

Lifelong Learning of Graph Neural Networks for Open-World Node Classification

Lukas Galke
Kiel University / ZBW, Germany
lga@informatik.uni-kiel.de

Benedikt Franke, Tobias Zielke, Ansgar Scherp
Ulm University, Germany
{benedikt.franke,tobias-1.zielke,ansgar.scherp}@uni-ulm.de

Abstract—Graph neural networks (GNNs) have emerged as the standard method for numerous tasks on graph-structured data such as node classification. However, real-world graphs are often evolving over time and even new classes may arise. We model these challenges as an instance of lifelong learning, in which a learner faces a sequence of tasks and may take over knowledge acquired in past tasks. Such knowledge may be stored explicitly as historic data or implicitly within model parameters. In this work, we systematically analyze the influence of implicit and explicit knowledge. Therefore, we present an incremental training method for lifelong learning on graphs and introduce a new measure based on k -neighborhood time differences to address variances in the historic data. We apply our training method to five representative GNN architectures and evaluate them on three new lifelong node classification datasets. Our results show that no more than 50% of the GNN’s receptive field is necessary to retain at least 95% accuracy compared to training over the complete history of the graph data. Furthermore, our experiments confirm that implicit knowledge becomes more important when fewer explicit knowledge is available.

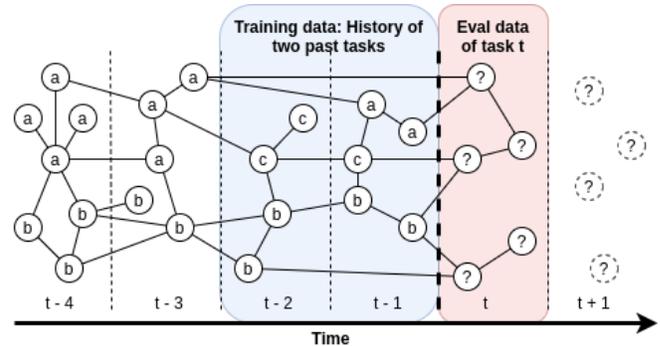


Fig. 1. Lifelong Open-World Node Classification. At each time t the learner has to classify new vertices of task \mathcal{T}_t (red). The learner may use knowledge from previous tasks to adapt to the current task, eventually cut off by a history size (blue). The current task might come with previously unseen classes, e. g., see class “c” that emerged only at task $t - 2$ and was subsequently added to the class set. After evaluating each task \mathcal{T}_t , we continue with task \mathcal{T}_{t+1} .

I. INTRODUCTION

Graph neural networks [1] (GNNs) have emerged as state-of-the-art methods in numerous tasks on graph-structured data such as vertex classification [2]–[5], graph classification [6], link prediction [7], and unsupervised vertex representation learning [8]. An intriguing property of GNNs is that they are capable of inductive learning. An inductive model for graph data only depends on the vertex features and the graph structure given by its edges. In many cases [3], [9], [10], this is a major advantage over models that rely on a static vertex embedding [11], which would need to be retrained [12] as soon as any new vertex appears. In contrast, inductively trained GNNs can be applied to new data – or even a different graph – without any retraining because of the property of only requiring vertex features and edges.

However, being able to apply the same model to unseen data also comes with challenges that have not been analyzed so far. Let us assume, we have a vertex classification model and new data streams in over time, i. e., new edges and vertices arrive; then even *new classes* may arise! This raises several new questions: Do we need to retrain the model? How much past data should be preserved for retraining? Is it helpful to preserve implicit knowledge within the model parameters or should we retrain from scratch?

To answer these questions, we frame the problem as an instance of lifelong machine learning [13]–[15]. In lifelong

learning, the learner has to perform a sequence of tasks $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_t$, and may use knowledge \mathcal{K} gained in previous tasks to perform task \mathcal{T}_t (see illustration in Figure 1). In our case, each task consists of classifying vertices given an attributed graph. Knowledge \mathcal{K} may be stored explicitly (the training data of past tasks) or implicitly within the model parameters. A particular challenge of lifelong learning in the context of graph data is that vertices cannot be processed independently because models typically take connected vertices into account. We also consider the challenge that the set of classes in task \mathcal{T}_t differs from classes in previous tasks, which is known as the *open-world classification* [15] problem.

We address these challenges by introducing a new incremental training method for lifelong learning on graph data. This training method can be applied on any dataset by the introduction of a new measure to harmonize temporal variances in the evolution of graph datasets with different characteristics and granularities. The method is capable to integrate various existing GNNs, both isotropic and anisotropic as well as scalable GNN methods. In our experiments, we thoroughly evaluate representative and scalable GNN architectures (and one graph-agnostic multi-layer perceptron) using our incremental training method that retrains the model for each task. We use an artificial *history size* that limits the amount of past data (called here: *explicit knowledge*) available for training

and compare *limited history-size* retraining against unlimited *full-history* retraining.

Furthermore, we compare reusing model parameters from previous tasks (*warm restart*) against retraining from scratch (*cold restarts*) to analyze the influence of *implicit knowledge*. In an ablation study, we compare incrementally training models against once-trained models. To facilitate our analyses, we contribute three new datasets for lifelong learning; one co-authorship and two citation graph datasets with different degrees of changes in the class set. In total, we have experimented with 48 different incremental training configurations, namely 6 architectures \times 4 history sizes \times cold restarts and warm restarts, which we evaluate on three new datasets. To enable a fair comparison, we tune the hyperparameters for each configuration separately. We repeat all experiments 10 times with different random seeds.

The results of our 1,440 experiments (48 configurations repeated 10 times on 3 datasets) reveal the following insights:

(i) *Data from only few past tasks are sufficient for retraining* Surprisingly, only the data from very few past tasks are sufficient to maintain a high level of accuracy that matches the accuracy of the same model retrained using all past data. In our experiments, merely 50% of the GNN's receptive field (corresponding to history sizes of 3 or 4 past tasks) are sufficient to match at least 95% of the accuracy of the same model trained on the full history.

(ii) *Reusing parameters requires less data* Tuning the learning rate may compensate for the knowledge lost when randomly re-initializing the model. However, when the history size is small, using warm restarts tends to be beneficial such that more knowledge from past tasks can be preserved implicitly. Our results further suggest that using warm restarts is generally preferable for lifelong learning.

(iii) *Incremental training is a simple yet effective technique to tackle lifelong learning on graphs* By comparing incrementally trained models with once-trained models, we find that the accuracy of once-trained models decreases quickly when new classes appear. We note that even when no new classes appear over time, incremental training still benefits from the increased amount of training data.

These insights have direct consequences for using GNNs in practical applications. It allows to decide how much historical data should be kept to maintain a certain performance versus having memory available in the GPUs. This is an important criterion that influences which GNN methods are applicable [10], [16]. We publicly provide the datasets as well as the evaluation framework to extend our experiments and accelerate research in lifelong learning on graphs.

Below, we discuss the related works and provide a problem statement with our training procedure and new measure in Section III. We describe the experimental apparatus and datasets in Section IV. The results of our experiments are reported in Section V as well as in Section VI for our ablation study. We discuss the results in Section VII before we conclude.

II. RELATED WORK

We discuss the works in lifelong machine learning and especially lifelong learning on graphs. Subsequently, we discuss methods for evolving graphs as well as methods regarding history sizes. GNN architectures will be discussed in more detail in Section III-D along with our selection of representative architectures.

Lifelong learning (sometimes called: continual learning [17]) is present in machine learning research since the mid 1990s [13], [18]–[20]. The goal of lifelong learning is to develop approaches that can adapt to new tasks. Thus, it differs from online learning [21], which focuses on efficiency. ELLA [22] introduces an efficient lifelong learning algorithm with convergence guarantees and employs multi-task learning such that future tasks can improve previous tasks. Fei et al. [14] analyze SVMs in a lifelong learning setting and introduce cumulative learning. Cumulative learning relates to our approach, as we consider that some data is shared among the tasks. However, we further investigate *how much* past data is necessary to retain accuracy compared to a fully-cumulative approach. Lopez-Paz & Ranzato [17] introduce a gradient episodic memory framework for the image domain, where examples can be processed independently, and tackle the catastrophic forgetting problem, i. e., the loss of previously learned information when new information is learned [23]. A recent overview of lifelong learning can be found in [15].

In the field of *open-world classification*, several methods have been proposed to dynamically detect novel classes for classic machine learning methods [14], [24] as well as graph neural networks [25]. For now, we assume that the model is retrained for each task. We lay the foundation for self-detection of new classes by providing our datasets along with a simple yet effective training scheme for lifelong learning on graphs.

Related works on *lifelong learning on graphs* are very limited. Concurrent with our research, there is one recent work by Wang et al. [26] that also explores GNNs on lifelong learning problems. The authors focus on catastrophic forgetting. To tackle this challenge, the authors explore a method to cast vertex classification as a graph classification task by transforming each vertex of the graph into a feature graph. This way, the vertices become independent such that they can follow the streaming setup from [17]. Apart from lifelong learning, also other methods have been developed to make vertices independent via preprocessing: SIGN [27] and Simplified GCN [16]. We include the latter in our experiments.

In the past, also different methods for *evolving graphs* have been proposed including dynamic embedding methods [28], [29], autoencoder-based methods [30], [31], GNNs for graphs with fixed vertex set [32]–[38], and inductive GNN methods that can deal with previously unseen vertices [9], [39]. However, these methods are not relevant to our work because we use the time information merely to induce the sequence of lifelong graph learning tasks (see Section IV-A). When a vertex is preserved from task t to $t + 1$, it will have also the same features and label(s). Thus, our sequence of tasks does

not require sophisticated long-horizon temporal information. Instead, in our work we require *the ability of the models to adapt to new classes* during a sequence of tasks.

Regarding finding the optimal *history size in data streams*, Fish and Caceres [40] treated the window size as a hyperparameter and proposed an optimization algorithm which requires multiple runs of the model. This is a rather costly procedure and the study does not yield insights on how much predictive power can be preserved when selecting a near-optimal but much smaller, and thus more efficient, window size. Other works, e. g., [41], indicate that smaller history sizes might be beneficial in some scenarios. However, there is no systematic study of the influence of history sizes in lifelong learning on graphs.

To summarize, lifelong learning on graphs is a surprisingly unexplored topic. In particular, none of discussed works analyzes the problem of *open-world classification* in graph data and how much past training data is necessary – or how few is enough – to maintain good predictive power.

III. PROBLEM STATEMENT AND METHODS

We first outline our problem formalization of lifelong learning on graphs. Subsequently, we introduce our incremental training procedure as well as our method to harmonize window sizes across different datasets. Finally, we briefly describe the base models that we incrementally train for our experiments.

A. Problem Formalization

We define our problem of open-world classification of graph vertices as a form of lifelong learning [14]. It is defined as:

Definition 1 (Lifelong Learning [14]). At any time t , the learner has performed a sequence of t learning tasks, $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_t$ and has accumulated the knowledge \mathcal{K} learned in these past tasks. At time $t+1$, it is faced with a new learning task \mathcal{T}_{t+1} . The learner is able to make use of past knowledge to help perform the new learning task \mathcal{T}_{t+1} .

We cast this definition into a lifelong graph learning problem by considering each task $\mathcal{T}_t := (\mathcal{G}_t, \mathbf{X}^{(t)}, \mathbf{y}^{(t)})$ to be a vertex classification task with graph $\mathcal{G}_t = (V_t, E_t)$, corresponding vertex features $\mathbf{X}^{(t)} \in \mathbb{R}^{|V_t| \times D}$, and vertex labels $\mathbf{y}^{(t)} \in \mathbb{N}^{|V_t|}$. We denote the set of all classes at time t as \mathbb{Y}_t . To ensure that past knowledge is helpful to perform \mathcal{T}_t , we impose $\mathcal{G}_{t-1} \cap \mathcal{G}_t \neq \emptyset$. We assume that the features and labels of the vertices do not change: $\mathbf{X}_u^{(t-1)} = \mathbf{X}_u^{(t)}$, $\mathbf{y}_u^{(t-1)} = \mathbf{y}_u^{(t)}$ if $u \in V_{t-1} \cap V_t$. Such changes can still be modeled by inserting a new vertex and removing the old one. The task is to predict the class labels for new vertices $V_t \setminus V_{t-1}$. Please note that these vertices may come with new, unseen classes as \mathbb{Y}_t may differ from \mathbb{Y}_{t-1} . Furthermore, we analyze the effect of a history size c , which limits the available past data. We call this past data *explicit knowledge*. In this case, we set $\tilde{\mathcal{T}}_t := (\tilde{\mathcal{G}}_t, \tilde{\mathbf{X}}^{(t)}, \tilde{\mathbf{y}}^{(t)})$ with $\tilde{\mathcal{G}}_t := \mathcal{G}_t \setminus (\mathcal{G}_1 \cup \mathcal{G}_2 \dots \cup \mathcal{G}_{t-c-1})$, and remove corresponding features and labels to construct $\tilde{\mathbf{X}}_t$ and $\tilde{\mathbf{y}}_t$. Still, *implicit knowledge* acquired in past tasks, e. g., within the model parameters, may be used for task $\tilde{\mathcal{T}}_t$.

B. Incremental Training for Lifelong Learning on Graphs

Without loss of generality, we assume to have a finite sequence of T tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ and a model f with parameters θ . Over the course of the tasks, the graph changes, including its vertices, edges, as well as the set of classes. To address these changes, we explore a simple yet effective incremental training technique for adapting neural networks to new graph-structured tasks. As preparation for task \mathcal{T}_{t+1} , we retrain f on the labels of \mathcal{T}_t to obtain $\theta^{(t)}$. Whenever l new classes appear in the training data, we add a corresponding amount of parameters to the output layer of $f^{(t)}$. Therefore, we have $|\theta_{\text{output weights}}^{(t)}| = |\theta_{\text{output weights}}^{(t-1)}| + l$ and $|\theta_{\text{output bias}}^{(t)}| = |\theta_{\text{output bias}}^{(t-1)}| + l$. Those parameters that are specific to new classes are newly initialized. For the other parameters, we consider two options in our incremental training procedure: *warm restarts* and *cold restarts*. With *cold restarts*, we re-initialize $\theta^{(t)}$ and retrain from scratch. In contrast, when using *warm restarts*, we initialize the parameters for training on task \mathcal{T}_t with the final parameters of the previous task $\theta^{(t-1)}$. Algorithm 1 outlines our incremental training procedure.

Algorithm 1: Incremental training for lifelong graph learning under cold-start vs. warm-start condition

Input : Sequence of tasks $\tilde{\mathcal{T}}_0, \dots, \tilde{\mathcal{T}}_T$, model f with parameters θ , flag for cold or warm restarts

Output: Predicted labels for new vertices of each task

```

1 known_classes  $\leftarrow \emptyset$ ;
2  $\theta \leftarrow \text{initialize\_parameters}()$ ;
3 for  $t \leftarrow 1$  to  $T$  do
4   new_classes  $\leftarrow \text{set}(\tilde{\mathbf{y}}^{(t-1)}) \setminus \text{known\_classes}$ ;
5   if new_classes  $\neq \emptyset$  then
6      $\theta' \leftarrow \text{expand\_output\_layer}(\theta, |\text{new\_classes}|)$ ;
7   end
8    $\theta' \leftarrow \text{initialize\_parameters}()$ ;
9   if  $t > 1$  and do_warm_restart = TRUE then
10     $\theta' \leftarrow \text{copy\_existing\_parameters}(\theta)$ ;
11  end
12   $\theta' \leftarrow \text{train}(\theta', \tilde{\mathcal{G}}_{t-1}, \tilde{\mathbf{X}}^{(t-1)}, \tilde{\mathbf{y}}^{(t-1)})$ ;
13   $\tilde{\mathbf{y}}_{\text{pred}} \leftarrow \text{predict}(\theta', \tilde{\mathcal{G}}_t, \tilde{\mathbf{X}}^{(t)})$  for vertices  $V_t \setminus V_{t-1}$ ;
14  known_classes  $\leftarrow \text{known\_classes} \cup \text{new\_classes}$ ;
15   $\theta \leftarrow \theta'$ ;
16 end

```

C. Compute the k -Neighborhood Time Difference Distribution

Real-world graphs grow and change at a different pace [42]. Some graphs change quickly within a few time steps like social networks while others evolve rather slowly such as citation networks. Furthermore, graphs show different change behaviour, i. e., different patterns in how vertices and edges are added and removed over time. Therefore, depending on the specific graph data, a different history of the data must be used for training to take these factors into account.

We propose a new measure of k -Neighborhood Time Difference Distribution tdiff_k , which enumerates the distribution

of time differences within the k -hop neighborhood of each vertex (corresponding to the *receptive field* [43] of a GNN with k graph convolutional layers). This measure ensures that the history sizes are comparable with respect to these factors.

Definition 2 (k -Neighborhood Time Difference Distribution). Given graph \mathcal{G} and let $\mathcal{N}^k(u)$ be the k -hop neighborhood of u , i. e., the set of vertices that are reachable from u by traversing at most k edges. Let $\text{time} : \mathcal{V} \rightarrow \mathbb{N}$ be a function that gives the time information for each vertex, e. g., the year of a publication. We define $\text{tdiff}_k(\mathcal{G})$ to be the *multiset of time differences*, computed over all vertices $u \in V$ to their k -distant neighboring vertices $v \in \mathcal{N}^k(u)$ that occurred before u .

$$\text{tdiff}_k(\mathcal{G}) := \{ \text{time}(u) - \text{time}(v) \mid u \in V \wedge v \in \mathcal{N}^k(u) \wedge \text{time}(v) \leq \text{time}(u) \}$$

We interpret the multiset tdiff_k as a distribution over time differences, that is used to further analyze a dataset’s temporal distribution (percentiles) and to make datasets comparable. In our experiments, we compare models trained with a *limited history size* against models trained with the *full history*. We use the 25th, 50th, and 75th percentiles of this distribution as history sizes versus the 100th percentile to model the full graph to analyze the influence of explicit knowledge.

D. Models

We select representative GNN architectures as well as scalable GNN techniques for our experiments. The goal is to understand how different approaches of GNNs react to situation of changing graphs and new classes. Dwivedi et al. [44] distinguish between isotropic and anisotropic GNN architectures. In isotropic GNNs, all edges are treated equally. Apart from graph convolutional networks [2], examples of isotropic GNNs are GraphSAGE-mean [3], DiffPool [6], and GIN [45]. In anisotropic GNNs, the weights for edges are computed dynamically. Instances of anisotropic GNNs include graph attention networks [4], GatedGCN [46], and MoNet [47]. There are further approaches, which have been specifically proposed to scale GNNs to large graphs. These approaches fall into two categories: sampling [3], [10], [48], [49], and separating neighborhood aggregation from the neural network [16], [27], [50]. From each of these four categories (anisotropic versus isotropic GNNs and preprocessing versus sampling), we select one representative for our experiments.

We select **Graph Attention Networks** (GATs) [4] as representative for the class of anisotropic GNNs. In GATs, the representations in layer $l+1$ for vertex i are computed as follows: $\hat{\mathbf{h}}_i^{l+1} = \alpha_{ii}^l \mathbf{h}_i^l + \sum_{j \in \mathcal{N}(i)} \alpha_{ij}^l \mathbf{h}_j^l$ and $\mathbf{h}_i^{l+1} = \sigma(\mathbf{U}^l \hat{\mathbf{h}}_i^{l+1})$, where $\mathcal{N}(i)$ is the set of adjacent vertices to vertex i , \mathbf{U}^l are learnable parameters, and σ is a nonlinearity. The edge weights α_{ij} are computed by a self-attention mechanism based on h_i and h_j , i. e., the softmax of $a(\mathbf{U}^l \mathbf{h}_i \parallel \mathbf{U}^l \mathbf{h}_j)$ over edges, where a is an MLP and \parallel is the concatenation operation.

We select **GraphSAGE-Mean** [3] as a representative for isotropic GNNs because its special treatment of the vertices’

self-information has shown to be beneficial [44]. The representations of self-connections are concatenated with averaged neighbors’ representations before multiplying the parameters. In GraphSAGE-Mean, the procedure to obtain representations in layer $l+1$ for vertex i is given by the equations: $\hat{\mathbf{h}}_i^{l+1} = \mathbf{h}_i^l \parallel \frac{1}{\text{deg}_i} \sum_{j \in \mathcal{N}(i)} \mathbf{h}_j^l$ and $\mathbf{h}_i^{l+1} = \sigma(\mathbf{U}^l \hat{\mathbf{h}}_i^{l+1})$,

We select **Simplified GCN** [16] as a representative for shifting the neighborhood aggregation into preprocessing. Simplified GCN is a scalable variant of Graph Convolutional Networks (GCN) [2] that admits regular mini-batch sampling. Simplified GCN removes nonlinearities and collapses consecutive weight matrices into a single one. Thus, simplified GCN can be described by the equation $\hat{\mathbf{Y}}_{\text{SGC}} = \text{softmax}(\mathbf{S}^K \mathbf{X} \Theta)$, where \mathbf{S} is the normalized adjacency matrix and Θ is the weight matrix. The hyperparameter K has a similar effect as the number of layers in regular GCNs. Instead of using multiple layers, the k -hop neighbourhood is computed by \mathbf{S}^K , such that $\mathbf{S}^K \mathbf{X}$ can be precomputed. This makes Simplified GCN efficient, while not necessarily harming the performance [16].

We use **GraphSAINT** [10] as state-of-the-art subgraph sampling technique. In GraphSAINT, entire subgraphs are sampled for training GNNs. Subgraph sampling introduces a bias which is counteracted by normalization coefficients for the loss function. The authors propose different sampling methods: vertex sampling, edge sampling, and random-walk sampling. We use the best-performing random-walk sampling for our experiments. The underlying GNN is exchangeable, yet the authors suggest to use **Jumping Knowledge networks** (JKNets) [51]. JKNets introduce skip-connection to GNNs: each hidden layer has a direct connection to the output layer, in which the representations are aggregated, e. g., by concatenation. To isolate the effect of GraphSAINT sampling, we also include JKNets in our experiments.

IV. EXPERIMENTAL APPARATUS

A. Datasets

Pre-compiled graph datasets for lifelong learning are rare [52]. Many commonly used graph datasets are stripped off from any temporal data. We contribute three new graph datasets for lifelong learning on the basis of scientific publications: one new co-authorship graph dataset (PharmaBio) as well as two newly compiled citation graph datasets based on DBLP (DBLP-easy and DBLP-hard). For PharmaBio, the classes are journal categories. For DBLP, we use the conferences and journals of the published papers as classes. When new conferences and journals emerge, as they do in computer science, new classes will be introduced to the data. The datasets were generated by imposing a minimum threshold of publications per class per year: 100 for DBLP-easy and 45 for DBLP-hard, 20 for PharmaBio. For the co-authorship graph PharmaBio we additionally require a minimum of two publications per author per year. In all datasets, vertex features are normalized TF-IDF representations of the publication title.

Table I summarizes the basic characteristics of the datasets and Figure 3 shows the distribution of time differences tdiff_k for different values of k . We select 1, 3, 6, 25 as history

TABLE I

GLOBAL DATASET CHARACTERISTICS: TOTAL NUMBER OF VERTICES $|V|$, EDGES $|E|$, FEATURES D , CLASSES $|\mathbb{Y}|$ ALONG WITH # OF NEWLY APPEARING CLASSES (IN BRACES) WITHIN THE T EVALUATION TASKS

Dataset	$ V $	$ E $	D	$ \mathbb{Y} $	T
DBLP-easy	45,407	112,131	2,278	12 (4 new)	12
DBLP-hard	198,675	643,734	4,043	73 (23 new)	12
PharmaBio	68,068	2,1M	4,829	7	18

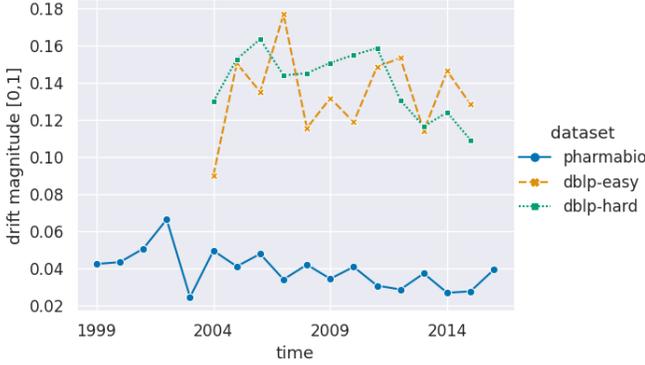


Fig. 2. Magnitude of the class drift per dataset. The drift within the PharmaBio dataset (no new classes) is lower than the drift of both DBLP variants. Independent and identically distributed data would have drift magnitude zero.

sizes for DBLP- $\{\text{easy,hard}\}$ and 1, 4, 8, 21 as history sizes for PharmaBio according to the 25,50,75,100-percentiles of tdiff_2 . For each dataset, we construct the sequence of tasks $\tilde{\mathcal{T}}_1, \dots, \tilde{\mathcal{T}}_T$ on the basis of the publication year along with a history size c . For each task $\tilde{\mathcal{T}}_t$, we construct a graph with respect to publications from time $[t-c, t]$, where publications from time t are the test vertices, and $t < c$ training vertices (transductive). In the inductive case that is used by GraphSAINT in our experiments, we train exclusively on $\tilde{\mathcal{T}}_{t-1}$, but still evaluate the test vertices of $\tilde{\mathcal{T}}_t$. We set the first evaluation task $\tilde{\mathcal{T}}_1$ to the time, in which 25% of the total number of publications are available. Thus, mapping the datasets to our problem statement (see Figure 1), our first evaluation task $t = 1$ corresponds to year 1999 in PharmaBio (total range: 1985–2016) and 2004 in DBLP- $\{\text{easy,hard}\}$ (1990–2015). We continue with the next years for subsequent tasks.

Regarding changes in the class set, DBLP-easy has 12 venues in total, including one bi-annual conference and four new venues appearing in 2005, 2006, 2007, and 2012. DBLP-hard has 73 venues, including one discontinued, nine bi-annual, six irregular venues, and 23 new venues. To quantify changes in the class set, we compute the magnitude of the class drift as total variation distance [53], [54]:

$$\sigma_{t-1,t} = \frac{1}{2} \sum_{y \in \mathbb{Y}_{t-1} \cup \mathbb{Y}_t} |P_{t-1}(y) - P_t(y)|$$

where $P_t(y)$ is the observed class probability at time t . The drift magnitudes are visualized per dataset in Figure 2.

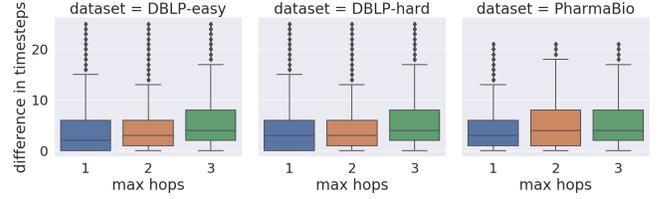


Fig. 3. Distributions of time differences tdiff_k (y-axis) for DBLP-easy (left), DBLP-hard (center) and PharmaBio (right) within the k -hop neighborhood for $k = \{1, 2, 3\}$ (x-axis).

B. Procedure

First, we *optimize the hyperparameters* by separately tuning the learning rate for each model, history size, and restart configuration on DBLP-easy (the effect of weight decay was negligible). All models are constrained to two graph convolutional layers, a comparable penultimate hidden dimension (2x32 GraphSAGE, 4x8 GAT, 2x2x16 JKNet, 64 MLP), and a 0.5 dropout rate. We fix an update step budget of 200 per task and use Adam [55] to optimize cross-entropy. We implemented GAT, GraphSAGE-mean, Simplified GCN, and JKNet with *dgl* [56] and use *torch-geometric* [57] for GraphSAINT. We had to disable GraphSAINT’s norm (re-)computation for each task such that our experiments could finish in reasonable time.

After hyperparameter optimization, we apply the models on all datasets using different history sizes. We run our incremental training method for graph learning from Section III-B for each of the models under both the warm restart and the cold restart configuration. Finally, we conduct an ablation study to isolate the effect of incremental training. All experiments are repeated 10 times with different random seeds.

C. Evaluation Measures

Our primary evaluation measure for our models f is accuracy. With $\text{acc}_t(f^{(t)})$, we denote the accuracy of model $f^{(t)}$ on task \mathcal{T}_t . We aggregate accuracy scores over the sequence of tasks $\mathcal{T}_1, \dots, \mathcal{T}_T$ by using their unweighted average [17]:

$$\text{acc}(f) = \frac{1}{T} \sum_{t \in \{1, \dots, T\}} \text{acc}_t(f^{(t)})$$

Following Lopez-Paz & Ranzato [17], we use Forward Transfer (FWT) to quantify the effect of reusing previous parameters. This is reflected by the accumulated differences in accuracy between the f_{warm} and f_{cold} models, defined below:

$$\text{FWT}(f_{\text{warm}}, f_{\text{cold}}) = \frac{1}{T-1} \sum_{t \in \{2, \dots, T\}} \text{acc}_t(f_{\text{warm}}^{(t)}) - \text{acc}_t(f_{\text{cold}}^{(t)})$$

V. EXPERIMENTAL RESULTS

Table II shows the aggregated results of 20,160 evaluation steps (48 configurations with 10 repetitions on two datasets with 12 tasks each and one dataset with 18 tasks). We consider method A better than method B when the mean accuracy of A is higher than the one of B and the 95% confidence intervals do not overlap [58]. In terms of absolute best methods per

setting (= dataset \times history size), we find that GraphSAGE consistently yields the highest scores except for DBLP-hard, where it is challenged by Simplified GCN.

Regarding the comparison of history sizes (i.e., *explicit knowledge*, see introduction), the highest scores are achieved in almost all cases by using an unlimited history size, i.e., using the full graph’s history. However, on all datasets, the scores for training with limited window sizes larger than 1 are close to the ones of full-graph training. With history sizes that cover 50% of the GNN’s receptive field, all methods achieve at least 95% relative accuracy compared to the same model under full-history training. When 75% of the receptive field is covered, the models yield at least 99% relative accuracy. To compute these percentages, we have selected the better one of either cold or warm restarts for each method.

Regarding the influence of *implicit knowledge*, we find that reusing parameters (warm restarts) leads to notably higher scores than re-training from scratch, when the history size comprises only one single previous task (see column Forward Transfer). The average Forward Transfer across all models and datasets with history size $c = 1$ is five accuracy points.

Regarding isotropic vs anisotropic GNNs, we find that GAT and GraphSAGE perform similarly well on DBLP-easy (on which the learning rate was tuned). However, GraphSAGE-mean yields higher scores on DBLP-hard and PharmaBio, which might indicate that GraphSAGE-mean is more robust to hyperparameters than GAT.

Regarding memory-efficient methods, we observe that the scores of Simplified GCN are among the highest of all methods on DBLP-hard. To understand this result, we recall that Simplified GCN uses only one single weight matrix of shape $n_{\text{features}} \times n_{\text{outputs}}$, which leads to 300k learnable parameters on DBLP-hard, but only 27k and 34k on DBLP-easy and PharmaBio, respectively. For comparison, GraphSAGE has 146k learnable parameters on DBLP-easy, 264k on DBLP-hard, and 310k on PharmaBio. On the other hand, GraphSAINT yields high scores on PharmaBio, comparable to GraphSAGE, but lower scores on both DBLP datasets.

VI. ABLATION STUDY

We isolate the effect of incremental training and compare once-trained (static) models against incrementally trained models. We train the static models for 400 epochs on the data before the first evaluation time step, which comprises 25% of the total vertices. We train incremental models for 200 epochs with history sizes of 3 time steps (4 on the PharmaBio dataset) before evaluating each task. We repeat each experiment 10 times with different random seeds.

In Figure 4, we see that the accuracy of the static models decreases over time on DBLP-easy and DBLP-hard, where new classes appear over time. On PharmaBio (fixed class set), the accuracy of the static models plateaus, while the accuracy of incrementally trained models increases.

VII. DISCUSSION

We have created a new experimental procedure for lifelong open-world node classification, for which we contribute three

newly compiled datasets. In this setup, we have evaluated five representative GNN architectures as well as an MLP baseline under an incremental training procedure. With the goal of generalizable results, we have introduced a new measure for k -neighborhood time differences tdiff_k , based on which we have selected the history sizes.

Our main result is that incremental training with limited history sizes is almost as good as using the full history of the graph. With window sizes of 3 or 4 (50% receptive field coverage), GNNs achieve at least 95% accuracy compared to using all past data for incremental training. With window sizes of 6 or 8 (75% receptive field coverage), at least 99% accuracy can be retained. This result holds for standard GNN architectures and also for scalable and sampling-based approaches. This has direct consequences for lifelong learning of GNNs on evolving graphs and, thus, impacts how GNNs can be employed for numerous real-world applications.

We have further investigated on reusing parameters from previous iterations (warm restarts). Using warm restarts is a viable strategy, i.e., reusing an “old” model, even though new classes appear during the sequence of tasks. We have identified a trend that reusing parameters from previous tasks becomes more important when the history sizes are small because less explicit knowledge is available. Even though it was not our main objective to compare the absolute performances of the models, it is noteworthy that Simplified GCNs perform surprisingly well on DBLP-hard. There, the model yields the highest absolute scores, on par with GraphSAGE-mean, despite the simplicity of the approach.

A limitation of the present work is that we assume that the labels, i.e., the ground truth of each task becomes available as training data for the next task. However, in practice only a small fraction of vertices might come with true labels, e.g., if the labels are professionally annotated subjects. Investigating whether only a small fraction of annotated nodes are sufficient for retraining would be an interesting direction of future work. Furthermore, our method does not yet include an unsupervised mechanism for self-detection of new classes [14], [24], [25]. However, our experimental apparatus lays the foundation for adding methods of self-detection of new classes.

To set our work in the broader context of lifelong learning, we reconsider the gradient episodic memory framework [17] for image data, in which examples are independent. To cast graph data into independent examples for vertex classification, certain preprocessing steps are required such as transforming each vertex into a graph [26]. This increases the number of inference steps by $\mathcal{O}(|V|)$ compared to our approach. Furthermore, we have shown that our approach of incremental training can be applied to various GNN architectures and is orthogonal to sampling and preprocessing approaches. In general, our incremental training procedure can be applied to any GNN architecture with few caveats: If the GNN architecture relies on transductive learning, the constraints also need to be satisfied during incremental training. Similarly, any pre-computation steps such as computing normalizing constants have to be performed again when adapting the model to a new task.

TABLE II
 ACCURACY (WITH 95% CONFIDENCE INTERVALS (CI) VIA 1.96 STANDARD ERROR OF THE MEAN) AND FORWARD TRANSFER (AVERAGED DIFFERENCE OF WARM AND COLD RESTARTS) ON OUR DATASETS WITH DIFFERENT HISTORY SIZES (COLUMN C). THE BEST METHOD PER CASE (= 1 DATASET + 1 HISTORY SIZE) IS MARKED IN BOLD, ALONG WITH METHODS WHERE THE 95% CI OVERLAPS.

Dataset	c	GAT			GraphSAGE-Mean			MLP (Baseline)		
		avg. accuracy cold	avg. accuracy warm	FWT	avg. accuracy cold	avg. accuracy warm	FWT	avg. accuracy cold	avg. accuracy warm	FWT
DBLP-easy	1	60.8 ± 0.5	64.9 ± 0.4	+4.5	60.4 ± 0.5	65.1 ± 0.4	+5.2	56.1 ± 0.4	62.2 ± 0.5	+6.6
	3	68.9 ± 0.3	69.3 ± 0.3	+0.2	68.7 ± 0.3	69.3 ± 0.3	+0.7	61.0 ± 0.5	62.9 ± 0.4	+2.0
	6	70.3 ± 0.4	70.2 ± 0.4	-0.1	71.1 ± 0.4	70.9 ± 0.4	-0.3	62.7 ± 0.3	62.7 ± 0.4	-0.2
	full	70.2 ± 0.4	70.2 ± 0.4	+0.1	71.6 ± 0.4	71.4 ± 0.3	-0.2	63.4 ± 0.3	61.9 ± 0.4	-1.2
DBLP-hard	1	39.4 ± 0.2	39.1 ± 0.2	-0.1	34.5 ± 0.4	40.0 ± 0.2	+5.9	31.6 ± 0.3	38.3 ± 0.3	+7.4
	3	44.0 ± 0.2	43.7 ± 0.2	-0.4	44.3 ± 0.2	45.1 ± 0.2	+0.8	33.7 ± 0.3	38.9 ± 0.2	+5.6
	6	45.1 ± 0.3	45.3 ± 0.3	+0.2	46.5 ± 0.3	46.7 ± 0.3	+0.2	39.2 ± 0.2	38.3 ± 0.2	-0.7
	full	45.6 ± 0.3	45.6 ± 0.3	-0.1	46.8 ± 0.2	47.1 ± 0.3	+0.4	38.2 ± 0.2	36.7 ± 0.2	-1.1
PharmaBio	1	61.6 ± 0.9	65.4 ± 0.9	+3.8	65.4 ± 0.9	68.6 ± 1.0	+3.3	62.7 ± 0.9	66.3 ± 0.9	+3.9
	4	64.5 ± 0.8	65.3 ± 0.9	+0.9	68.0 ± 0.8	69.0 ± 0.8	+1.1	66.3 ± 0.7	65.7 ± 0.8	-0.7
	8	65.1 ± 0.8	65.4 ± 0.8	+0.3	68.8 ± 0.7	69.0 ± 0.8	+0.2	64.2 ± 0.8	65.3 ± 0.7	+0.9
	full	64.3 ± 0.8	65.4 ± 0.8	+0.2	69.0 ± 0.7	68.4 ± 0.7	-0.7	65.4 ± 0.8	64.4 ± 0.6	-1.1

Dataset	c	Simplified GCN			GraphSAINT			Jumping Knowledge		
		avg. accuracy cold	avg. accuracy warm	FWT	avg. accuracy cold	avg. accuracy warm	FWT	avg. accuracy cold	avg. accuracy warm	FWT
DBLP-easy	1	57.1 ± 0.4	63.7 ± 0.4	+7.2	62.1 ± 0.3	63.2 ± 0.4	+1.2	56.2 ± 0.5	61.4 ± 0.5	+5.6
	3	66.4 ± 0.3	67.4 ± 0.3	+1.2	66.4 ± 0.4	65.3 ± 0.5	-0.9	65.2 ± 0.3	65.9 ± 0.5	+1.0
	6	69.3 ± 0.4	69.3 ± 0.4	+0.1	68.1 ± 0.4	65.5 ± 0.7	-2.1	68.0 ± 0.4	66.9 ± 0.6	-0.7
	full	71.0 ± 0.4	70.0 ± 0.4	-1.0	68.4 ± 0.5	65.7 ± 0.5	-2.8	68.7 ± 0.4	66.3 ± 0.4	-2.5
DBLP-hard	1	34.5 ± 0.3	41.0 ± 0.3	+7.0	35.9 ± 0.3	35.6 ± 0.4	+0.5	33.0 ± 0.2	35.3 ± 0.3	+2.9
	3	44.1 ± 0.2	44.8 ± 0.3	+0.8	39.3 ± 0.3	38.1 ± 0.5	-0.6	39.1 ± 0.3	38.8 ± 0.4	+0.3
	6	46.9 ± 0.3	46.2 ± 0.3	-0.4	40.6 ± 0.3	38.8 ± 0.6	-1.2	41.0 ± 0.3	40.1 ± 0.5	-0.3
	full	48.8 ± 0.4	47.5 ± 0.3	-1.2	41.0 ± 0.4	40.7 ± 0.4	-0.3	41.6 ± 0.3	40.8 ± 0.2	-0.9
PharmaBio	1	62.3 ± 0.9	64.5 ± 0.8	+2.3	65.7 ± 0.8	68.6 ± 0.8	+3.0	64.1 ± 0.9	68.3 ± 0.9	+4.3
	4	64.4 ± 0.8	64.4 ± 0.8	-0.0	67.3 ± 0.8	68.4 ± 0.7	+1.0	67.1 ± 0.8	68.2 ± 0.8	+1.1
	8	65.3 ± 0.8	64.0 ± 0.7	-1.4	68.1 ± 0.8	68.0 ± 0.7	-0.1	67.8 ± 0.8	67.7 ± 0.7	-0.3
	full	62.4 ± 0.8	61.7 ± 0.6	-0.8	68.2 ± 0.8	66.1 ± 0.8	-2.2	66.8 ± 0.8	64.5 ± 0.7	-2.6

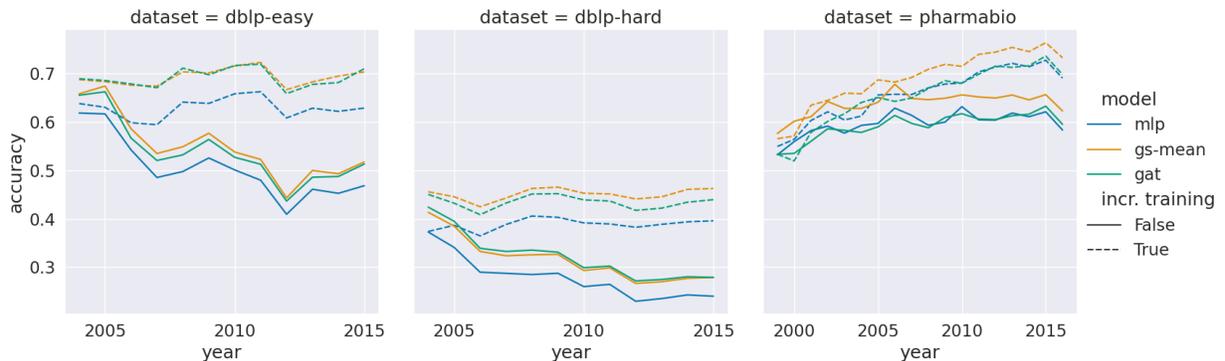


Fig. 4. Results of Ablation Study: Accuracy scores of once-trained, static models (solid lines) are lower than incrementally trained models (dashed lines).

VIII. CONCLUSION

We have introduced an incremental training scheme for lifelong learning on evolving graphs where new classes appear over time. We have proposed a method to select the history sizes for training the models agnostic to different granularities of temporal intervals. Our analysis of implicit and explicit knowledge in lifelong learning on graphs shows that high

levels of accuracy can be preserved, even when training only on a fraction of past data. Furthermore, using warm restarts becomes more important when few past data is available.

Data Availability and Reproducibility

We published our lifelong graph learning datasets (zenodo.org/record/3764770#.YCQUOHWYXmg) along with an implementation of our experimental procedure (github.com/lgalke/lifelong-learning).

Acknowledgments: Parts of this research were carried out on the computing infrastructure *bwCloud* of the State of Baden-Württemberg and the Bioinformatics and Systems Biology group at Ulm University.

REFERENCES

- [1] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini, "The graph neural network model," *IEEE Trans. Neural Networks*, vol. 20, no. 1, 2009.
- [2] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.
- [3] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NeurIPS*, 2017.
- [4] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *ICLR*, 2018.
- [5] J. Klicpera, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank," in *ICLR*, 2019.
- [6] Z. Ying, J. You, C. Morris, X. Ren, W. L. Hamilton, and J. Leskovec, "Hierarchical graph representation learning with differentiable pooling," in *NeurIPS*, 2018.
- [7] M. Zhang and Y. Chen, "Link prediction based on graph neural networks," in *NeurIPS*, 2018.
- [8] P. Veličković, W. Fedus, W. L. Hamilton, P. Liò, Y. Bengio, and R. D. Hjelm, "Deep graph infomax," in *ICLR*, 2019.
- [9] X. Da, R. Chuanwei, K. Evren, K. Sushant, and A. Kannan, "Inductive representation learning on temporal graphs," in *ICLR*, 2020.
- [10] H. Zeng, H. Zhou, A. Srivastava, R. Kannan, and V. K. Prasanna, "Graphsaint: Graph sampling based inductive learning method," in *ICLR*, 2020.
- [11] Z. Yang, W. W. Cohen, and R. Salakhutdinov, "Revisiting semi-supervised learning with graph embeddings," in *ICML*, 2016.
- [12] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: online learning of social representations," in *KDD*. ACM, 2014.
- [13] S. Thrun and T. M. Mitchell, "Learning one more thing," in *IJCAI*. Morgan Kaufmann, 1995.
- [14] G. Fei, S. Wang, and B. Liu, "Learning cumulatively to become more knowledgeable," in *KDD*. ACM, 2016.
- [15] Z. Chen and B. Liu, *Lifelong Machine Learning, Second Edition*, ser. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2018.
- [16] F. Wu, A. H. S. Jr., T. Zhang, C. Fifty, T. Yu, and K. Q. Weinberger, "Simplifying graph convolutional networks," vol. 97, 2019.
- [17] D. Lopez-Paz and M. Ranzato, "Gradient episodic memory for continual learning," in *NIPS*, 2017.
- [18] S. Thrun, "Lifelong learning algorithms," in *Learning to learn*. Springer, 1998.
- [19] D. L. Silver, Q. Yang, and L. Li, "Lifelong machine learning systems: Beyond learning algorithms," in *AAAI Spring Symposium: Lifelong Machine Learning*, vol. SS-13-05. AAAI, 2013.
- [20] B. Liu, "Lifelong machine learning: a paradigm for continuous learning," *Frontiers of Computer Science*, vol. 11, no. 3, 2017.
- [21] M. Herbster, M. Pontil, and L. Wainer, "Online learning over graphs," in *ICML*, vol. 119. ACM, 2005, pp. 305–312.
- [22] P. Ruvolo and E. Eaton, "ELLA: an efficient lifelong learning algorithm," in *ICML*, vol. 28. JMLR.org, 2013.
- [23] A. V. Robins, "Catastrophic forgetting, rehearsal and pseudorehearsal," *Connect. Sci.*, vol. 7, no. 2, 1995.
- [24] M. M. Masud, J. Gao, L. Khan, J. Han, and B. M. Thuraisingham, "Classification and novel class detection in concept-drifting data streams under time constraints," *IEEE Trans. Knowl. Data Eng.*, 2011.
- [25] M. Wu, S. Pan, and X. Zhu, "Openwgl: Open-world graph learning," in *ICDM*. IEEE, 2020.
- [26] C. Wang, Y. Qiu, and S. A. Scherer, "Lifelong graph learning," *CoRR*, vol. abs/2009.00647, 2020.
- [27] E. Rossi, F. Frasca, B. Chamberlain, D. Eynard, M. Bronstein, and F. Monti, "Sign: Scalable inception graph neural networks," 2020.
- [28] G. H. Nguyen, J. B. Lee, R. A. Rossi, N. K. Ahmed, E. Koh, and S. Kim, "Continuous-time dynamic network embeddings," in *WWW 2018*.
- [29] J. B. Lee, G. Nguyen, R. A. Rossi, N. K. Ahmed, E. Koh, and S. Kim, "Dynamic node embeddings from edge streams," *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2020.
- [30] P. Goyal, N. Kamra, X. He, and Y. Liu, "Dyngem: Deep embedding method for dynamic graphs," *CoRR*, vol. abs/1805.11273, 2018.
- [31] P. Goyal, S. R. Chhetri, and A. Canedo, "dyngraph2vec: Capturing network dynamics using dynamic graph representation learning," *Knowl.-Based Syst.*, vol. 187, 2020.
- [32] R. Trivedi, H. Dai, Y. Wang, and L. Song, "Know-evolve: Deep temporal reasoning for dynamic knowledge graphs," in *ICML*, 2017.
- [33] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson, "Structured sequence modeling with graph convolutional recurrent networks," in *ICONIP*. Springer, 2018.
- [34] S. Kumar, X. Zhang, and J. Leskovec, "Learning dynamic embeddings from temporal interactions," *arXiv preprint arXiv:1812.02289*, 2018.
- [35] R. Trivedi, M. Farajtabar, P. Biswal, and H. Zha, "Dyrep: Learning representations over dynamic graphs," in *ICLR*, 2019.
- [36] F. Manessi, A. Rozza, and M. Manzo, "Dynamic graph convolutional networks," *Pattern Recognition*, vol. 97, 2020.
- [37] A. Sankar, Y. Wu, L. Gou, W. Zhang, and H. Yang, "Dysat: Deep neural representation learning on dynamic graphs via self-attention networks," in *WSDM*, 2020.
- [38] E. Rossi, B. Chamberlain, F. Frasca, D. Eynard, F. Monti, and M. Bronstein, "Temporal graph networks for deep learning on dynamic graphs," *arXiv preprint arXiv:2006.10637*, 2020.
- [39] A. Pareja, G. Domeniconi, J. Chen, T. Ma, T. Suzumura, H. Kanezashi, T. Kaler, T. B. Schardl, and C. E. Leiserson, "Evolvegcn: Evolving graph convolutional networks for dynamic graphs," in *AAAI*, 2020.
- [40] B. Fish and R. S. Caceres, "A task-driven approach to time scale detection in dynamic networks," in *Workshop on Mining and Learning with Graphs*, 2017.
- [41] D. Ersan, C. Nishioka, and A. Scherp, "Comparison of machine learning methods for financial time series forecasting at the examples of over 10 years of daily and hourly data of DAX 30 and S&P 500," *J Comput Soc Sc*, vol. 3, 2020.
- [42] C. Aggarwal and K. Subbian, "Evolutionary network analysis: A survey," *ACM Comput. Surv.*, vol. 47, no. 1, May 2014.
- [43] J. Chen, J. Zhu, and L. Song, "Stochastic training of graph convolutional networks with variance reduction," in *ICML*, 2018.
- [44] V. P. Dwivedi, C. K. Joshi, T. Laurent, Y. Bengio, and X. Bresson, "Benchmarking graph neural networks," *CoRR*, 2020, abs/2003.00982.
- [45] K. Xu, W. Hu, J. Leskovec, and S. Jegelka, "How powerful are graph neural networks?" in *ICLR*, 2019.
- [46] X. Bresson and T. Laurent, "Residual gated graph convnets," *CoRR*, vol. abs/1711.07553, 2017.
- [47] F. Monti, D. Boscaini, J. Masci, E. Rodolà, J. Svoboda, and M. M. Bronstein, "Geometric deep learning on graphs and manifolds using mixture model cnns," in *CVPR*. IEEE Computer Society, 2017.
- [48] W. Huang, T. Zhang, Y. Rong, and J. Huang, "Adaptive sampling towards fast graph representation learning," in *NeurIPS*, 2018.
- [49] W. Chiang, X. Liu, S. Si, Y. Li, S. Bengio, and C. Hsieh, "Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks," in *KDD*. ACM, 2019.
- [50] A. Bojchevski, J. Klicpera, B. Perozzi, A. Kapoor, M. Blais, B. Rózemberczki, M. Lukasik, and S. Günnemann, "Scaling graph neural networks with approximate pagerank," in *KDD*. ACM, 2020.
- [51] K. Xu, C. Li, Y. Tian, T. Sonobe, K. Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks," in *ICML*, 2018.
- [52] 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," in *NeurIPS*, 2020.
- [53] G. I. Webb, R. Hyde, H. Cao, H. L. Nguyen, and F. Petitjean, "Characterizing concept drift," *Data Mining and Knowledge Discovery*, vol. 30, no. 4, 2016.
- [54] G. I. Webb, L. K. Lee, B. Goethals, and F. Petitjean, "Analyzing concept drift and shift from sample data," *Data Mining and Knowledge Discovery*, vol. 32, no. 5, 2018.
- [55] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *ICLR (Poster)*, 2015.
- [56] M. Wang, *et al.*, "Deep graph library: Towards efficient and scalable deep learning on graphs," *arXiv preprint arXiv:1909.01315*, 2019.
- [57] M. Fey and J. E. Lenssen, "Fast graph representation learning with PyTorch Geometric," in *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [58] I. J. Goodfellow, Y. Bengio, and A. C. Courville, *Deep Learning*, ser. Adaptive computation and machine learning. MIT Press, 2016.