new patterns so that the model or memory buffer can expand. Second, since the data from a certain task may not only emerge once, even if the distribution shift is detected, the incoming 'new' patterns may have already been learnt previously. In this case, if the model always treats the incoming data with distribution shift as a new task, memory space will be wasted to maintain duplicated knowledge. Therefore, how to properly manage the learnt knowledge so that the learnt knowledge is preserved effectively and efficiently is also a challenge to tackle.

## 7.3 Concept Drift of Nodes in Existing Tasks

As pointed out in CGLB [1] and multiple technical works [84], [121], since the inter-task edges connecting nodes from old and new tasks appear after learning on the old tasks, the affected nodes in the old tasks would experience concept drift. Currently, the only solution in the existing works [84], [88] is to retrain the model on these affected nodes within the updated large graph, so that the nodes after concept drift are learnt by the model. This approach assumes that all previous nodes including the affected nodes and their

TABLE 2: Commonly used datasets for CGL.

| Dataset | Related works | Link |
|---|---|---|
| CoraFull [109] | TWP [2],SSM [77],PDGNNs-TEM [79],RLC-CN [85],CGLB [1],SEM [78],CaT [73],UGCL [72] | github.com/shchur/gnn-benchmark#datasets |
| OGB-Arxiv [108] | RieGrace [75],SSM [77],PDGNNs-TEM [79],HPNs [3],CGLB [1],SEM [78],CaT [73],UGCL [72] | ogb.stanford.edu/docs/nodeprop/#ogbn-arxiv |
| Reddit [110] | TWP [2],RieGrace [75],DyGRAIN [89],ER-GNN [83],SSM [77],PDGNNs-TEM [79],SEA-ER [88],CGLB [1],SEM [78],CaT [73],UGCL [72] | snap.stanford.edu/graphsage/ |
| OGB-Products [108] | DyGRAIN [89],SSM [77],PDGNNs-TEM [79],HPNs [3],CGLB [1],SEM [78],CaT [73] | ogb.stanford.edu/docs/nodeprop/#ogbn-products |
| Cora [113] | RieGrace [75],ER-GNN [83],RLC-CN [85],SGNN-ER [86],SEA-ER [88],HPNs [3],PI-GNN [74] | github.com/tkipf/gcn/tree/master/gcn/data |
| Citeseer [113] | RieGrace [75],ER-GNN [83],RLC-CN [85],SGNN-ER [86],HPNs [3],PI-GNN [74] | github.com/tkipf/gcn/tree/master/gcn/data |
| Amazon Computers [114] | TWP [2],RLC-CN [85],PI-GNN [74] | nijianmo.github.io/amazon/index.html |
| Actor [115] | RieGrace [75],HPNs [3] | github.com/graphdml-uiuc-jlu/geom-gcn/tree/master/new_data/film |
| Gowalla [116] | GraphSAIL [80],SAER [87] | snap.stanford.edu/data/loc-gowalla.html |

neighboring nodes are accessible during the learning of new tasks. However, due to reasons like privacy issues, nodes from the previous tasks may not always be accessible, which is why continual learning works typically forbid to access historical data. In light of this, how to tackle the concept drift on the affected previous nodes without accessing previous nodes is a crucial problem, especially in dense networks with a large amount of newly emerging inter-task edges.

### 7.4 Exploring Task-wise Knowledge Transfer

Currently, the focus of the CGL research is still on alleviating the forgetting problem. However, as pointed out by [3], [97], besides avoiding the negative transfer (forgetting), we also expect the task-wise transfer to be positive. In other words, the knowledge learnt from one task is expected to be transferable and can benefit the performance of other tasks. By learning transferable features, the catastrophic forgetting problem can also be overcome as a side effect since the task-wise influence becomes beneficial instead of detrimental. Works like HPNs [3] and RieGrace [75] have made preliminary attempts to learn transferable features via decomposing each node into basic features or self-supervised learning, respectively. But more future works are still desired not only for enhanced performance but also for better explainability of the transferable knowledge.

### 7.5 Extension to Multimodal Input and Heterogeneous Graphs

Since graph data is related to various research areas, information from different modalities may accompany the given graphs. For example, images may be attached to products along with their textual name in a co-purchasing network [81]. Therefore, a natural problem is how to incorporate the information from different modalities into the continual learning process. As an early attempt mentioned in Section 4.1, MSCGL [81] approaches the problem by accommodating different modalities with separate GNNs. In the following, we list three promising directions for exploring the multi-modal CGL problem. First, during the learning process, information from all modalities may not

be available during each task. Therefore, if the modality composition of different tasks is largely different, whether more severe forgetting would be caused and how to avoid it is a practical problem. Second, since each same object (*e.g.* a graph node) is described with information of different modalities, information overlapping may exist across different modalities. Therefore, when trying to preserve the representative data, *e.g.* through a memory buffer, how to avoid storing redundant information and preserve the information efficiently is also practically important. Finally, in real-world applications, heterogeneous graphs (characterized by different types of nodes and edges) are more prevalent as they can provide a more complex and nuanced representation of relationships and entities. Therefore, how to extend CGL with multimodal information over heterogeneous graphs is also worthy of investigation.

### 7.6 Beyond Graph Classification Tasks

Currently, existing CGL works are mainly applied to sequences of classification tasks, which can be easily constructed by grouping the data according to the classes. However, there are also other important graph learning scenarios, like the graph regression tasks [122], [123], graph anomaly detection tasks [124], [125], graph generation tasks [126], [127], *etc.*. Different from classification tasks with explicit classes corresponding to data from different distributions, the boundaries between different distributions are implicit, posing challenges for task construction. It is also valuable to explore CGL over those tasks and provide more practical solutions in different scenarios.

## 8 CGL AND GRAPH FOUNDATION MODELS

Large language models (LLMs) have achieved remarkable success in serving as the foundation model for various natural language processing (NLP) downstream tasks. However, graph foundation models are still in their infancy. In this section, we will introduce how CGL techniques can be related to building effective graph foundation models.

## 8.1 Facilitating in Graph Foundation Model Training

To construct a graph foundation model, a promising approach is to design strategies to fine-tune pre-trained LLMs with tokenized graph input, so that the graph structural knowledge is aligned with the knowledge contained in language data. However, such consecutive training would inevitably trigger the catastrophic forgetting problem, since the graph data lie in significantly different domains from the language data used to pre-train the LLMs. Therefore, CGL techniques and the general continual learning techniques become indispensable in this process.

Besides, an ideal graph foundation model, whether developed on top of pre-trained LLMs or trained from scratch, should maintain an up-to-date knowledge base by continuously integrating newly emerged data. However, this continual learning process may also incur the forgetting problem, which is also faced by current LLMs. In this case, CGL techniques should be developed to resolve this issue. Moreover, when the new graph data correspond to different domains, *e.g.* a model may encounter both biological network data and molecule graph data, the cross-domain learning can also trigger the forgetting issue. In this scenario, the domain-IL CGL techniques could be developed to resolve the problem.

## 8.2 A Potential Efficient Way to Develop Large Graph Models

Unlike LLMs, which have demonstrated impressive performance in NLP tasks, the effectiveness of large models for graph data remains a question. Parameter-isolation based CGL models learn on expanding datasets that may significantly increase in size, with the model's size scaling accordingly. Therefore, it provides a potential solution to develop large graph models from scratch by incrementally expanding both the model and the dataset. On the one hand, gradually learning over the tasks instead of optimizing the model for all tasks simultaneously may reduce the optimization difficulty. On the other hand, CGL methods will only allocate new parameters when necessary [3], which could help alleviate the parameter redundancy issue in large models. Besides, the existing large pre-trained graph models are typically domain specific [128], greatly limiting their capability. Fortunately, this challenge may potentially be tackled by leveraging the CGL techniques working under the domain-IL, which aims to train a given graph learning model consecutively across multiple different domains.

## REFERENCES

[1] X. Zhang, D. Song, and D. Tao, "Cglb: Benchmark tasks for continual graph learning," in *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.

[2] H. Liu, Y. Yang, and X. Wang, "Overcoming catastrophic forgetting in graph neural networks," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 10, 2021, pp. 8653–8661.

[3] X. Zhang, D. Song, and D. Tao, "Hierarchical prototype networks for continual graph representation learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4622–4636, 2023.

[4] L. Galke, I. Vagliano, B. Franke, T. Zielke, M. Hoffmann, and A. Scherp, "Lifelong learning on evolving graphs under the constraints of imbalanced classes and new classes," *Neural Networks*, vol. 164, pp. 156–176, 2023.

[5] Q. Shen, W. Ren, and W. Qin, "Graph relation aware continual learning," *arXiv preprint arXiv:2308.08259*, 2023.

[6] G. Lombardo, A. Poggi, and M. Tomaiuolo, "Continual representation learning for node classification in power-law graphs," *Future Generation Computer Systems*, vol. 128, pp. 420–428, 2022.

[7] X. Kou, Y. Lin, S. Liu, P. Li, J. Zhou, and Y. Zhang, "Disentangle-based continual graph representation learning," *arXiv preprint arXiv:2010.02565*, 2020.

[8] A. Carta, A. Cossu, F. Errica, and D. Bacciu, "Catastrophic forgetting in deep graph networks: A graph classification benchmark," *Frontiers in artificial intelligence*, vol. 5, p. 824655, 2022.

[9] M. Mirtaheri, M. Rostami, and A. Galstyan, "History repeats: Overcoming catastrophic forgetting for event-centric temporal knowledge graph completion," *arXiv preprint arXiv:2305.18675*, 2023.

[10] Q. Yuan, S.-U. Guan, P. Ni, T. Luo, K. L. Man, P. Wong, and V. Chang, "Continual graph learning: A survey," *arXiv preprint arXiv:2301.12230*, 2023.

[11] F. G. Febrinanto, F. Xia, K. Moore, C. Thapa, and C. Aggarwal, "Graph lifelong learning: A survey," *IEEE Computational Intelligence Magazine*, vol. 18, no. 1, pp. 32–51, 2023.

[12] H.-J. Song and S.-B. Park, "Enriching translation-based knowledge graph embeddings through continual learning," *IEEE Access*, vol. 6, pp. 60 489–60 497, 2018.

[13] Y. Ren, L. Ke, D. Li, H. Xue, Z. Li, and S. Zhou, "Incremental graph classification by class prototype construction and augmentation," in *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, 2023, pp. 2136–2145.

[14] H. Lin, R. Jia, and X. Lyu, "Gated attention with asymmetric regularization for transformer-based continual graph learning," in *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2023, pp. 2021–2025.

[15] K. Javed and M. White, "Meta-learning representations for continual learning," *Advances in neural information processing systems*, vol. 32, 2019.

[16] X. Chen, J. Wang, and K. Xie, "Trafficstream: A streaming traffic flow forecasting framework based on graph neural networks and continual learning," *arXiv preprint arXiv:2106.06273*, 2021.

[17] Z. Tan, M. Hu, Y. Wang, L. Wei, and B. Liu, "Futures quantitative investment with heterogeneous continual graph neural network," *arXiv preprint arXiv:2303.16532*, 2023.

[18] M. Perini, G. Ramponi, P. Carbone, and V. Kalavri, "Learning on streaming graphs with experience replay," in *Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing*, 2022, pp. 470–478.

[19] S. Gupta, S. Manchanda, S. Ranu, and S. J. Bedathur, "Grafenne: learning on graphs with heterogeneous and dynamic feature sets," in *International Conference on Machine Learning*. PMLR, 2023, pp. 12 165–12 181.

[20] X. Chen, J. Zhang, X. Wang, T. Wu, S. Deng, Y. Wang, L. Si, H. Chen, and N. Zhang, "Continual multimodal knowledge graph construction," *arXiv preprint arXiv:2305.08698*, 2023.

[21] B. Wang, Y. Zhang, X. Wang, P. Wang, Z. Zhou, L. Bai, and Y. Wang, "Pattern expansion and consolidation on evolving graphs for continual traffic prediction," in *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2023, pp. 2223–2232.

[22] J. Omeliyanenko, A. Zehe, A. Hotho, and D. Schlör, "Capskg: Enabling continual knowledge integration in language models for automatic knowledge graph completion," in *International Semantic Web Conference*. Springer, 2023, pp. 618–636.

[23] Y. Wang, Y. Cui, W. Liu, Z. Sun, Y. Jiang, K. Han, and W. Hu, "Facing changes: Continual entity alignment for growing knowledge graphs," in *International Semantic Web Conference*. Springer, 2022, pp. 196–213.

[24] G. M. Van de Ven and A. S. Tolias, "Three scenarios for continual learning," *arXiv preprint arXiv:1904.07734*, 2019.

[25] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, and S. Wermter, "Continual lifelong learning with neural networks: A review," *Neural networks*, vol. 113, pp. 54–71, 2019.

[26] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.

[27] M. Masana, X. Liu, B. Twardowski, M. Menta, A. D. Bagdanov, and J. van de Weijer, "Class-incremental learning: Survey and performance evaluation on image classification," *IEEE Transactions*

*on Pattern Analysis and Machine Intelligence*, vol. 45, no. 5, pp. 5513–5533, 2023.

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

[29] C. Zheng, B. Zong, W. Cheng, D. Song, J. Ni, W. Yu, H. Chen, and W. Wang, "Robust graph representation learning via neural sparsification," in *Proceedings of the 37th International Conference on Machine Learning (ICML)*, July 2020.

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

[31] L. Wang, B. Zong, Q. Ma, W. Cheng, J. Ni, W. Yu, Y. Liu, D. Song, H. Chen, and Y. Fu, "Inductive and unsupervised representation learning on graph structured objects," in *Proceedings of International Conference on Learning Representations (ICLR)*, April 2020.

[32] C. Zhang*, D. Song, C. Huang, A. Swami, and N. V. Chawla, "Heterogeneous graph neural network," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, Anchorage, AK, USA, August 2019, pp. 793–803.

[33] X. Zhang, C. Xu, and D. Tao, "Context aware graph convolution for skeleton-based action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 333–14 342.

[34] W. Yu, C. Zhen, W. Cheng, C. Aggarwal, D. Song, B. Zong, H. Chen, and W. Wang, "Learning deep network representations with adversarially regularized autoencoders," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, London, UK, August 2018, pp. 2,663–2,671.

[35] H. Lei, N. Akhtar, and A. Mian, "Spherical kernel for efficient graph convolution on 3d point clouds," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[36] Y. Xie, Z. Xu, J. Zhang, Z. Wang, and S. Ji, "Self-supervised learning of graph neural networks: A unified review," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 2, pp. 2412–2429, 2023.

[37] R. Zhang, Y. Zhang, C. Lu, and X. Li, "Unsupervised graph embedding via adaptive graph learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 5329–5336, 2023.

[38] C. Zheng, B. Zong, W. Cheng, D. Song, J. Ni, W. Yu, H. Chen, and W. Wang, "Node classification in temporal graphs through stochastic sparsification and temporal structural convolution," in *Proceedings of European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML-PKDD)*, September 2020.

[39] L. Lü and T. Zhou, "Link prediction in complex networks: A survey," *Physica A: statistical mechanics and its applications*, vol. 390, no. 6, pp. 1150–1170, 2011.

[40] C. Zhang, H. Yao, L. Yu, C. Huang, D. Song, M. Jiang, H. Chen, and N. V. Chawla, "Inductive contextual relation learning for personalization," *ACM Transactions on Information Systems (TOIS)*, vol. 39, no. 35, pp. 1–22, 2021.

[41] L. Cai, J. Li, J. Wang, and S. Ji, "Line graph neural networks for link prediction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 9, pp. 5103–5113, 2021.

[42] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Neural message passing for quantum chemistry," in *International Conference on Machine Learning*. PMLR, 2017, pp. 1263–1272.

[43] J. T. Vogelstein, W. G. Roncal, R. J. Vogelstein, and C. E. Priebe, "Graph classification using signal-subgraphs: Applications in statistical connectomics," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 7, pp. 1539–1551, 2012.

[44] T. T. Mueller, J. C. Paetzold, C. Prabhakar, D. Usynin, D. Rueckert, and G. Kaissis, "Differentially private graph neural networks for whole-graph classification," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 6, pp. 7308–7318, 2023.

[45] K. Du, F. Lyu, L. Li, F. Hu, W. Feng, F. Xu, X. Xi, and H. Cheng, "Multi-label continual learning using augmented graph convolutional network," *IEEE Transactions on Multimedia*, 2023.

[46] H. Bo, R. McConville, J. Hong, and W. Liu, "Ego-graph replay based continual learning for misinformation engagement prediction," in *2022 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2022, pp. 01–08.

[47] L. Hedegaard, N. Heidari, and A. Iosifidis, "Continual spatio-temporal graph convolutional networks," *Pattern Recognition*, vol. 140, p. 109528, 2023.

[48] A. Zaman, F. Yangyu, M. S. Ayub, M. Irfan, L. Guoyun, and L. Shiya, "Cmdgat: Knowledge extraction and retention based continual graph attention network for point cloud registration," *Expert Systems with Applications*, vol. 214, p. 119098, 2023.

[49] L.-P. Xhonneux, M. Qu, and J. Tang, "Continuous graph neural networks," in *International Conference on Machine Learning*. PMLR, 2020, pp. 10 432–10 441.

[50] Y. Luo, Z. Huang, Z. Zhang, Z. Wang, M. Baktashmotlagh, and Y. Yang, "Learning from the past: continual meta-learning with bayesian graph neural networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 5021–5028.

[51] B. Tang and D. S. Matteson, "Graph-based continual learning," *arXiv preprint arXiv:2007.04813*, 2020.

[52] B. Das and E. Isufi, "Graph filtering over expanding graphs," in *2022 IEEE Data Science and Learning Workshop (DSLW)*. IEEE, 2022, pp. 1–8.

[53] ——, "Learning expanding graphs for signal interpolation," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 5917–5921.

[54] ——, "Online filtering over expanding graphs," in *2022 56th Asilomar Conference on Signals, Systems, and Computers*. IEEE, 2022, pp. 43–47.

[55] T. Wu, Q. Liu, Y. Cao, Y. Huang, X.-M. Wu, and J. Ding, "Continual graph convolutional network for text classification," *arXiv preprint arXiv:2304.04152*, 2023.

[56] L. Galke, B. Franke, T. Zielke, and A. Scherp, "Lifelong learning of graph neural networks for open-world node classification," in *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2021, pp. 1–8.

[57] J. Wang, G. Song, Y. Wu, and L. Wang, "Streaming graph neural networks via continual learning," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 1515–1524.

[58] Y. Han, S. Karunasekera, and C. Leckie, "Graph neural networks with continual learning for fake news detection from social media," *arXiv preprint arXiv:2007.03316*, 2020.

[59] W. Yu, W. Cheng, C. C. Aggarwal, K. Zhang, H. Chen, and W. Wang, "Netwalk: A flexible deep embedding approach for anomaly detection in dynamic networks," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2672–2681.

[60] G. H. Nguyen, J. B. Lee, R. A. Rossi, N. K. Ahmed, E. Koh, and S. Kim, "Continuous-time dynamic network embeddings," in *Companion Proceedings of the The Web Conference 2018*, 2018, pp. 969–976.

[61] L. Zhou, Y. Yang, X. Ren, F. Wu, and Y. Zhuang, "Dynamic network embedding by modeling triadic closure process," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[62] Y. Ma, Z. Guo, Z. Ren, J. Tang, and D. Yin, "Streaming graph neural networks," in *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 719–728.

[63] Y. Feng, J. Jiang, and Y. Gao, "Incremental learning on growing graphs," 2020.

[64] P. Bielak, K. Tagowski, M. Falkiewicz, T. Kajdanowicz, and N. V. Chawla, "Fildne: A framework for incremental learning of dynamic networks embeddings," *Knowledge-Based Systems*, vol. 236, p. 107453, 2022.

[65] B. He, X. He, Y. Zhang, R. Tang, and C. Ma, "Dynamically expandable graph convolution for streaming recommendation," *arXiv preprint arXiv:2303.11700*, 2023.

[66] F. Zhou, C. Cao, K. Zhang, G. Trajcevski, T. Zhong, and J. Geng, "Meta-gnn: On few-shot node classification in graph meta-learning," in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*, 2019, pp. 2357–2360.

[67] Z. Guo, C. Zhang, W. Yu, J. Herr, O. Wiest, M. Jiang, and N. V. Chawla, "Few-shot graph learning for molecular property prediction," in *Proceedings of the Web Conference 2021*, 2021, pp. 2559–2567.

[68] H. Yao, C. Zhang, Y. Wei, M. Jiang, S. Wang, J. Huang, N. Chawla, and Z. Li, "Graph few-shot learning via knowledge transfer," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 04, 2020, pp. 6656–6663.

[69] Z. Tan, K. Ding, R. Guo, and H. Liu, "Graph few-shot class-incremental learning," in *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, 2022, pp. 987–996.

[70] V. Garcia and J. Bruna, "Few-shot learning with graph neural networks," *arXiv preprint arXiv:1711.04043*, 2017.

[71] F. G. Febrinanto, F. Xia, K. Moore, C. Thapa, and C. Aggarwal, "Graph lifelong learning: A survey," *IEEE Computational Intelligence Magazine*, vol. 18, no. 1, pp. 32–51, 2023.

[72] T. D. Hoang, D. V. Tung, D.-H. Nguyen, B.-S. Nguyen, H. H. Nguyen, and H. Le, "Universal graph continual learning," *Transactions on Machine Learning Research*, 2023. [Online]. Available: https://openreview.net/forum?id=wzRE5kTnl3

[73] Y. Liu, R. Qiu, and Z. Huang, "Cat: Balanced continual graph learning with graph condensation," *arXiv preprint arXiv:2309.09455*, 2023.

[74] P. Zhang, Y. Yan, C. Li, S. Wang, X. Xie, G. Song, and S. Kim, "Continual learning on dynamic graphs via parameter isolation," *arXiv preprint arXiv:2305.13825*, 2023.

[75] L. Sun, J. Ye, H. Peng, F. Wang, and S. Y. Philip, "Self-supervised continual graph learning in adaptive riemannian spaces," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 4, 2023, pp. 4633–4642.

[76] X. Zhang, D. Song, and D. Tao, "Hierarchical prototype networks for continual graph representation learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[77] ——, "Sparsified subgraph memory for continual graph representation learning," in *2022 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2022.

[78] ——, "Ricci curvature-based graph sparsification for continual graph representation learning," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[79] X. Zhang, D. Song, Y. Chen, and D. Tao, "Topology-aware embedding memory for learning on expanding networks," *arXiv preprint arXiv:2401.13200*, 2024.

[80] Y. Xu, Y. Zhang, W. Guo, H. Guo, R. Tang, and M. Coates, "Graphsail: Graph structure aware incremental learning for recommender systems," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 2861–2868.

[81] J. Cai, X. Wang, C. Guan, Y. Tang, J. Xu, B. Zhong, and W. Zhu, "Multimodal continual graph learning with neural architecture search," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 1292–1300.

[82] Y. Cui, Y. Wang, Z. Sun, W. Liu, Y. Jiang, K. Han, and W. Hu, "Lifelong embedding learning and transfer for growing knowledge graphs," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 4, 2023, pp. 4217–4224.

[83] F. Zhou and C. Cao, "Overcoming catastrophic forgetting in graph neural networks with experience replay," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 5, 2021, pp. 4714–4722.

[84] D. Wei, Y. Gu, Y. Song, Z. Song, F. Li, and G. Yu, "Incregnn: Incremental graph neural network learning by considering node and parameter importance," in *Database Systems for Advanced Applications: 27th International Conference, DASFAA 2022, Virtual Event, April 11–14, 2022, Proceedings, Part I*. Springer, 2022, pp. 739–746.

[85] A. Rakaraddi, L. Siew Kei, M. Pratama, and M. De Carvalho, "Reinforced continual learning for graphs," in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, 2022, pp. 1666–1674.

[86] J. Wang, W. Zhu, G. Song, and L. Wang, "Streaming graph neural networks with generative replay," in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 1878–1888.

[87] K. Ahrabian, Y. Xu, Y. Zhang, J. Wu, Y. Wang, and M. Coates, "Structure aware experience replay for incremental learning in graph-based recommender systems," in *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, 2021, pp. 2832–2836.

[88] J. Su and C. Wu, "Towards robust inductive graph incremental learning via experience replay," *arXiv preprint arXiv:2302.03534*, 2023.

[89] S. Kim, S. Yun, and J. Kang, "Dygrain: An incremental learning framework for dynamic graphs," in *31st International Joint Conference on Artificial Intelligence, IJCAI 2022*. International Joint Conferences on Artificial Intelligence, 2022, pp. 3157–3163.

[90] A. Daruna, M. Gupta, M. Sridharan, and S. Chernova, "Continual learning of knowledge graph embeddings," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1128–1135, 2021.

[91] J. Kirkpatrick, R. Pascanu, N. Rabinowitz, J. Veness, G. Desjardins, A. A. Rusu, K. Milan, J. Quan, T. Ramalho, A. Grabska-Barwinska *et al.*, "Overcoming catastrophic forgetting in neural networks," *Proceedings of the National Academy of Sciences*, vol. 114, no. 13, pp. 3521–3526, 2017.

[92] R. Aljundi, F. Babiloni, M. Elhoseiny, M. Rohrbach, and T. Tuytelaars, "Memory aware synapses: Learning what (not) to forget," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 139–154.

[93] J. Topping, F. Di Giovanni, B. P. Chamberlain, X. Dong, and M. M. Bronstein, "Understanding over-squashing and bottlenecks on graphs via curvature," *arXiv preprint arXiv:2111.14522*, 2021.

[94] Z. Ye, K. S. Liu, T. Ma, J. Gao, and C. Chen, "Curvature graph network," in *ICLR*, 2019.

[95] R. Pasunuru and M. Bansal, "Continual and multi-task architecture search," *arXiv preprint arXiv:1906.05226*, 2019.

[96] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," *Advances in neural information processing systems*, vol. 26, 2013.

[97] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *Advances in Neural Information Processing Systems*, 2017, pp. 6467–6476.

[98] F. R. Hampel, E. M. Ronchetti, P. Rousseeuw, and W. A. Stahel, *Robust statistics: the approach based on influence functions*. Wiley-Interscience; New York, 1986.

[99] A. Gretton, A. Smola, J. Huang, M. Schmittfull, K. Borgwardt, and B. Schölkopf, "Covariate shift by kernel mean matching," *Dataset shift in machine learning*, vol. 3, no. 4, p. 5, 2009.

[100] P. Buzzega, M. Boschini, A. Porrello, D. Abati, and S. Calderara, "Dark experience for general continual learning: a strong, simple baseline," *Advances in neural information processing systems*, vol. 33, pp. 15 920–15 930, 2020.

[101] A. A. Rusu, N. C. Rabinowitz, G. Desjardins, H. Soyer, J. Kirkpatrick, K. Kavukcuoglu, R. Pascanu, and R. Hadsell, "Progressive neural networks," *arXiv preprint arXiv:1606.04671*, 2016.

[102] M. Wortsman, V. Ramanujan, R. Liu, A. Kembhavi, M. Rastegari, J. Yosinski, and A. Farhadi, "Supermasks in superposition," *arXiv preprint arXiv:2006.14769*, 2020.

[103] C. Wang, Y. Qiu, D. Gao, and S. Scherer, "Lifelong graph learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 13 719–13 728.

[104] H. Shin, J. K. Lee, J. Kim, and J. Kim, "Continual learning with deep generative replay," in *Advances in Neural Information Processing Systems*, 2017, pp. 2990–2999.

[105] A. Chaudhry, P. K. Dokania, T. Ajanthan, and P. H. Torr, "Riemannian walk for incremental learning: Understanding forgetting and intransigence," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 532–547.

[106] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, and X. Bresson, "Benchmarking graph neural networks," 2020.

[107] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, and S. Süsstrunk, "Slic superpixels compared to state-of-the-art superpixel methods," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 11, pp. 2274–2282, 2012.

[108] W. Hu, M. Fey, M. Zitnik, Y. Dong, H. Ren, B. Liu, M. Catasta, and J. Leskovec, "Open graph benchmark: Datasets for machine learning on graphs," *Advances in neural information processing systems*, vol. 33, pp. 22 118–22 133, 2020.

[109] A. K. McCallum, K. Nigam, J. Rennie, and K. Seymore, "Automating the construction of internet portals with machine learning," *Information Retrieval*, vol. 3, no. 2, pp. 127–163, 2000.

[110] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.

[111] 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.

[112] Z. Xiong, D. Wang, X. Liu, F. Zhong, X. Wan, X. Li, Z. Li, X. Luo, K. Chen, H. Jiang *et al.*, "Pushing the boundaries of molecular representation for drug discovery with the graph attention mechanism," *Journal of medicinal chemistry*, vol. 63, no. 16, pp. 8749–8760, 2019.

[113] 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.

[114] O. Shchur, M. Mumme, A. Bojchevski, and S. Günnemann, "Pitfalls of graph neural network evaluation," *arXiv preprint arXiv:1811.05868*, 2018.

[115] J. Tang, J. Sun, C. Wang, and Z. Yang, "Social influence analysis in large-scale networks," in *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge discovery and data mining*, 2009, pp. 807–816.

[116] E. Cho, S. A. Myers, and J. Leskovec, "Friendship and mobility: user movement in location-based social networks," in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2011, pp. 1082–1090.

[117] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2001–2010.

[118] A. Chaudhry, M. Rohrbach, M. Elhoseiny, T. Ajanthan, P. K. Dokania, P. H. Torr, and M. Ranzato, "On tiny episodic memories in continual learning," *arXiv preprint arXiv:1902.10486*, 2019.

[119] A. Prabhu, P. H. Torr, and P. K. Dokania, "Gdumb: A simple approach that questions our progress in continual learning," in *European conference on computer vision*. Springer, 2020, pp. 524–540.

[120] J. Yoon, E. Yang, J. Lee, and S. J. Hwang, "Lifelong learning with dynamically expandable networks," *arXiv preprint arXiv:1708.01547*, 2017.

[121] J. Su, D. Zou, Z. Zhang, and C. Wu, "Towards robust graph incremental learning on evolving graphs," in *International Conference on Machine Learning*. PMLR, 2023, pp. 32 728–32 748.

[122] J. Jia and A. R. Benson, "Residual correlation in graph neural network regression," in *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 588–598.

[123] H. Maron, H. Ben-Hamu, H. Serviansky, and Y. Lipman, "Provably powerful graph networks," *Advances in neural information processing systems*, vol. 32, 2019.

[124] L. Akoglu, H. Tong, and D. Koutra, "Graph based anomaly detection and description: a survey," *Data mining and knowledge discovery*, vol. 29, pp. 626–688, 2015.

[125] X. Ma, J. Wu, S. Xue, J. Yang, C. Zhou, Q. Z. Sheng, H. Xiong, and L. Akoglu, "A comprehensive survey on graph anomaly detection with deep learning," *IEEE Transactions on Knowledge and Data Engineering*, 2021.

[126] X. Guo and L. Zhao, "A systematic survey on deep generative models for graph generation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 5, pp. 5370–5390, 2023.

[127] X. Chang, P. Ren, P. Xu, Z. Li, X. Chen, and A. Hauptmann, "A comprehensive survey of scene graphs: Generation and application," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 1–26, 2023.

[128] J. Xia, Y. Zhu, Y. Du, and S. Z. Li, "A survey of pretraining on graphs: Taxonomy, methods, and applications," *arXiv preprint arXiv:2202.07893*, 2022.