

# ONLINE CONTINUAL GRAPH LEARNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

The aim of Continual Learning (CL) is to learn new tasks incrementally while avoiding catastrophic forgetting. Online Continual Learning (OCL) specifically focuses on learning efficiently from a continuous stream of data with shifting distribution. While recent studies explore Continual Learning on graphs exploiting Graph Neural Networks (GNNs), only few of them focus on a streaming setting. Many real-world graphs evolve over time and timely (online) predictions could be required. However, current approaches are not well aligned with the standard OCL literature, partly due to the lack of a clear definition of online Continual Learning on graphs. In this work, we propose a general formulation for online Continual Learning on graphs, emphasizing the efficiency of batch processing while accounting for graph topology, providing a grounded setting to analyze different methods. We present a set of benchmark datasets for online continual graph learning, together with the results of several methods in CL literature, adapted to our setting. Additionally, we address the challenge of GNN memory usage, as considering multiple hops of neighborhood aggregation can require access to the entire growing graph, resulting in prohibitive costs for the setting. We thus propose solutions to maintain bounded complexity for efficient online learning.

## 1 INTRODUCTION

In traditional machine learning, models are at first trained in an environment where all training data is simultaneously available to the learning algorithm, and only at a later time model predictions are produced on new input data. Importantly, data observations are assumed to be mutually independent and identically distributed. Real-world environments however often generate data in chunks or streams which often entail shifts in the data distribution or even variations in tasks to be solved. In turn, previously trained models may require expensive retraining or model reconfiguration to stay up to date. In this setting, Continual Learning (CL) (Parisi et al., 2019; De Lange et al., 2022), lifelong learning (Chen & Liu, 2018) and incremental learning (Chaudhry et al., 2018a) are similar machine learning paradigms sharing the same goal of adapting models to incrementally learn as soon as new data and new tasks are presented.

Even further, in the *online* learning setting, training data points are collected sequentially, processed in real-time by the learning method, and immediately discarded (Chaudhry et al., 2018b; Mai et al., 2022). Such strict environments are found in monitoring and control problems (Zliobaite et al., 2016; Gunasekara et al., 2023), such as traffic management, activity recognition, fraud detection, with other applications such as on data generated by optical sensors (Souza et al., 2020) or prediction of power production considering environmental conditions (Lobo et al., 2020). Online CL is therefore an extremely challenging setting that requests models to quickly adapt to new conditions, allowing for anytime inference with minimal latency and without forgetting already acquired knowledge.

Research on CL extends to graph-structured data as well, often referred to as Continual Graph Learning (CGL) (Yuan et al., 2023). Indeed, many machine learning tasks naturally represent data as graphs and address predictions of properties of single nodes or edges, or entire graphs. Common examples include social networks, citation networks, and biological systems, where relationships between entities are modeled as graphs. However, most graphs in the real world are not static: they continuously evolve over time experiencing the addition/removal of nodes and modifications to the topology. Examples are seen in dynamic environments such as social networks, citation networks, and traffic systems, where new users join, papers are published, and road conditions may change (Liu et al., 2021; Zhou & Cao, 2021). Countless deep learning models for graph-

054 structured data, such as Graph Neural Networks (GNNs), make predictions by relying on node-  
 055 level representation computed from neighboring nodes. While this enables rich representations that  
 056 condition the predictions on related observations, it poses unique challenges associated with the  
 057 online setting, due to the need for neighboring, past information to make predictions for new nodes,  
 058 as well as memory and computational issues due to the rapid growth of certain graph families.

059 In this paper, we introduce three main contributions to address these challenges. (1) We formalize  
 060 the Online Continual Graph Learning (OCGL) framework, establishing a foundation for Contin-  
 061 ual Learning on graphs in a node-level streaming environment. By doing so, we bridge the gap  
 062 between Online CL and CGL settings. (2) We present a benchmarking environment for OCGL.  
 063 We introduce four benchmark datasets and conduct a first set of extended experiments to evaluate  
 064 various CL methods, that we suitably adapted to operate under the OCGL setting. We use a hyper-  
 065 parameter selection protocol tailored for online learning that ensures a fair comparison between CL  
 066 techniques. Insightful conclusions are drawn from the analyses, highlighting the higher performance  
 067 of replay-based methods, especially A-GEM. (3) We bring the attention of the research community  
 068 to the issue of neighborhood expansion which could undermine the online computational efficiency  
 069 due to the complexity of multi-hop aggregation in GNNs. To address it, we propose a first, simple  
 070 solution to ensure that models can efficiently scale as the graph evolves. We conduct a second set of  
 071 experiments with the proposed neighborhood sampling solution. This work leaves room for further  
 072 research to develop more effective approaches to tackle the problem.

## 074 2 BACKGROUND AND RELATED WORKS

076 **Continual Learning.** Depending on the type of shift in the data distribution, CL has been catego-  
 077 rized into three main scenarios (Van De Ven et al., 2022): in *task-incremental* learning, the model  
 078 sequentially learns distinct tasks, which requires availability of task identifiers to make predictions;  
 079 *class-incremental* learning consists in classifying instances with an increasing number of classes,  
 080 without task identifiers; finally, *domain-incremental* learning requires solving the same problem in  
 081 different contexts. In the past, the main applications of CL included reinforcement learning (Kirk-  
 082 patrick et al., 2017; Rolnick et al., 2019) and especially computer vision (Rebuffi et al., 2017; Lopez-  
 083 Paz & Ranzato, 2017; Aljundi et al., 2018; Li & Hoiem, 2018; Masana et al., 2022; Mai et al., 2022),  
 084 even though most of the methods that have been developed to address these problem domains are  
 085 not limited to these fields and can be used for a wide range of other machine learning tasks. Three  
 086 main broad categories of CL approaches to mitigate forgetting have been proposed in the literature  
 087 (De Lange et al., 2022): *regularization* methods, *replay* methods and *architectural* methods. Regu-  
 088 larization methods (Kirkpatrick et al., 2017; Zenke et al., 2017; Aljundi et al., 2018; Li & Hoiem,  
 089 2018; Chaudhry et al., 2018a) introduce additional loss terms to preserve important parameters to  
 090 retain previously acquired knowledge. Replay methods (Rebuffi et al., 2017; Lopez-Paz & Ranzato,  
 091 2017; Chaudhry et al., 2019; 2018b) use a memory buffer to store some representative samples from  
 092 old tasks, to use them jointly with new samples to update model parameters. Architectural methods  
 093 (Fernando et al., 2017; Masse et al., 2018; Rusu et al., 2022) avoid changes to model parameters  
 094 either by gating mechanisms or by introducing new parameters, allowing the model to grow in size.

095 **Online Continual Learning.** In the usual Continual Learning scenarios described above, data arrive  
 096 one task at a time, allowing for *offline* training with multiple passes over the data for the current  
 097 task (De Lange et al., 2022). A more realistic yet challenging scenario is that of *Online Continual*  
 098 *Learning (OCL)* (Chaudhry et al., 2018b; Mai et al., 2022; Soutif-Cormerais et al., 2023), where  
 099 data arrive in small batches of only few samples, without the possibility for the model to store all the  
 100 data for the current task, either for privacy reasons or memory limitations. In this setting the learning  
 101 algorithm must be able to use efficiently the mini-batches that arrive in the non-stationary stream.  
 102 Additionally, whereas for CL we assume to know the task boundaries, OCL can be performed in  
 103 a boundary-agnostic setting, or task-free, potentially allowing for more varied distribution shifts  
 104 compared to the three settings described before (Koh et al., 2021). Many CL methods though are  
 105 not suited to this setting and require modifications. An additional characteristic of OCL is the ability  
 106 to perform anytime inference: the model should always be up-to-date and ready to make predictions  
 107 online after each training batch, reacting quickly to distribution shifts (Koh et al., 2021).

**Learning on graphs.** *Graph Neural Networks (GNN)* (Sperduti & Starita, 1997; Scarselli et al.,  
 2009; Micheli, 2009; Kipf & Welling, 2017) have emerged as the state-of-the-art approach for

dealing with network data, generalizing convolution to graph structures. The core mechanism of most GNNs is message passing (Gilmer et al., 2017): at each layer, the hidden embedding  $h_v^{(k)}$  of each node  $v$  is updated using information from its neighborhood  $\mathcal{N}(v)$  as  $h_v^{(l+1)} = \text{UPDATE}(h_v^{(l)}, \text{AGGREGATE}(\{h_u^{(l)} : u \in \mathcal{N}(v)\}))$ . Here AGGREGATE and UPDATE are differentiable functions specified by the particular model. Specifically, as at each step each node updates its embedding using the information (message) coming from its neighbors, after  $l$  layers it will depend on its  $l$ -hop neighborhood. Graph-based processing of temporal data has a relatively short history, primarily encompassing the study of temporal graphs (Kazemi et al., 2020; Gravina & Bacciu, 2024) and time series data (Cini et al., 2023a; Jin et al., 2024) with dedicated adaptation strategies to deal with evolving graphs (Cini et al., 2023b) and benchmarks (Huang et al., 2023).

**Continual Graph Learning.** In the past few years researcher started to develop CL strategies tailored to graph data (Wang et al., 2020), with applications such as recommender systems (Xu et al., 2020) and traffic prediction (Chen et al., 2021). The main approaches resemble the general CL strategies, with a particular focus on preserving topological information with a loss term on neighborhood aggregation parameters (Liu et al., 2021), or specific node selection policies to retain informative nodes as memory buffer (Zhou & Cao, 2021). Recently a number of surveys have been published on the topic (Febrinanto et al., 2023; Yuan et al., 2023; Zhang et al., 2024; Tian et al., 2024), and a benchmark has been proposed (Zhang et al., 2022). *Continual Graph Learning (CGL)* possesses some peculiarities that differentiate it from other problem domains with independent samples, as graph structure requires careful consideration. Specifically, we can distinguish between *graph-level CGL* and *node-level CGL* (Zhang et al., 2022). In the first, the model needs to make predictions about entire graphs, and thus each sample in the training data is an independent graph. Standard Continual Learning settings and methods can thus be applied with relative ease (Carta et al., 2022). Instead, in node-level CGL predictions are made about the single nodes in usually one single graph. In this node-level CGL setting each new task therefore consists a new subgraph of the overall graph, with new classes or type of nodes. Specifically, the task subgraph arrives all at once, and offline training is performed on it. Here, an additional specification needs to be made about the availability of inter-task edges (Tian et al., 2024): since message passing GNNs aggregate information from neighboring nodes, the presence of inter-task edges towards nodes of previous tasks would inevitably lead to the use of past data. In practice, inter-task edges are often kept, but without access to the labels from past tasks (Zhang et al., 2024). Finally, the addition of new edges creates an additional source of backward interference: when evaluating the model on nodes belonging to past tasks, their neighborhood composition and topology will be changed compared to when the task was observed.

The Online Continual Learning setting described before has been explored in domains like computer vision (Mai et al., 2022; Soutif-Cormerais et al., 2023) and on sequences (Parisi & Lomonaco, 2020), but to the best of our knowledge it has not yet been applied to graphs. Some papers on CGL consider a setting referred to as *streaming* (Wang et al., 2020; Perini et al., 2022), yet the approaches can be categorized as offline CL as the streams consist of graph snapshots (often corresponding to entire tasks or nonetheless not classifiable as mini-batches), on which models are trained with multiple passes.

### 3 ONLINE CONTINUAL GRAPH LEARNING

This section introduces *Online Continual Graph Learning (OCGL)*, a framework that ports CGL to the online problem setting. OCGL is particularly applicable to dynamic real-world scenarios such as social networks or recommender systems, where sudden distribution changes occur, and quick model adjustments are essential for making anytime predictions.

#### 3.1 PROBLEM FORMULATION

**A growing network.** We model the data associated with an OCGL problem as an evolving graph  $\mathcal{G}$  induced by a stream of nodes  $v_1, v_2, \dots, v_t, \dots \in \mathbb{N}$ . At every time-step  $t$ , the monitored system is represented as a graph  $\mathcal{G}^t = (\mathbb{V}^t, \mathbb{E}^t, \mathbf{X}^t)$  defined by node set  $\mathbb{V}^t = \{v_1, \dots, v_t\}$ , edge set  $\mathbb{E}^t \subseteq \mathbb{V}^t \times \mathbb{V}^t$ , and a set of node attributes  $\mathbf{X}^t = \{\mathbf{x}^i\}_{i \leq t} \subset \mathbb{R}^F$ . Edge attributes, e.g., accounting for edge directions or defining the type of node-node relations, can be considered likewise, however, they are here excluded to ease the presentation. The graph nodes  $v_i$  can be associated with

class labels  $y_i \in \{1, \dots, C\}$  collected in set  $Y^t = \{y_i\}_{i \leq t}$  to be predicted and/or used as training samples to learn the model. Graph  $\mathcal{G}^t$  is a snapshot of the temporal graph  $\mathcal{G}$  which, together with  $Y^t$ , collects all information available up to time step  $t$ . At every time step  $t$  a new node is acquired from the monitored system and populates the current snapshot graph  $\mathcal{G}^{t-1}$ . Specifically, tuple  $(v_t, \mathcal{N}(v_t), \mathbf{x}^t)$  containing a new node index  $v_t \notin \mathbb{V}^{t-1}$ , associated node features  $\mathbf{x}_t$ , and a set of neighbors  $\mathcal{N}(v_t) \subseteq \mathbb{V}^{t-1}$  is presented and connected to graph  $\mathcal{G}^{t-1}$  according to the relations contained in  $\mathcal{N}(v_t)$ . Finally, target class label  $y_t$  of node  $v_t$  may or may not be acquired contextually to  $(v_t, \mathcal{N}(v_t), \mathbf{x}^t)$ ; for instance, a prediction for node  $v_t$  might be requested at time  $t$  while the true class label  $y_t$  is observed only at a later time.

**Problem statement.** The goal of OCGL is to learn a model  $f_\theta$  to predict class label  $y_t$  only from the subgraph of  $\mathcal{G}_t$  associated with  $v_t$  and its neighbors of order 1 or more. Parameters  $\theta$  have a predefined (maximum) dimension and model  $f_\theta$  is trained incrementally as soon as new nodes and the associated labels are provided. New nodes are acquired either individually or in small minibatches, slightly weakening the online setting as commonly done in the literature (Chaudhry et al., 2018b). Moreover, minibatches are seen only once and, after prediction and/or training is performed the mini-batch is discarded. As per the CGL problem setting, we allow the task to change over time. This requirement on small mini-batches is set to meet the need to perform anytime predictions Koh et al. (2021) where dealing with the entire graph  $\mathcal{G}^t$  – or a large subgraph of it – is unfeasible either memory-wise or computationally. Specifically, we require the training on each mini-batch to have bounded compute and memory budgets. Although the size of the node mini-batch may vary between applications, we assume it to be small enough that using only edges within the mini-batch would not provide sufficient context for effective learning. In contrast to some CGL settings where training can be done on a task-specific subgraph, in OCGL the mini-batch is too small to be treated as a meaningful graph on its own, requiring to access neighborhood information.

**Mini-batching.** Within the introduced OCGL framework, and depending on the application, minibatches can include  $L$ -hop neighbors with  $L \geq 1$ . Considering  $L > 1$  does not conflict with the growing nature of the graph. Instead, it reflects the requirements of the predictive model  $f_\theta$ , which may rely on aggregating multi-hop information. As such, to construct  $L$ -hop neighborhoods, an up-to-date snapshot  $\mathcal{G}^t$  is assumed to be stored in a *Past Information Store (PIS)* system, separate from an eventual memory buffer associated with predictive model  $f_\theta$ , as in the more general lifelong-learning setup (Chen & Liu, 2018). We do not impose memory limitations on the PIS to allow graph growth, but we require the training on each mini-batch to have bounded computational time and memory cost. That is, we assume to make limited use of information from the PIS for each mini-batch, and to only have access to the labels of nodes in the current batch. The definition of the evolving graph is general as it does not make assumptions on the distribution shifts happening in the node stream. This general, task free setting can be easily adapted or made more specific depending on the node stream: while a real-world stream could be induced by a time-stamp on the nodes, this setting can be derived from any static graph by establishing an ordering on the nodes. The three CL scenarios of task-, class- and domain-incremental can thus be easily adapted to an online setting by ordering nodes by task, similarly to what is done in other domains (Mai et al., 2022; Soutif–Cormerais et al., 2023).

### 3.2 NEIGHBORHOOD EXPANSION PROBLEM

The efficiency requirement for online learning poses non-negligible issues associated with reiterated message passing within multi-layer GNNs. As commented above, at each layer of the  $L$  layer, the GNN aggregates the embeddings of the neighboring nodes, thus requiring access to  $L$ -hop neighborhoods. The size of the  $L$ -hop neighborhood however scales as  $O(d^L)$  where  $d$  is the average degree. Moreover,  $d$  is time-dependent and can increase as the graph grows. Depending on how well-connected a graph is, very few hops may be required to go from any node to almost any other one; empirical evidence on the neighborhood growth is depicted in Figure 2 and those in Appendix E. Thus, having high  $d$  or  $L$  would require processing a number of nodes in the order of the entire (growing) graph for each mini-batch, going against efficiency in our definition of OCGL. It is therefore of paramount importance to keep a low  $L$  or to introduce a strategy to mitigate  $d$ .

In cases where the topology of the full graph or the maximum node degree are known a-priori (such as in citation networks, where we expect the number of references for an article not to explode even if the body of literature increases), fixing a low number of layers can be a good solution. This is

the first strategy we evaluate in Section 6. Yet, in many real applications we may not have prior knowledge on the evolution of average degree. In such cases, even choosing a model with a low number of layers may be problematic, as the average degree could grow indefinitely and with it the computational complexity and memory usage for batch training. This issue is similar to the problem of scaling static GNNs for large graphs, for which mini-batch training is required both for memory and efficiency reasons. Numerous approaches have been developed to address this problem, such as fixing a number of neighborhood sampled for aggregation (Hamilton et al., 2017; Chen et al., 2018) or training on partitions of the graph (Chiang et al., 2019). In our context, the simplest solution seems to fix the number of neighbors for aggregation through sampling, which guarantees an upper bound on the size of the computational graph for each batch. Results with neighborhood sampling are reported in Section 7, after an empirical assessment of the problem of neighborhood expansion.

## 4 METHODS

Having defined the Online Continual Graph Learning setting, we consider and evaluate some popular CL techniques, most of which are agnostic with respect to the type of the input data. Some CGL learning strategies are not applicable to the online setting, such as ER-GNN (Zhou & Cao, 2021), which stores representative nodes according to metrics computed offline on an entire graph snapshot, and thus we resorted to a simplified version which performs reservoir sampling. Additionally, baselines that require expensive fine-tuning steps such as GDumb (Prabhu et al., 2020) are excluded, as it would violate the online setting. Overall, most strategies natively require task boundaries, and have been modified for the task free setting as described below.

1. **ER**. Experience replay (Chaudhry et al., 2019) is a simple yet powerful replay-based method, which selects samples to be stored in a memory buffer by reservoir sampling (Vitter, 1985). New incoming batches for training are then augmented with nodes sampled uniformly from the buffer.
2. **EWC**. Elastic Weight Consolidation (Kirkpatrick et al., 2017) adds a quadratic term to the loss to penalize the modification of important parameters. Parameter importance is approximated by the diagonal of the Fisher information matrix, which needs to be computed offline for each task. We therefore modify the algorithm to keep one single Fisher information matrix updated with a running average over the batches, similarly to the MAS approach detailed later. Another approach would be to keep a moving average, as done in EWC++ (Chaudhry et al., 2018a).
3. **A-GEM**. Averaged GEM (Chaudhry et al., 2018b) is a more efficient version of GEM (Lopez-Paz & Ranzato, 2017), which ensures that the average loss for past tasks does not increase. It achieves this by projecting the gradient of the incoming batch in the orthogonal space of the gradient computed on samples from a memory buffer, if their scalar product is negative. We select nodes for the buffer with reservoir sampling (Vitter, 1985).
4. **LwF**. Learning without Forgetting (Li & Hoiem, 2018) uses distillation (Hinton et al., 2015) to regularize the loss with logits from a previous version of the model (teacher) on the current batch. To use it in a task free setting, we introduce an additional hyperparameter: the number of batches after which the teacher is updated with the current model.
5. **MAS**. Memory Aware Synapses (Aljundi et al., 2018) is a quadratic regularization similar to EWC, but it calculates importance as the sensitivity of the output on parameters. MAS is natively an online method, as the importance scores are accumulated with each new data point.
6. **TWP**. Topology-aware Weight Preserving (Liu et al., 2021) is another regularization method, which preserves important weights for topological aggregation in GAT (Veličković et al., 2018), generalized also to other GNNs. We modify it for the online setting as EWC.

## 5 EXPERIMENTAL SETUP

In this section we introduce the specific experimental setup used, describing the construction of the node streams from benchmark datasets, and the details of model training and evaluation<sup>1</sup>.

**Benchmarks.** Four node classification graph datasets are used in our experiments: CoraFull (Bojchevski & Günnemann, 2018), Arxiv (Hu et al., 2021), Reddit (Hamilton et al., 2017) and Amazon Computer (Shchur et al., 2019). The datasets are described in Appendix A. In order to position our experiments close to the rest of the Continual Learning literature, we devise a node stream derived

<sup>1</sup>Code available at: (supplementary material, link to be added after double-blind review)

270 from the class-incremental CL setting, which is considered as the most challenging one for catas-  
 271 trophic forgetting (Masana et al., 2022). We divide the nodes in the graph into groups with fixed  
 272 order consisting of 2 classes: this would be the sequence of tasks in class-incremental learning (re-  
 273 sulting in 35 tasks for CoraFull, 20 tasks for Arxiv and Reddit, and 5 for Amazon Computer), with  
 274 task boundaries between pairs of classes. Then, we fix and ordering on the nodes of each task, and  
 275 we stream the nodes accordingly. Therefore, the graph will gradually grow with mini-batches of  
 276 nodes from two new classes at a time, which are processed in an online fashion. This allows us to  
 277 consider metrics from the CL literature which require task boundaries, even though in our experi-  
 278 ments the learning algorithm itself is task agnostic and simply adds a new output neuron when an  
 279 instance of a new class is observed. For each dataset, we split the graph into 60% for training, 20%  
 280 for validation and 20% for testing. A transductive setting is used: validation and test nodes are not  
 281 used for loss computation, but they are still used for message passing.

282 **Performance assessment.** We consider three metrics widely adopted in the literature: *Average Ac-*  
 283 *curacy (AA)*, *Average Forgetting (AF)* (Lopez-Paz & Ranzato, 2017), and *Average Anytime Accuracy*  
 284 *(AAA)* (Caccia et al., 2021). AAA is obtained by evaluating the model on the validation set after each  
 285 training batch, which we refer to as anytime evaluation. More details are reported in Appendix B.

286 **Training details.** In our experiments the backbone for all the Continual Learning strategies is the  
 287 *Graph Convolutional Network (GCN)* (Kipf & Welling, 2017). For CoraFull, Arxiv and Amazon  
 288 Computer datasets we use a 2-layer GCN with a fixed hidden dimension of 256 units as done by  
 289 Zhang et al. (2022). On the Reddit dataset a single layer of GCN was used due to requirement of  
 290 efficiency for OCGL (Section 3): with an average degree of 984, considering even two layers would  
 291 require to use almost the entire graph for each mini-batch (this is further discussed in Section 7,  
 292 where we provide results with neighborhood sampling instead). We use Adam optimizer (Kingma  
 293 & Ba, 2017) without weight decay, tuning the learning rate as an hyperparameter with the protocol  
 294 defined below. We consider the batch size to be fixed, as it could be dependent on the real world  
 295 problem, and we experiment with two different sizes for each dataset. For the smaller datasets  
 296 CoraFull and Amazon Computer we consider batches of 10 and 50 nodes, while for the larger Arxiv  
 297 and Reddit we use sizes 50 and 250 simply for computational reasons. As suggested by Aljundi  
 298 et al. (2019), multiple passes on the same mini-batch before passing to the next can be beneficial.  
 299 We therefore considered as an additional tuned hyperparameter whether to perform multiple passes  
 300 (5) on each batch. As a baseline, we use a *bare* model which is simply fine-tuned on the incoming  
 301 stream without any CL strategy applied. Additionally, we provide an upper bound in the form of a  
 model that is jointly trained offline on the entire graph.

302 **Hyperparameter selection.** Many works in the Continual Learning literature use a learning proto-  
 303 col that is akin to the classic machine learning setting, selecting hyperparameters by performing as  
 304 many full passes over the task sequence as required by a grid search. This protocol violates stricter  
 305 definitions of Lifelong Learning, where the stream is observed only once, and is indeed unrealistic  
 306 for real applications where a model needs to quickly adapt to changes in data distribution. Chaudhry  
 307 et al. (2018b) therefore proposed a more sensible hyperparameter selection protocol, which has now  
 308 been used in several works (Xu et al., 2020; Mai et al., 2022; Soutif-Cormerais et al., 2023) and that  
 309 we use for our experiments. With this protocol, only the first few tasks are used for hyperparameter  
 310 selection, allowing the model to perform multiple passes, with the same online setting, over them  
 311 to select the hyperparameters that lead to the best performance on validation nodes. In our case,  
 312 due to the different number of tasks in the datasets used, we considered 20% of the tasks for this  
 313 validation, with the exception of Amazon Computer where it was set to 2 as there are only 5 total  
 314 tasks. Hyperparameters are then selected based on the AA on the validation set at this validation  
 315 boundary. Details on the search space for the various methods are reported in Appendix C.

## 316 6 RESULTS ON FULL-NEIGHBORHOOD MINI-BATCHING

317 We start by discussing empirical results obtained on mini-batches containing the full neighborhoods  
 318 of newly presented nodes, in contrast with next Section 7 where neighborhood sampling is analyzed.  
 319 Although the neighboring expansion problem (Section 3.2) is present, the results of this section pro-  
 320 vide reference performance to be compared against the more realistic setting studied in the next  
 321 section. Moreover, we discuss forgetting issues, the impact of the batch size, and the sensitivity to  
 322 hyperparameters setting aside potential biases introduced by the neighborhood sampling. A com-  
 323

Table 1: Performance comparison on CoraFull with full neighborhood.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	15.97 $\pm$ 0.99	24.52 $\pm$ 0.87	-41.84 $\pm$ 1.89	11.86 $\pm$ 3.11	24.88 $\pm$ 0.67	-77.67 $\pm$ 3.37
ER	28.32 $\pm$ 3.20	35.33 $\pm$ 0.65	-63.29 $\pm$ 3.59	13.74 $\pm$ 2.04	30.00 $\pm$ 0.21	-75.10 $\pm$ 2.62
EWC	29.34 $\pm$ 2.82	43.96 $\pm$ 0.78	-20.06 $\pm$ 5.40	29.29 $\pm$ 1.45	46.55 $\pm$ 2.31	-17.56 $\pm$ 4.28
A-GEM	31.25 $\pm$ 1.58	44.30 $\pm$ 1.95	-57.88 $\pm$ 2.45	30.22 $\pm$ 2.25	39.97 $\pm$ 0.64	-55.56 $\pm$ 3.31
LWF	19.88 $\pm$ 2.43	33.37 $\pm$ 1.34	-48.64 $\pm$ 3.53	20.70 $\pm$ 1.87	28.28 $\pm$ 0.66	-55.87 $\pm$ 2.84
MAS	29.56 $\pm$ 1.58	44.46 $\pm$ 1.28	-15.35 $\pm$ 1.69	30.08 $\pm$ 1.00	52.01 $\pm$ 1.30	-10.62 $\pm$ 1.33
TWP	20.33 $\pm$ 2.36	28.70 $\pm$ 0.64	-64.81 $\pm$ 2.22	26.61 $\pm$ 2.65	36.45 $\pm$ 1.42	-60.79 $\pm$ 3.14
JOINT	67.55 $\pm$ 0.05	-	-	67.55 $\pm$ 0.05	-	-

Table 2: Performance comparison on Arxiv with full neighborhood.

METHOD	BATCH SIZE 50			BATCH SIZE 250		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	3.19 $\pm$ 0.50	11.09 $\pm$ 0.10	-52.49 $\pm$ 2.06	4.66 $\pm$ 0.90	11.45 $\pm$ 0.35	-41.36 $\pm$ 1.44
ER	5.92 $\pm$ 1.53	16.96 $\pm$ 0.87	-49.19 $\pm$ 2.97	5.30 $\pm$ 0.76	18.61 $\pm$ 0.67	-59.95 $\pm$ 1.83
EWC	4.43 $\pm$ 0.33	10.66 $\pm$ 0.12	-65.51 $\pm$ 0.22	2.19 $\pm$ 1.62	10.96 $\pm$ 1.13	-26.93 $\pm$ 0.79
A-GEM	16.14 $\pm$ 0.90	26.10 $\pm$ 0.34	-41.13 $\pm$ 0.53	10.61 $\pm$ 1.28	22.73 $\pm$ 0.43	-44.83 $\pm$ 1.40
LWF	3.72 $\pm$ 0.71	11.37 $\pm$ 0.46	-51.09 $\pm$ 1.60	4.55 $\pm$ 1.14	10.88 $\pm$ 0.17	-50.46 $\pm$ 2.69
MAS	5.25 $\pm$ 0.50	12.48 $\pm$ 0.80	-3.07 $\pm$ 0.49	4.51 $\pm$ 1.14	13.97 $\pm$ 0.67	-34.68 $\pm$ 1.03
TWP	2.30 $\pm$ 1.10	8.77 $\pm$ 1.29	-21.25 $\pm$ 6.08	4.43 $\pm$ 0.96	13.23 $\pm$ 1.88	-30.94 $\pm$ 3.59
JOINT	46.85 $\pm$ 0.44	-	-	46.85 $\pm$ 0.44	-	-

Table 3: Performance comparison on Reddit with full neighborhood.

METHOD	BATCH SIZE 50			BATCH SIZE 250		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	22.16 $\pm$ 1.26	39.12 $\pm$ 3.15	-62.68 $\pm$ 1.59	21.00 $\pm$ 1.61	44.66 $\pm$ 3.30	-73.06 $\pm$ 1.50
ER	33.39 $\pm$ 2.12	64.11 $\pm$ 0.51	-64.99 $\pm$ 2.18	36.93 $\pm$ 1.67	61.85 $\pm$ 0.46	-60.66 $\pm$ 1.66
EWC	22.16 $\pm$ 1.26	39.12 $\pm$ 3.15	-62.68 $\pm$ 1.59	18.51 $\pm$ 2.80	37.65 $\pm$ 4.90	-67.89 $\pm$ 3.09
A-GEM	57.71 $\pm$ 2.61	68.22 $\pm$ 0.35	-36.45 $\pm$ 2.56	35.54 $\pm$ 1.27	51.03 $\pm$ 3.54	-54.60 $\pm$ 1.28
LWF	18.31 $\pm$ 2.34	41.08 $\pm$ 3.20	-59.92 $\pm$ 3.50	21.63 $\pm$ 1.99	43.73 $\pm$ 2.32	-68.57 $\pm$ 2.81
MAS	30.06 $\pm$ 1.72	50.31 $\pm$ 2.18	-46.72 $\pm$ 1.33	21.00 $\pm$ 1.61	44.66 $\pm$ 3.30	-73.06 $\pm$ 1.50
TWP	22.16 $\pm$ 1.26	39.12 $\pm$ 3.15	-62.68 $\pm$ 1.59	17.64 $\pm$ 2.29	40.63 $\pm$ 3.17	-72.79 $\pm$ 2.37
JOINT	90.02 $\pm$ 0.12	-	-	90.02 $\pm$ 0.12	-	-

Table 4: Performance comparison on Amazon Computer with full neighborhood.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	13.39 $\pm$ 7.46	40.19 $\pm$ 2.21	-62.38 $\pm$ 7.69	18.54 $\pm$ 0.27	40.76 $\pm$ 0.52	-74.32 $\pm$ 2.26
ER	27.00 $\pm$ 4.06	47.48 $\pm$ 0.82	-65.36 $\pm$ 4.28	19.74 $\pm$ 0.48	43.29 $\pm$ 0.59	-76.45 $\pm$ 0.61
EWC	18.97 $\pm$ 0.16	41.84 $\pm$ 0.40	-74.89 $\pm$ 2.68	18.61 $\pm$ 0.27	40.78 $\pm$ 0.53	-74.64 $\pm$ 2.38
A-GEM	35.50 $\pm$ 0.58	57.22 $\pm$ 0.92	-57.90 $\pm$ 3.41	21.38 $\pm$ 0.18	48.94 $\pm$ 1.03	-74.58 $\pm$ 1.00
LWF	6.55 $\pm$ 6.47	36.29 $\pm$ 2.88	-48.19 $\pm$ 16.77	3.24 $\pm$ 0.36	24.60 $\pm$ 1.56	-21.84 $\pm$ 5.97
MAS	21.85 $\pm$ 4.84	46.46 $\pm$ 2.32	-38.97 $\pm$ 9.99	18.62 $\pm$ 6.09	44.84 $\pm$ 4.74	-41.79 $\pm$ 11.41
TWP	18.49 $\pm$ 0.66	41.27 $\pm$ 0.55	-72.43 $\pm$ 3.74	7.78 $\pm$ 6.21	35.65 $\pm$ 5.47	-46.67 $\pm$ 16.90
JOINT	72.06 $\pm$ 9.24	-	-	72.06 $\pm$ 9.24	-	-

parison of the considered CL methods for the four benchmark datasets is reported in Tables 1-4. All experiments were repeated five times with different initializations, and the metrics are reported as mean and standard deviation across runs. Additionally, we plot in Figure 3 of Appendix D the performance measured using anytime evaluation.

**Final Average Accuracy.** Across all datasets, we observe that none of the compared strategies comes close to the upper bound consisting of joint training. In general, the considered replay methods A-GEM and ER achieve higher final AA compared to the baseline and regularization methods. This can be expected, as rehearsal methods in Continual Learning are generally known to achieve most of the state-of-the-art results (Van De Ven et al., 2022); this holds true also for CGL (Zhang et al., 2022). In particular, A-GEM appears to be the best performing strategy in all cases (except on Reddit when using batch size 250 where it is surpassed by ER), often with a large margin, such as on Arxiv, where it is the only approach with a considerable improvement on the fine-tuned bare model. Regularization methods struggle more, as only on CoraFull EWC and MAS are competitive with A-GEM, while on other datasets their performances are closer to the lower bound.

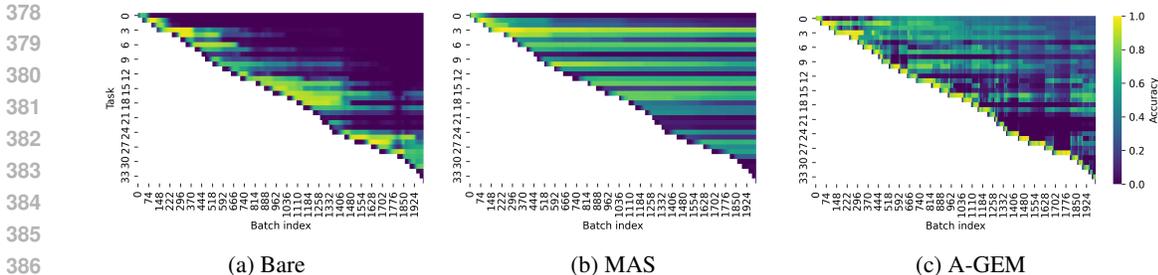


Figure 1: Anytime evaluation by task: a breakdown of model performance at the end of each training batch for three selected methods on CoraFull with batch size 10.

**Anytime Average Accuracy.** Looking at Average Accuracy gives us an easy way to compare the performance after the entire learning process. However, since in the online setting we expect the model to be ready to make predictions at any time, it is arguably more useful to look at Average Anytime Accuracy as a metric of performance over time. For this purpose Figure 3 (Appendix D) can give us some additional insights into the performance trends through the node stream. We note how it is natural and expected that accuracy tends to decrease with the batch index, as new classes are introduced and the classification task gets increasingly complex. In particular, from Figures 3a and 3b we see that thanks to regularization, on CoraFull the performance of EWC and MAS is much more stable compared to A-GEM, which is much more sensitive to sudden distribution shifts. In this case these two methods might thus be preferable for the stability of their performance. On the other three datasets though the higher stability does not offset their poorer performance. Another phenomenon we note is that while the performance of regularization methods is almost monotonically decreasing, the performance of ER and A-GEM shows much higher variations. We observe abrupt falls at some task boundaries, unexplained by the simple introduction of two new classes and thus due to catastrophic forgetting caused by interference of the current task with old ones. On the other hand, there are also instances where performance increases much more than what could be done with perfect accuracy on the current task, indicating successful backward transfer.

**Forgetting and strategy comparison.** Average Forgetting by itself is not a reliable metric of performance, and must be considered together with accuracy to assess the results of a Continual Learning methods: indeed, not learning anything would lead to 0 forgetting, but this is not our main goal. Indeed, Overall the regularization methods achieve a lower AF score. However, due to how AF is defined, this can be due to lower performances also right after task observation. To better understand the different results in light of the choice of strategy, we report in Figure 1 a more detailed breakdown of accuracy by task for three representative methods. With this comparison we clearly see how the fine-tuned baseline (1a), while it can learn new tasks, is not able to retain past knowledge, which is relatively quickly forgotten. MAS on the other hand (1b), thanks to its regularization is able to preserve quite well the knowledge it has learned, but it struggles to acquire new knowledge as it is presented with more and more classes. Instead A-GEM (1c) seems to achieve a better trade-off between stability and plasticity, maintaining the capacity to learn new tasks while generally preserving acquired knowledge, even showing signs of backward transfer at some task boundaries.

**Impact of batch size.** With regards to the dimensions of the node batches in the stream, we do not observe clear general patterns, as the two considered sizes have somewhat similar results. The only strategy that consistently benefits from a smaller batch size is A-GEM, with considerable increase of performance on all datasets except for CoraFull. For this method (and to a lower extent also for ER), a smaller batch size could thus be a regularizing factor. The fact that very small batch sizes in an online setting can lead to relatively good performance is also encouraging for the future development of OCGL techniques, despite the challenges of catastrophic forgetting (and the online setting for graphs itself, as we discuss in Section 7) that still need to be further addressed.

**Sensitivity to hyperparameters.** In our experiments the backbone model architecture, including number of layers and hidden units, is kept fixed. We conducted an ablation study on CoraFull to assess the impact of these choices, with full results reported in Appendix F. Decreasing the number of hidden units to 128 lowers performances overall, while increasing it to 512 has mixed results, with ER and A-GEM showing small changes, LwF, MAS and TWP scoring lower, but EWC reaches 40% AA. Using 3 GCN layers all results get worse, while with 1 layer there is a general smaller decrease,

432 except for the bare baseline and ER which improve compared with 2 layers. We note how the  
 433 difference in results might also be due to a sort of *butterfly effect* caused by the hyperparameter  
 434 selection policy: since they are selected early on, they could sub-optimal for the entire stream.  
 435

436  
 437 **7 RESULTS ON NEIGHBORHOOD SAMPLING**

438  
 439 **Neighborhood expansion.** To assess the neighborhood expansion phenomenon on the four datasets  
 440 we used in our experiments, we plot in Figure 2 the size of the  $l$ -hop neighborhood in the evolving  
 441 graph of each mini-batch in the node stream, for some values of  $l$ . The Reddit graph in particular is  
 442 very well connected, with two hops containing the majority of the graph, and three hops practically  
 443 all nodes. This is the motivation for our use of only one GCN layer on this dataset. CoraFull and  
 444 Arxiv have a much more contained neighborhood expansion, while Amazon Computer with two  
 445 hops covers about half of the nodes in the worst cases. Additional plots are shown in Appendix E,  
 446 with the addition of the number of edges connecting the various hops, which can be used as a proxy  
 447 of computational complexity.

448  
 449 **Neighborhood sampling.** Following the same experimental setup outlined in Section 5, we con-  
 450 ducted experiments with neighborhood sampling instead of using full neighborhood information.  
 451 As in this case we use sampling to address the problem of neighborhood expansion, we use 2 GCN  
 452 layers on all datasets. We choose the number of nodes to be sampled with a double rationale: we  
 453 want to guarantee that processing each mini-batch requires much less than the full graph to conform  
 454 with the requirements of the online setting, and we want to sample significantly less nodes than the  
 455 average degree for our analysis of the sampling strategy to be meaningful. Therefore, we fixed the  
 456 number of neighbors to 5 for CoraFull and Amazon Computer, 10 for Arxiv and 15 for Reddit. Full  
 457 results are reported in Tables 5-8 and plots with anytime evaluation are shown in Figure 4 (Appendix  
 458 D). Specifically, on CoraFull only EWC and A-GEM maintain results similar to those obtained with  
 459 full neighborhood aggregation, while the performance of other models degrades sharply. An Arxiv,  
 460 where the catastrophic forgetting problem lead to low accuracy scores, the use of sampling does  
 461 not appear to particularly worsen performance. In fact, the results of ER benefit from it, surpassing  
 462 A-GEM. The same higher performance of ER is shown on Amazon Computer, where additionally

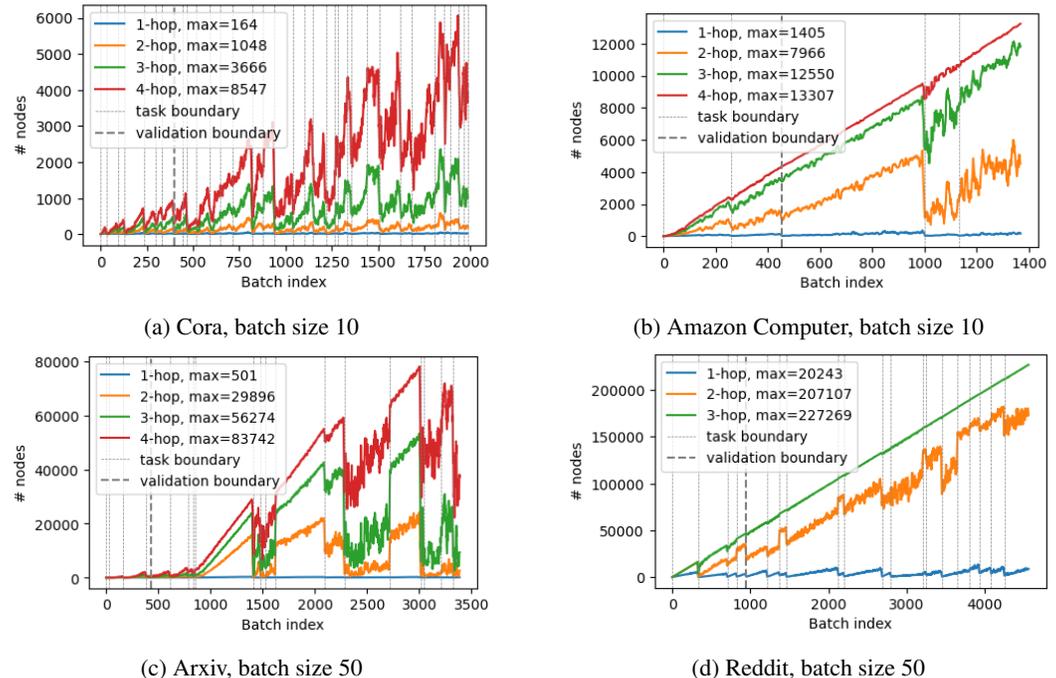


Figure 2: Number of nodes in the union of  $l$ -hop neighborhoods of each training batch. Smoothed with rolling average over windows of 10 batches for readability, maximum is reported in the legend.

Table 5: Performance comparison on CoraFull with neighborhood sampling (5 nodes).

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	8.76 $\pm$ 1.53	23.47 $\pm$ 2.32	-30.88 $\pm$ 2.79	9.80 $\pm$ 1.57	23.33 $\pm$ 1.32	-36.59 $\pm$ 2.87
ER	19.17 $\pm$ 1.79	34.55 $\pm$ 0.69	-71.14 $\pm$ 2.71	10.26 $\pm$ 1.79	25.74 $\pm$ 0.58	-37.05 $\pm$ 3.22
EWC	26.30 $\pm$ 3.40	41.46 $\pm$ 1.97	-21.71 $\pm$ 5.19	30.30 $\pm$ 4.25	44.34 $\pm$ 2.89	-21.07 $\pm$ 4.75
A-GEM	33.08 $\pm$ 0.84	33.60 $\pm$ 5.88	-58.78 $\pm$ 0.67	27.21 $\pm$ 1.26	37.81 $\pm$ 2.38	-34.14 $\pm$ 1.97
LWF	14.87 $\pm$ 0.48	25.92 $\pm$ 0.81	-49.70 $\pm$ 1.68	12.01 $\pm$ 0.32	27.20 $\pm$ 3.49	-12.97 $\pm$ 1.18
MAS	12.64 $\pm$ 2.32	39.71 $\pm$ 1.08	-35.23 $\pm$ 2.27	26.66 $\pm$ 2.40	46.70 $\pm$ 0.41	-32.75 $\pm$ 2.05
TWP	11.73 $\pm$ 1.73	23.75 $\pm$ 1.75	-28.58 $\pm$ 4.52	12.87 $\pm$ 3.14	25.05 $\pm$ 0.94	-34.96 $\pm$ 5.11

Table 6: Performance comparison on Arxiv with neighborhood sampling (10 nodes).

METHOD	BATCH SIZE 50			BATCH SIZE 250		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	4.74 $\pm$ 0.08	11.88 $\pm$ 0.06	-82.97 $\pm$ 1.93	4.67 $\pm$ 0.83	11.93 $\pm$ 0.08	-59.34 $\pm$ 2.33
ER	18.41 $\pm$ 2.31	36.34 $\pm$ 0.26	-72.96 $\pm$ 2.62	16.96 $\pm$ 1.45	32.91 $\pm$ 0.57	-73.34 $\pm$ 1.60
EWC	4.81 $\pm$ 0.08	12.35 $\pm$ 0.19	-77.99 $\pm$ 1.52	8.49 $\pm$ 3.20	14.84 $\pm$ 0.38	-64.19 $\pm$ 4.70
A-GEM	16.43 $\pm$ 3.20	27.99 $\pm$ 0.52	-73.36 $\pm$ 3.45	12.39 $\pm$ 1.14	21.42 $\pm$ 0.57	-74.42 $\pm$ 2.18
LWF	4.79 $\pm$ 0.08	11.85 $\pm$ 0.15	-79.49 $\pm$ 1.13	4.48 $\pm$ 1.05	11.88 $\pm$ 0.15	-61.81 $\pm$ 1.90
MAS	3.35 $\pm$ 0.99	12.49 $\pm$ 0.44	-32.58 $\pm$ 1.73	4.56 $\pm$ 1.06	13.48 $\pm$ 0.72	-49.05 $\pm$ 2.86
TWP	4.74 $\pm$ 0.05	12.17 $\pm$ 0.14	-80.22 $\pm$ 0.41	3.65 $\pm$ 0.94	11.81 $\pm$ 0.06	-61.81 $\pm$ 1.90

Table 7: Performance comparison on Reddit with neighborhood sampling (15 nodes).

METHOD	BATCH SIZE 50			BATCH SIZE 250		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	12.86 $\pm$ 2.56	36.72 $\pm$ 2.20	-85.00 $\pm$ 2.54	20.44 $\pm$ 3.20	42.20 $\pm$ 3.85	-58.20 $\pm$ 2.50
ER	17.84 $\pm$ 2.89	46.24 $\pm$ 0.54	-81.01 $\pm$ 2.85	19.18 $\pm$ 3.80	45.18 $\pm$ 3.94	-60.14 $\pm$ 5.59
EWC	4.29 $\pm$ 2.73	20.40 $\pm$ 7.00	-12.64 $\pm$ 1.54	5.36 $\pm$ 3.16	26.56 $\pm$ 6.43	-14.26 $\pm$ 1.46
A-GEM	43.24 $\pm$ 4.08	63.44 $\pm$ 3.05	-55.09 $\pm$ 4.18	21.97 $\pm$ 4.03	59.51 $\pm$ 2.36	-71.29 $\pm$ 3.28
LWF	12.77 $\pm$ 1.69	37.42 $\pm$ 1.46	-83.82 $\pm$ 1.89	16.64 $\pm$ 1.32	43.40 $\pm$ 1.98	-76.59 $\pm$ 1.76
MAS	9.91 $\pm$ 0.84	35.86 $\pm$ 3.21	-88.15 $\pm$ 0.70	15.44 $\pm$ 2.54	37.62 $\pm$ 3.97	-80.23 $\pm$ 2.94
TWP	12.60 $\pm$ 2.13	36.10 $\pm$ 0.26	-85.46 $\pm$ 2.23	20.16 $\pm$ 4.52	40.96 $\pm$ 3.84	-59.25 $\pm$ 7.10

Table 8: Performance comparison on Amazon Computer with neighborhood sampling (5 nodes).

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	19.34 $\pm$ 0.26	43.03 $\pm$ 0.15	-78.66 $\pm$ 0.35	18.47 $\pm$ 0.76	41.65 $\pm$ 0.33	-77.04 $\pm$ 1.71
ER	24.37 $\pm$ 2.03	52.34 $\pm$ 2.30	-59.93 $\pm$ 7.94	48.28 $\pm$ 7.83	67.00 $\pm$ 0.71	-45.77 $\pm$ 9.46
EWC	17.60 $\pm$ 2.53	40.12 $\pm$ 1.60	-35.37 $\pm$ 9.96	19.39 $\pm$ 0.18	43.32 $\pm$ 0.19	-78.15 $\pm$ 1.30
A-GEM	20.05 $\pm$ 0.59	50.36 $\pm$ 1.35	-77.13 $\pm$ 0.96	19.95 $\pm$ 1.16	50.28 $\pm$ 1.24	-77.50 $\pm$ 0.52
LWF	19.29 $\pm$ 0.48	42.95 $\pm$ 0.19	-78.35 $\pm$ 0.50	18.12 $\pm$ 0.62	41.09 $\pm$ 0.67	-76.37 $\pm$ 2.36
MAS	18.82 $\pm$ 1.19	42.92 $\pm$ 1.83	-60.99 $\pm$ 7.52	18.66 $\pm$ 0.18	43.17 $\pm$ 0.61	-77.33 $\pm$ 0.72
TWP	19.13 $\pm$ 0.45	42.98 $\pm$ 0.18	-78.54 $\pm$ 0.39	17.82 $\pm$ 0.45	41.78 $\pm$ 1.01	-71.46 $\pm$ 3.48

we do not see the same low performance of LwF as without sampling. Finally, on Reddit the performance degradation is more pronounced, as the number of neighbors was cut more substantially. In summary, as expected due to ignoring some neighborhood information, most of the performances are lowered, indicating that more research is required to properly address this issue in OGL.

## 8 CONCLUSIONS

In this paper, we introduced the formulation of the Online Continual Graph Learning setting, closing the gap between the Continual Graph Learning and Online Continual Learning literature. We adapted four node classification datasets to the proposed framework, constructing node streams starting from the class-incremental learning scenario. Our evaluation of suitably adapted Continual Learning methods highlights the higher performance of replay-based methods. Finally, we raise the issue of neighborhood expansion for GNNs, proposing neighborhood sampling as a straight-forward solution to bound the computational cost of training on each mini-batch. In future works, we plan to further tackle the issue of neighborhood expansion, developing tailored strategies that can ensure computational efficiency while better addressing the catastrophic forgetting problem. We further intend to consider more diverse node stream construction and additional tasks such as link prediction.

## REFERENCES

- 540  
541  
542 Rahaf Aljundi, Francesca Babiloni, Mohamed Elhoseiny, Marcus Rohrbach, and Tinne Tuytelaars.  
543 Memory Aware Synapses: Learning what (not) to forget. In *Proceedings of the European Con-*  
544 *ference on Computer Vision (ECCV)*, pp. 139–154, 2018.
- 545 Rahaf Aljundi, Eugene Belilovsky, Tinne Tuytelaars, Laurent Charlin, Massimo Caccia, Min Lin,  
546 and Lucas Page-Caccia. Online Continual Learning with Maximal Interfered Retrieval. In *Ad-*  
547 *vances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- 548  
549 Aleksandar Bojchevski and Stephan Günnemann. Deep Gaussian Embedding of Graphs: Unsuper-  
550 vised Inductive Learning via Ranking. In *International Conference on Learning Representations*,  
551 February 2018.
- 552 Lucas Caccia, Rahaf Aljundi, Nader Asadi, Tinne Tuytelaars, Joelle Pineau, and Eugene Belilovsky.  
553 New Insights on Reducing Abrupt Representation Change in Online Continual Learning. In *In-*  
554 *ternational Conference on Learning Representations*, October 2021.
- 555  
556 Antonio Carta, Andrea Cossu, Federico Errica, and Davide Bacciu. Catastrophic Forgetting in  
557 Deep Graph Networks: A Graph Classification Benchmark. *Frontiers in Artificial Intelligence*, 5:  
558 824655, February 2022. doi: 10.3389/frai.2022.824655.
- 559 Arslan Chaudhry, Puneet K. Dokania, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Riemannian  
560 Walk for Incremental Learning: Understanding Forgetting and Intransigence. In *Proceedings of*  
561 *the European Conference on Computer Vision (ECCV)*, pp. 532–547, 2018a.
- 562  
563 Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient  
564 Lifelong Learning with A-GEM. In *ICLR*, September 2018b.
- 565  
566 Arslan Chaudhry, Marcus Rohrbach, Mohamed Elhoseiny, Thalaiyasingam Ajanthan, Puneet K.  
567 Dokania, Philip H. S. Torr, and Marc’Aurelio Ranzato. On Tiny Episodic Memories in Continual  
568 Learning, June 2019. arXiv:1902.10486 [cs, stat].
- 569 Jie Chen, Tengfei Ma, and Cao Xiao. FastGCN: Fast Learning with Graph Convolutional Networks  
570 via Importance Sampling. In *International Conference on Learning Representations*, February  
571 2018.
- 572  
573 Xu Chen, Junshan Wang, and Kunqing Xie. Trafficstream: A streaming traffic flow forecasting  
574 framework based on graph neural networks and continual learning. In *Proceedings of the Thirtieth*  
575 *International Joint Conference on Artificial Intelligence, IJCAI-21*, pp. 3620–3626. International  
576 Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/498.
- 577 Zhiyuan Chen and Bing Liu. *Lifelong Machine Learning*. Synthesis Lectures on Artificial Intel-  
578 ligence and Machine Learning. Springer International Publishing, Cham, 2018. doi: 10.1007/  
579 978-3-031-01581-6.
- 580  
581 Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-GCN: An  
582 Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. In *Proceedings*  
583 *of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*,  
584 pp. 257–266, July 2019. doi: 10.1145/3292500.3330925.
- 585  
586 Andrea Cini, Ivan Marisca, Daniele Zambon, and Cesare Alippi. Graph Deep Learning for Time  
587 Series Forecasting, October 2023a. arXiv:2310.15978 [cs].
- 588  
589 Andrea Cini, Ivan Marisca, Daniele Zambon, and Cesare Alippi. Taming Local Effects in Graph-  
590 based Spatiotemporal Forecasting. In *Advances in Neural Information Processing Systems*, vol-  
591 ume 36, pp. 55375–55393. Curran Associates, Inc., 2023b.
- 592  
593 Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah Parisot, Xu Jia, Aleš Leonardis, Gregory  
594 Slabaugh, and Tinne Tuytelaars. A Continual Learning Survey: Defying Forgetting in Classifica-  
595 tion Tasks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3366–3385,  
596 July 2022. doi: 10.1109/TPAMI.2021.3057446.

- 594 Falih Gozi Febrinanto, Feng Xia, Kristen Moore, Chandra Thapa, and Charu Aggarwal. Graph  
595 Lifelong Learning: A Survey. *IEEE Computational Intelligence Magazine*, 18(1):32–51, Febru-  
596 ary 2023. doi: 10.1109/MCI.2022.3222049.
- 597
- 598 Chrisantha Fernando, Dylan Banarse, Charles Blundell, Yori Zwols, David Ha, Andrei A. Rusu,  
599 Alexander Pritzel, and Daan Wierstra. PathNet: Evolution Channels Gradient Descent in Super  
600 Neural Networks, January 2017. arXiv:1701.08734 [cs].
- 601 Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural  
602 Message Passing for Quantum Chemistry. In *Proceedings of the 34th International Conference*  
603 *on Machine Learning*, pp. 1263–1272. PMLR, July 2017. ISSN: 2640-3498.
- 604
- 605 Alessio Gravina and Davide Bacciu. Deep learning for dynamic graphs: Models and benchmarks.  
606 *IEEE Transactions on Neural Networks and Learning Systems*, 35(9):11788–11801, 2024. doi:  
607 10.1109/TNNLS.2024.3379735.
- 608 Nuwan Gunasekara, Bernhard Pfahringer, Heitor Murilo Gomes, and Albert Bifet. Survey on Online  
609 Streaming Continual Learning. In *Proceedings of the Thirty-Second International Joint Confer-*  
610 *ence on Artificial Intelligence*, pp. 6628–6637, Macau, SAR China, August 2023. International  
611 Joint Conferences on Artificial Intelligence Organization. doi: 10.24963/ijcai.2023/743.
- 612
- 613 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive Representation Learning on Large  
614 Graphs. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates,  
615 Inc., 2017.
- 616 Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the Knowledge in a Neural Network,  
617 March 2015. arXiv:1503.02531 [cs, stat].
- 618
- 619 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,  
620 and Jure Leskovec. Open Graph Benchmark: Datasets for Machine Learning on Graphs, February  
621 2021. arXiv:2005.00687 [cs, stat].
- 622 Shenyang Huang, Farimah Poursafaei, Jacob Danovitch, Matthias Fey, Weihua Hu, Emanuele Rossi,  
623 Jure Leskovec, Michael Bronstein, Guillaume Rabusseau, and Reihaneh Rabbany. Temporal  
624 graph benchmark for machine learning on temporal graphs. In A. Oh, T. Naumann, A. Globerson,  
625 K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*,  
626 volume 36, pp. 2056–2073. Curran Associates, Inc., 2023.
- 627
- 628 Ming Jin, Huan Yee Koh, Qingsong Wen, Daniele Zambon, Cesare Alippi, Geoffrey I. Webb, Irwin  
629 King, and Shirui Pan. A survey on graph neural networks for time series: Forecasting, classifi-  
630 cation, imputation, and anomaly detection. *IEEE Transactions on Pattern Analysis and Machine*  
631 *Intelligence*, 2024. doi: 10.1109/TPAMI.2024.3443141.
- 632 Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, Peter Forsyth, and  
633 Pascal Poupert. Representation learning for dynamic graphs: A survey. *Journal of Machine*  
634 *Learning Research*, 21(70):1–73, 2020.
- 635
- 636 Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization, January 2017.  
637 arXiv:1412.6980 [cs].
- 638
- 639 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional net-  
640 works. In *International Conference on Learning Representations*, 2017.
- 641 James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A.  
642 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hass-  
643 abis, Claudia Clopath, Dharshan Kumaran, and Raia Hadsell. Overcoming catastrophic forgetting  
644 in neural networks. *Proceedings of the National Academy of Sciences*, 114(13):3521–3526, March  
645 2017. doi: 10.1073/pnas.1611835114.
- 646
- 647 Hyunseo Koh, Dahyun Kim, Jung-Woo Ha, and Jonghyun Choi. Online Continual Learning on Class  
Incremental Blurry Task Configuration with Anytime Inference. In *International Conference on*  
*Learning Representations*, October 2021.

- 648 Zhizhong Li and Derek Hoiem. Learning without Forgetting. *IEEE Transactions on Pattern Anal-*  
649 *ysis and Machine Intelligence*, 40(12):2935–2947, December 2018. doi: 10.1109/TPAMI.2017.  
650 2773081.
- 651 Huihui Liu, Yiding Yang, and Xinchao Wang. Overcoming Catastrophic Forgetting in Graph Neural  
652 Networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(10):8653–8661,  
653 May 2021. doi: 10.1609/aaai.v35i10.17049.
- 654 Jesus L. Lobo, Igor Ballesteros, Izaskun Oregi, Javier Del Ser, and Sancho Salcedo-Sanz. Stream  
655 Learning in Energy IoT Systems: A Case Study in Combined Cycle Power Plants. *Energies*, 13  
656 (3):740, January 2020. doi: 10.3390/en13030740.
- 657 David Lopez-Paz and Marc’ Aurelio Ranzato. Gradient Episodic Memory for Continual Learning.  
658 In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 659 Zheda Mai, Ruiwen Li, Jihwan Jeong, David Quispe, Hyunwoo Kim, and Scott Sanner. Online  
660 continual learning in image classification: An empirical survey. *Neurocomputing*, 469:28–51,  
661 January 2022. doi: 10.1016/j.neucom.2021.10.021.
- 662 Marc Masana, Xialei Liu, Bartłomiej Twardowski, Mikel Menta, Andrew D. Bagdanov, and Joost  
663 Van De Weijer. Class-Incremental Learning: Survey and Performance Evaluation on Image Clas-  
664 sification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2022. doi:  
665 10.1109/TPAMI.2022.3213473.
- 666 Nicolas Y. Masse, Gregory D. Grant, and David J. Freedman. Alleviating catastrophic forgetting  
667 using context-dependent gating and synaptic stabilization. *Proceedings of the National Academy*  
668 *of Sciences*, 115(44), October 2018. doi: 10.1073/pnas.1803839115.
- 669 Alessio Micheli. Neural Network for Graphs: A Contextual Constructive Approach. *IEEE Trans-*  
670 *actions on Neural Networks*, 20(3):498–511, March 2009. doi: 10.1109/TNN.2008.2010350.
- 671 German I. Parisi and Vincenzo Lomonaco. Online Continual Learning on Sequences. In *Recent*  
672 *trends in learning from data: tutorials from the INNS Big Data and Deep Learning Conference*  
673 *(INNSBDDL2019)*, volume 896 of *Studies in computational intelligence*, pp. 197–221. Springer,  
674 Cham, 2020. doi: 10.1007/978-3-030-43883-8\_8. arXiv:2003.09114 [cs].
- 675 German I. Parisi, Ronald Kemker, Jose L. Part, Christopher Kanan, and Stefan Wermter. Continual  
676 lifelong learning with neural networks: A review. *Neural Networks*, 113:54–71, May 2019. doi:  
677 10.1016/j.neunet.2019.01.012.
- 678 Massimo Perini, Giorgia Ramponi, Paris Carbone, and Vasiliki Kalavri. Learning on streaming  
679 graphs with experience replay. In *Proceedings of the 37th ACM/SIGAPP Symposium on Applied*  
680 *Computing*, pp. 470–478, Virtual Event, April 2022. ACM. doi: 10.1145/3477314.3507113.
- 681 Ameya Prabhu, Philip H. S. Torr, and Puneet K. Dokania. GDumb: A Simple Approach that Ques-  
682 tions Our Progress in Continual Learning. In *Computer Vision – ECCV 2020*, pp. 524–540, Cham,  
683 2020. Springer International Publishing. doi: 10.1007/978-3-030-58536-5\_31.
- 684 Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H. Lampert. iCaRL:  
685 Incremental Classifier and Representation Learning. In *2017 IEEE Conference on Computer*  
686 *Vision and Pattern Recognition (CVPR)*, pp. 5533–5542, Honolulu, HI, July 2017. IEEE. doi:  
687 10.1109/CVPR.2017.587.
- 688 David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experi-  
689 ence Replay for Continual Learning. In *Advances in Neural Information Processing Systems*,  
690 volume 32. Curran Associates, Inc., 2019.
- 691 Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray  
692 Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive Neural Networks, October 2022.  
693 arXiv:1606.04671 [cs].
- 694 Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.  
695 The Graph Neural Network Model. *IEEE Transactions on Neural Networks*, 20(1):61–80, January  
696 2009. doi: 10.1109/TNN.2008.2005605.

- 702 Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls  
703 of Graph Neural Network Evaluation, June 2019. arXiv:1811.05868 [cs, stat].  
704
- 705 Albin Soutif-Cormerais, Antonio Carta, Andrea Cossu, Julio Hurtado, Vincenzo Lomonaco, Joost  
706 Van De Weijer, and Hamed Hemati. A Comprehensive Empirical Evaluation on Online Continual  
707 Learning. In *2023 IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*,  
708 pp. 3510–3520, October 2023. doi: 10.1109/ICCVW60793.2023.00378.
- 709 Vinicius M. A. Souza, Denis M. Dos Reis, André G. Maletzke, and Gustavo E. A. P. A. Batista.  
710 Challenges in benchmarking stream learning algorithms with real-world data. *Data Mining and*  
711 *Knowledge Discovery*, 34(6):1805–1858, November 2020. doi: 10.1007/s10618-020-00698-5.  
712
- 713 A. Sperduti and A. Starita. Supervised neural networks for the classification of structures. *IEEE*  
714 *Transactions on Neural Networks*, 8(3):714–735, May 1997. doi: 10.1109/72.572108.
- 715 Zonggui Tian, Du Zhang, and Hong-Ning Dai. Continual Learning on Graphs: A Survey, February  
716 2024. arXiv:2402.06330 [cs].  
717
- 718 Guido M. Van De Ven, Tinne Tuytelaars, and Andreas S. Tolias. Three types of incremental  
719 learning. *Nature Machine Intelligence*, 4(12):1185–1197, December 2022. doi: 10.1038/  
720 s42256-022-00568-3.
- 721 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua  
722 Bengio. Graph Attention Networks. In *International Conference on Learning Representations*,  
723 February 2018.
- 724 Jeffrey S. Vitter. Random sampling with a reservoir. *ACM Trans. Math. Softw.*, 11(1):37–57, March  
725 1985. doi: 10.1145/3147.3165.  
726
- 727 Junshan Wang, Guojie Song, Yi Wu, and Liang Wang. Streaming Graph Neural Networks via  
728 Continual Learning. In *Proceedings of the 29th ACM International Conference on Information*  
729 *& Knowledge Management*, pp. 1515–1524, Virtual Event Ireland, October 2020. ACM. doi:  
730 10.1145/3340531.3411963.
- 731 Yishi Xu, Yingxue Zhang, Wei Guo, Huifeng Guo, Ruiming Tang, and Mark Coates. GraphSAIL:  
732 Graph Structure Aware Incremental Learning for Recommender Systems. In *Proceedings of the*  
733 *29th ACM International Conference on Information & Knowledge Management, CIKM '20*, pp.  
734 2861–2868, New York, NY, USA, 2020. Association for Computing Machinery. doi: 10.1145/  
735 3340531.3412754.
- 736 Qiao Yuan, Sheng-Uei Guan, Pin Ni, Tianlun Luo, Ka Lok Man, Prudence Wong, and Victor Chang.  
737 Continual Graph Learning: A Survey, January 2023. arXiv:2301.12230 [cs].  
738
- 739 Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual Learning Through Synaptic Intelli-  
740 gence. In *Proceedings of the 34th International Conference on Machine Learning*, pp. 3987–3995.  
741 PMLR, July 2017.
- 742 Xikun Zhang, Dongjin Song, and Dacheng Tao. CGLB: Benchmark Tasks for Continual Graph  
743 Learning. *Advances in Neural Information Processing Systems*, 35:13006–13021, December  
744 2022.  
745
- 746 Xikun Zhang, Dongjin Song, and Dacheng Tao. Continual Learning on Graphs: Challenges, Solu-  
747 tions, and Opportunities, February 2024. arXiv:2402.11565 [cs].
- 748 Fan Zhou and Chengtai Cao. Overcoming Catastrophic Forgetting in Graph Neural Networks with  
749 Experience Replay. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(5):4714–  
750 4722, May 2021. doi: 10.1609/aaai.v35i5.16602.
- 751 Indre Zliobaite, Mykola Pechenizkiy, and Joao Gama. An Overview of Concept Drift Applications.  
752 In Nathalie Japkowicz and Jerzy Stefanowski (eds.), *Big Data Analysis, Studies in Big Data*, pp.  
753 91–114. Springer International Publishing AG, Cham, 2016. doi: 10.1007/978-3-319-26989-4.  
754  
755

## A DATASETS

In the experiments for this paper, we used four node-level classification graph datasets. The CoraFull dataset (Bojchevski & Günnemann, 2018) is a citation network where nodes represent research papers and edges indicate citations between them, with labels based on paper topics. Arxiv (Hu et al., 2021) is a larger citation network derived from arXiv papers in the Computer Science category. The Reddit dataset (Hamilton et al., 2017) consists of posts from different communities of the Reddit platform, where nodes represent posts, and edges connect posts commented on by the same user, forming a large interaction graph. Finally, Amazon Computer (Shchur et al., 2019) is a co-purchase network, where nodes are products and edges indicate frequently co-purchased items within the computers category on Amazon. Summary statistics for the four graphs are reported in Table 9.

Table 9: Dataset statistics.

DATASET	CORAFULL	ARXIV	REDDIT	AMAZON COMPUTER
# NODES	19,793	169,343	227,853	13,752
# EDGES	130,622	1,166,243	114,615,892	491,722
# CLASSES	70	40	40	10

## B METRICS

Thanks to the construction of the node stream starting from the class-incremental setting, we can use two widely used metrics in CL: *Average Accuracy (AA)* and *Average Forgetting (AF)* (Lopez-Paz & Ranzato, 2017). The most comprehensive metric for CL, from which AA and AF are derived, is the performance matrix  $M \in \mathbb{R}^{T \times T}$ , where  $T$  is the number of tasks and  $M_{i,j}$  is the test classification accuracy on task  $j$  after the model has observed task  $i$ . AA is then defined as  $AA = \frac{1}{T} \sum_{i=1}^T M_{T,i}$ , and average forgetting as  $AF = \frac{1}{T-1} \sum_{i=1}^{T-1} M_{T,i} - M_{i,i}$ . AA serves as a single value to quantify the performance of the model after having observed the entire sequence of tasks, or stream in our case. AF measures the performance degradation (forgetting), that occurs from when a task was just observed to the end of training.

To assess the performance of the model throughout the node stream, we also perform anytime evaluation, meaning that we evaluate the model on validation nodes after training on each mini-batch (Koh et al., 2021). This allows us to capture the performance at any point in time, and observe also graphically how the model reacts to changes in data distribution. We measure this with *Average Anytime Accuracy (AAA)* (Caccia et al., 2021), which is a generalization of average incremental accuracy for the online setting. Indicating with  $AA_t$  the average accuracy after training on batch  $t$ , and having  $n$  batches in total, AAA is defined as  $AAA = \frac{1}{n} \sum_{t=1}^n AA_t$ . This can be interpreted as an Area Under the Curve accuracy score (Koh et al., 2021).

## C HYPERPARAMETERS

A standard grid search was performed to select training hyperparameters for the models used in all experiments. We detail here the specific search space for each of the methods used in our comparisons. Two hyperparameters are common for all techniques: the learning rate, selected in the set  $\{0.01, 0.001, 0.0001, 0.00001\}$ , and the number of passes on each batch before passing to the next one, chosen between 1 and 5. No weight decay or dropout were used. Method specific hyperparameters are reported in Table 10, and specific details can be found in the original papers. In particular, the hyperparameters of regularization methods regulate the strenght of the regularization. For LwF a new hyperparameter has been introduced to adapt to the online setting: the number of batches after which to update the teacher model. For replay based methods we consider memory size (budget) and the proportion of memories to use with respect to each training batch.

Table 10: Method specific hyperparameters.

METHOD	HYPERPARAMETER CANDIDATES
ER	BUDGET: $\{100, 1000\}$ ; MEMORY_PROPORTION: $\{1, 2, 3\}$
EWC	LAMBDA: $\{10^0, 10^2, 10^4, 10^6, 10^8, 10^{10}\}$
A-GEM	BUDGET: $\{100, 1000\}$ ; MEMORY_PROPORTION: $\{1, 2, 3\}$
LwF	LAMBDA_DIST: $\{0.1, 1, 10\}$ ; T: $\{0.2, 2, 20\}$ , UPDATE_EVERY: $\{1, 10, 100\}$
MAS	LAMBDA: $\{10^0, 10^2, 10^4, 10^6, 10^8, 10^{10}\}$
TWP	LAMBDA_L: $\{10^2, 10^4, 10^6\}$ ; LAMBDA_T: $\{10^2, 10^4, 10^6\}$ ; BETA: $\{0.001, 0.01, 0.1\}$

## D ANYTIME EVALUATION PLOTS

We show here the plots with anytime evaluation, on all four datasets and with both choices of batch size. In Figure 3 we show the results using full neighborhood information (see Section 6 for all the results and comments), and in Figure 4 the anytime evaluation when using neighborhood sampling (more details in Section 7 of the main paper).

## E NEIGHBORHOOD EXPANSION

Neighborhood expansion has been identified as the main issue with the online setting fro graphs. In Figure 2 we showed the phenomenon, and we report here additional plots for all datasets and batch sizes. In particular, in Figures 5-7 we provide measure neighborhood expansion non only in term of number of nodes in the  $l$ -hop neighborhood, as done in Section 7, but we additionally count the number of edges between hops, which gives us a good proxy of computational time. Interestingly, the number of edges increases more drastically than the number of nodes when incrementing the number of layers, further confirming the need for a strategy such as sampling to maintain the computational complexity bounded within the limits of the online learning setting.

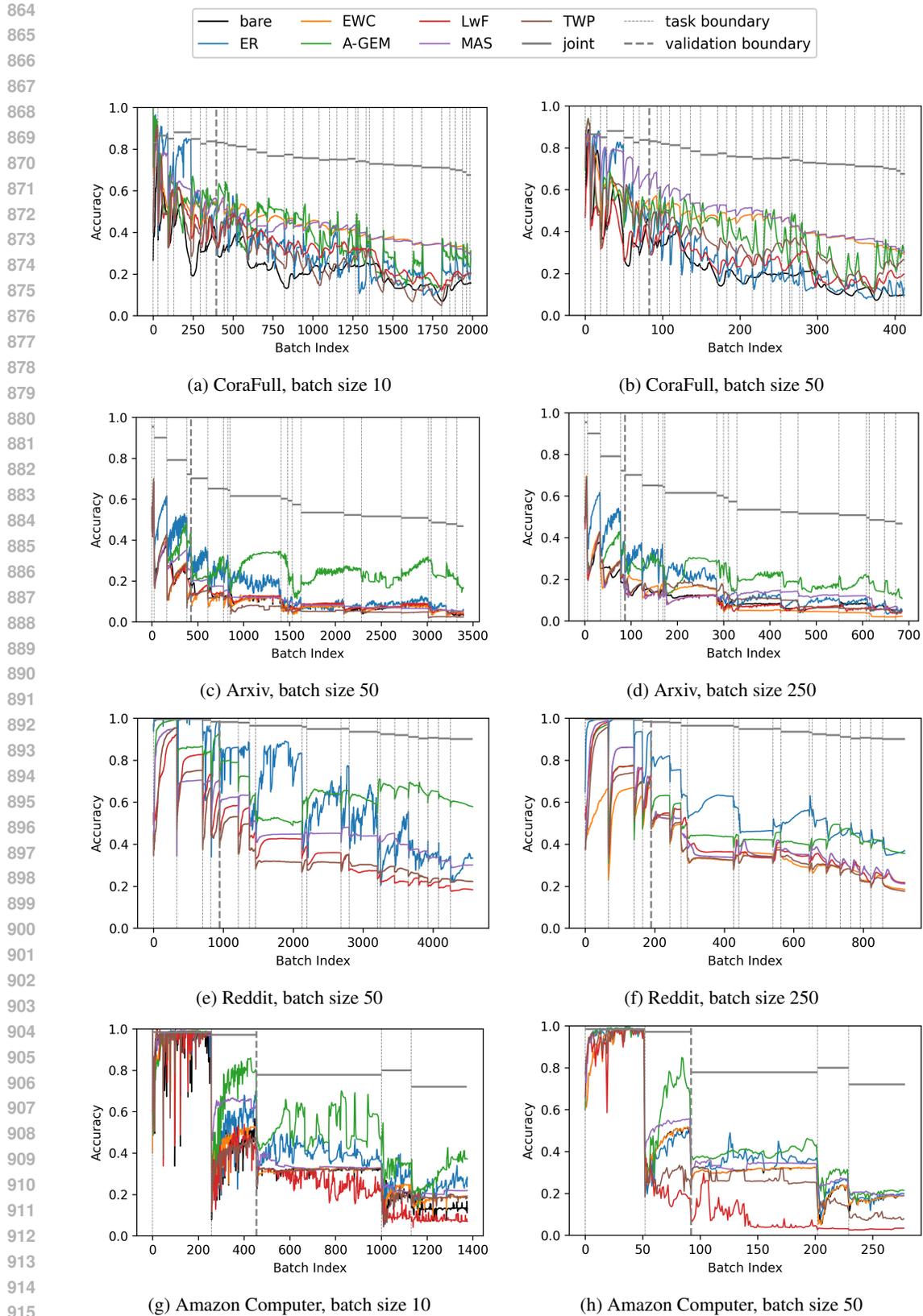


Figure 3: Anytime evaluation, showing AA on validation nodes after training on each mini-batch. We highlight the boundaries between tasks and the threshold up to which hyperparameter selection is performed.

918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971

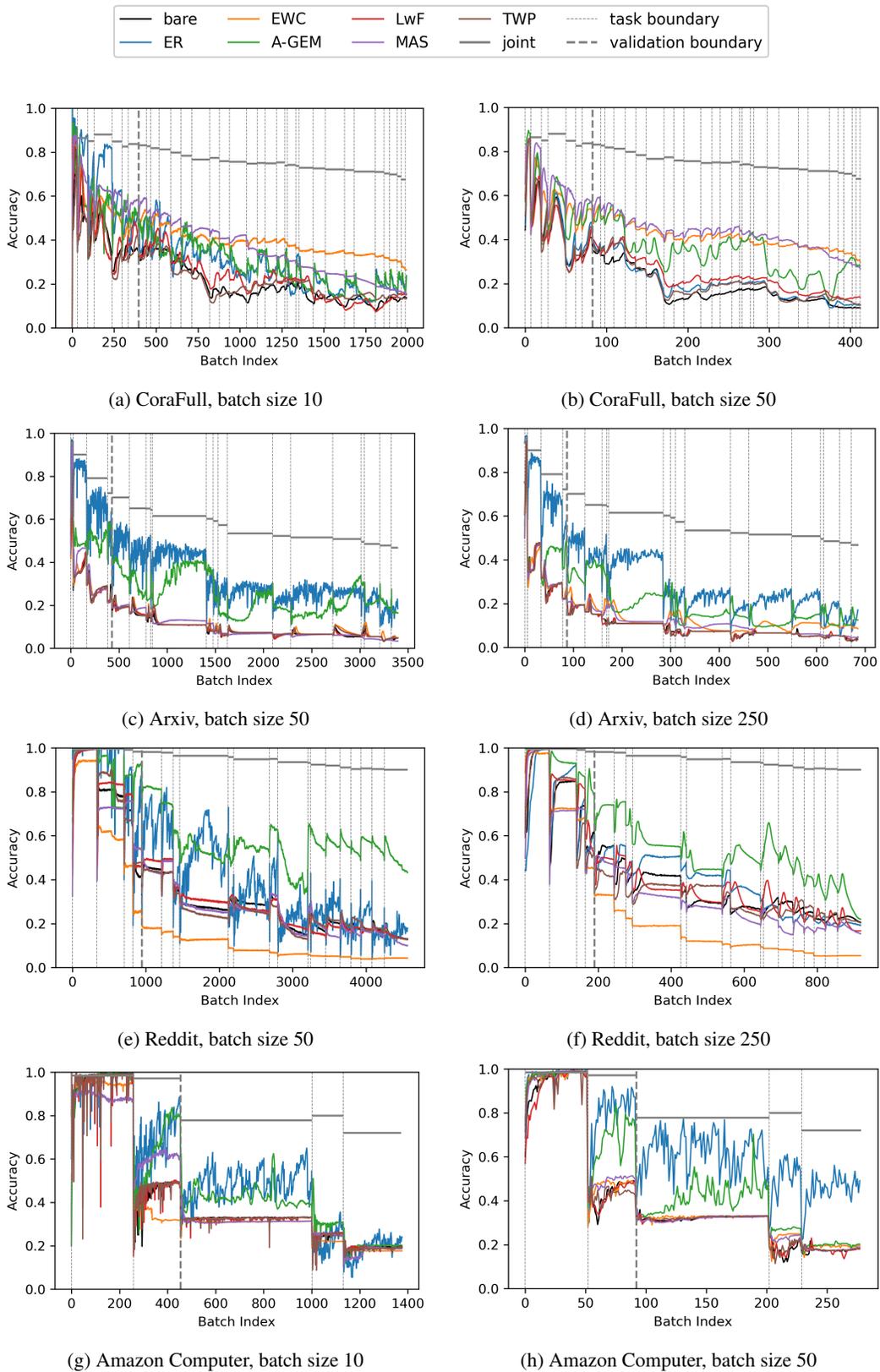


Figure 4: Anytime evaluation, showing AA on validation nodes after training on each mini-batch, when performing neighborhood sampling. We highlight the boundaries between tasks and the threshold up to which hyperparameter selection is performed.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

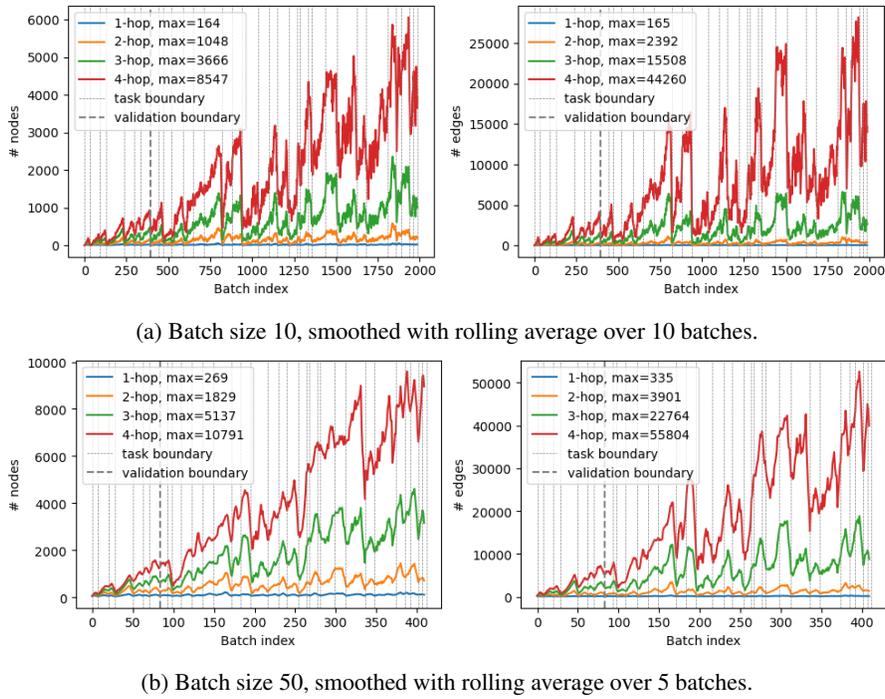


Figure 5: Neighborhood Expansion on CoraFull dataset (total nodes: 19,793). On the left: number of nodes in the l-hop neighborhood of the training batch; on the right: number of edges connecting the l-1 to l-hop neighborhood of the training batch.

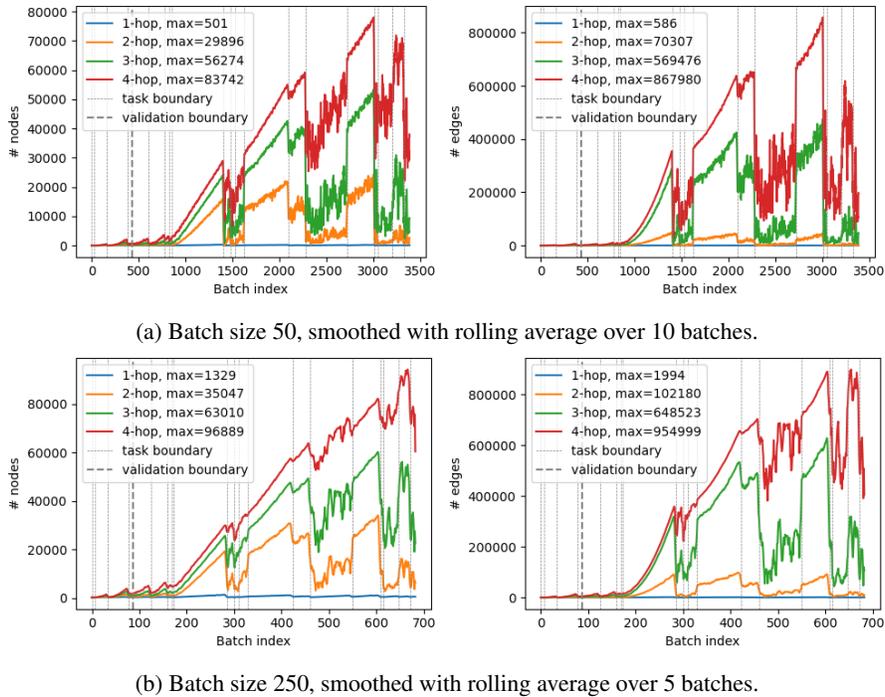
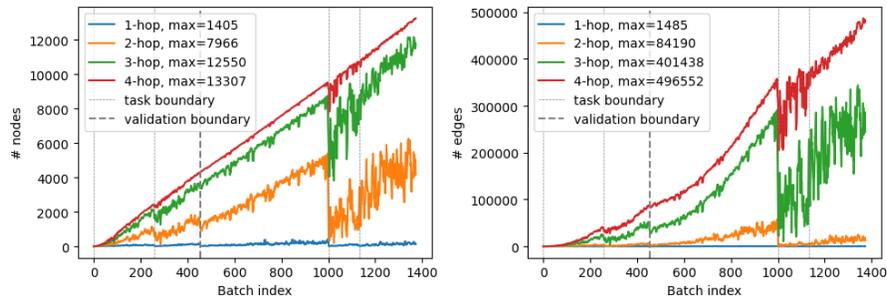
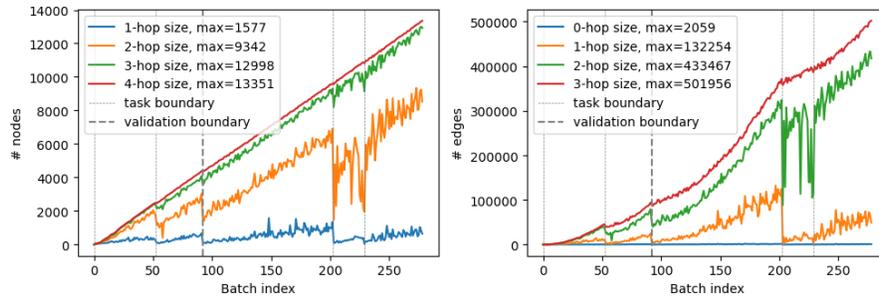


Figure 6: Neighborhood Expansion on Arxiv dataset (total nodes: 169,343). On the left: number of nodes in the l-hop neighborhood of the training batch; on the right: number of edges connecting the l-1 to l-hop neighborhood of the training batch.

1026  
 1027  
 1028  
 1029  
 1030  
 1031  
 1032  
 1033  
 1034  
 1035  
 1036  
 1037  
 1038  
 1039  
 1040  
 1041  
 1042  
 1043  
 1044  
 1045  
 1046  
 1047  
 1048  
 1049  
 1050  
 1051  
 1052  
 1053  
 1054  
 1055  
 1056  
 1057  
 1058  
 1059  
 1060  
 1061  
 1062  
 1063  
 1064  
 1065  
 1066  
 1067  
 1068  
 1069  
 1070  
 1071  
 1072  
 1073  
 1074  
 1075  
 1076  
 1077  
 1078  
 1079

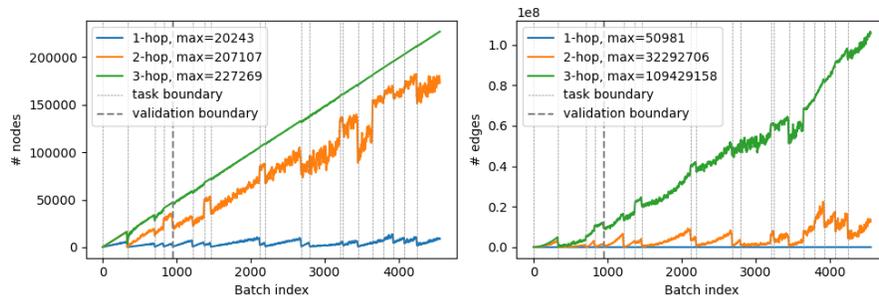


(a) Batch size 10, smoothed with rolling average over 5 batches.

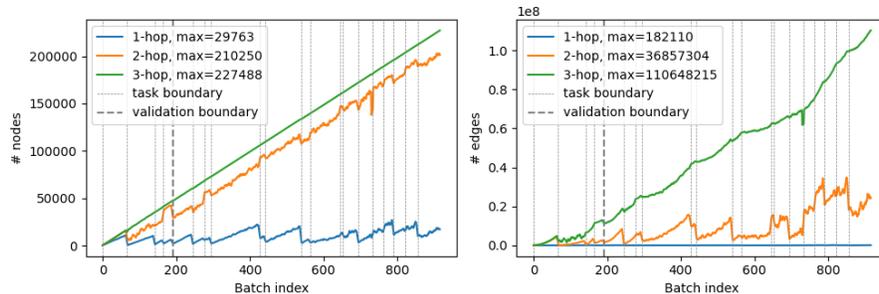


(b) Batch size 50 (no smoothing).

Figure 7: Neighborhood Expansion on Amazon Computer dataset (total nodes: 13,752). On the left: number of nodes in the l-hop neighborhood of the training batch; on the right: number of edges connecting the l-1 to l-hop neighborhood of the training batch.



(a) Batch size 50, smoothed with rolling average over 10 batches.



(b) Batch size 250, smoothed with rolling average over 5 batches.

Figure 8: Neighborhood Expansion on Reddit dataset (total nodes: 227,853). On the left: number of nodes in the l-hop neighborhood of the training batch; on the right: number of edges connecting the l-1 to l-hop neighborhood of the training batch.

## F ABLATION STUDY ON CORAFULL

In our experiments, the backbone model for all Continual Learning strategies was fixed as a 2-layer GCN with 256 hidden units (with the exception of the Reddit dataset where with full neighborhood only one layer was used), following Zhang et al. (2022). Thus, we performed a small ablation study on the CoraFull dataset to assess the impact of the number of hidden units and number of GCN layers. This is conducted only with full neighborhood aggregation.

### F.1 NUMBER OF HIDDEN UNITS

Still maintaining a 2-layer GCN, we changed the number of hidden units, from the original 256 to alternatively 128 and 512. We report the results in Tables 11-12, and we show anytime evaluation in Figure 9. We observe how overall performance are lower than with 256 hidden units, possibly as this number may have been validated in previous works. A smaller network seems to consistently damage performance metrics, while with 512 units A-GEM scores similarly as before, and EWC actually improves greatly, reaching 40% AA. If the increase in computational complexity with additional units is not significant, it may thus be worthwhile to consider validating the number of hidden units as well.

Table 11: Performance comparison on CoraFull with 512 hidden units.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	18.16 $\pm$ 0.57	26.30 $\pm$ 0.53	-45.41 $\pm$ 0.98	9.94 $\pm$ 0.40	25.30 $\pm$ 0.30	-35.66 $\pm$ 1.58
ER	27.09 $\pm$ 2.16	35.20 $\pm$ 0.29	-64.72 $\pm$ 2.34	10.28 $\pm$ 0.51	25.49 $\pm$ 0.41	-36.63 $\pm$ 1.47
EWC	38.56 $\pm$ 1.80	48.70 $\pm$ 1.01	-22.13 $\pm$ 2.76	40.76 $\pm$ 1.05	52.15 $\pm$ 0.80	-19.52 $\pm$ 2.10
A-GEM	34.57 $\pm$ 2.03	37.73 $\pm$ 0.35	-56.61 $\pm$ 2.19	28.59 $\pm$ 1.92	42.51 $\pm$ 0.38	-16.95 $\pm$ 1.35
LWF	19.61 $\pm$ 0.84	26.41 $\pm$ 0.71	-43.19 $\pm$ 1.26	12.93 $\pm$ 1.42	27.73 $\pm$ 0.26	-36.52 $\pm$ 1.35
MAS	3.10 $\pm$ 0.31	26.61 $\pm$ 1.38	-27.12 $\pm$ 1.50	11.52 $\pm$ 1.62	37.41 $\pm$ 1.42	-23.65 $\pm$ 1.91
TWP	18.18 $\pm$ 0.57	26.30 $\pm$ 0.53	-45.40 $\pm$ 0.99	9.90 $\pm$ 0.41	25.29 $\pm$ 0.29	-35.69 $\pm$ 1.48

Table 12: Performance comparison on CoraFull with 128 hidden units.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	15.12 $\pm$ 2.27	22.75 $\pm$ 1.64	-35.08 $\pm$ 2.38	10.74 $\pm$ 0.97	24.55 $\pm$ 1.10	-25.33 $\pm$ 1.83
ER	15.05 $\pm$ 2.38	29.55 $\pm$ 0.40	-66.55 $\pm$ 1.51	13.96 $\pm$ 1.18	32.48 $\pm$ 0.33	-70.49 $\pm$ 1.33
EWC	5.73 $\pm$ 1.47	21.50 $\pm$ 3.44	-6.50 $\pm$ 2.53	5.90 $\pm$ 1.73	28.90 $\pm$ 2.14	-20.11 $\pm$ 4.39
A-GEM	15.94 $\pm$ 3.23	28.41 $\pm$ 0.65	-71.24 $\pm$ 2.38	26.62 $\pm$ 1.48	39.37 $\pm$ 0.47	-34.64 $\pm$ 1.77
LWF	15.30 $\pm$ 1.61	26.03 $\pm$ 0.28	-69.40 $\pm$ 0.69	10.06 $\pm$ 2.76	24.78 $\pm$ 0.85	-29.47 $\pm$ 2.63
MAS	6.25 $\pm$ 2.74	28.81 $\pm$ 1.89	-24.70 $\pm$ 2.94	10.85 $\pm$ 3.74	33.50 $\pm$ 1.97	-31.71 $\pm$ 3.37
TWP	15.12 $\pm$ 2.27	22.75 $\pm$ 1.64	-35.10 $\pm$ 2.35	10.70 $\pm$ 0.96	24.56 $\pm$ 1.10	-25.33 $\pm$ 1.81

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

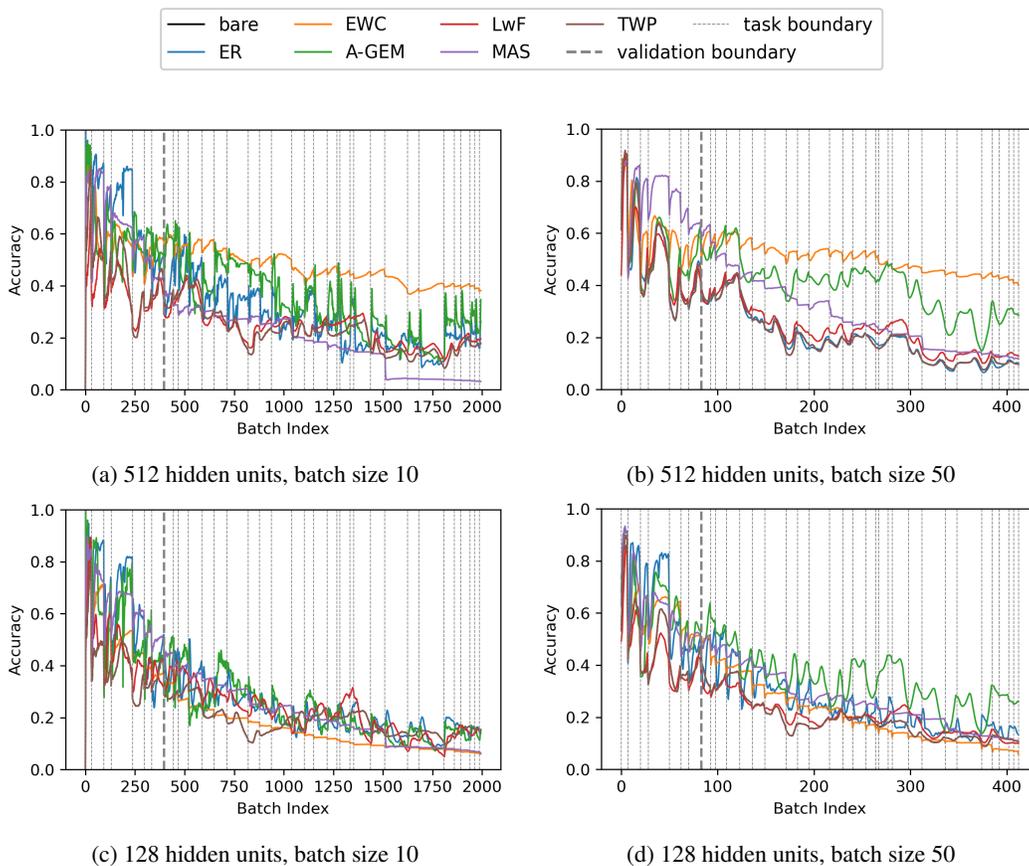


Figure 9: Anytime evaluation on CoraFull, showing AA on validation nodes after training on each mini-batch. We highlight the boundaries between tasks and the threshold up to which hyperparameter selection is performed.

## F.2 NUMBER OF GCN LAYERS

We considered here a different number of GCN layers compared to our main results, 1 and 3 specifically, as on CoraFull also the 3-hop neighborhood has a relatively limited expansion. Results are reported in Tables 13-14, and anytime evaluation plots in Figure 10. We observe overall lower results specifically with 3 layers, while ER improves using only 1 GCN layer.

Additionally, as a sort of “0-layer” baseline, we test the usage of a 2 layer MLP, equivalent to our main backbone GCN applied using the identity as adjacency matrix. This is the same as discarding graph topology information. While most results are consequently greatly reduced, ER, GEM and EWC are able to maintain some performance even without using the graph topology.

Table 13: Performance comparison on CoraFull with 1 GCN layer.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	18.39 $\pm$ 0.06	24.02 $\pm$ 0.02	-76.21 $\pm$ 0.08	15.00 $\pm$ 0.10	22.94 $\pm$ 0.15	-27.09 $\pm$ 0.19
ER	38.68 $\pm$ 0.67	57.97 $\pm$ 0.40	-53.69 $\pm$ 0.73	33.38 $\pm$ 0.83	55.66 $\pm$ 0.18	-60.26 $\pm$ 0.72
EWC	18.40 $\pm$ 0.11	24.01 $\pm$ 0.03	-76.14 $\pm$ 0.08	15.01 $\pm$ 0.15	22.87 $\pm$ 0.11	-27.30 $\pm$ 0.22
A-GEM	28.90 $\pm$ 1.15	54.16 $\pm$ 0.26	-64.51 $\pm$ 1.42	29.21 $\pm$ 0.86	52.70 $\pm$ 0.11	-64.92 $\pm$ 1.00
LWF	24.09 $\pm$ 0.34	28.77 $\pm$ 0.10	-64.59 $\pm$ 0.41	15.08 $\pm$ 0.13	22.88 $\pm$ 0.19	-27.26 $\pm$ 0.19
MAS	34.49 $\pm$ 0.08	31.93 $\pm$ 0.03	-58.54 $\pm$ 0.08	15.01 $\pm$ 0.15	22.87 $\pm$ 0.11	-27.30 $\pm$ 0.22
TWP	17.01 $\pm$ 0.14	23.36 $\pm$ 0.01	-76.52 $\pm$ 0.11	15.01 $\pm$ 0.05	22.85 $\pm$ 0.13	-26.97 $\pm$ 0.34

Table 14: Performance comparison on CoraFull with 3 GCN layers.

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	6.66 $\pm$ 1.83	18.24 $\pm$ 0.13	-83.82 $\pm$ 1.84	5.99 $\pm$ 0.49	17.70 $\pm$ 0.82	-52.80 $\pm$ 1.36
ER	9.32 $\pm$ 1.41	24.08 $\pm$ 0.56	-77.20 $\pm$ 1.76	2.74 $\pm$ 0.21	16.40 $\pm$ 1.04	-32.44 $\pm$ 14.08
EWC	18.20 $\pm$ 2.25	38.86 $\pm$ 2.95	-25.36 $\pm$ 2.71	20.89 $\pm$ 3.20	43.46 $\pm$ 2.09	-23.17 $\pm$ 4.05
A-GEM	3.85 $\pm$ 0.13	19.97 $\pm$ 0.62	-61.79 $\pm$ 11.32	23.91 $\pm$ 3.51	34.53 $\pm$ 0.90	-29.78 $\pm$ 3.63
LWF	16.18 $\pm$ 1.45	30.47 $\pm$ 1.40	-53.27 $\pm$ 2.98	15.69 $\pm$ 1.90	28.96 $\pm$ 1.26	-42.63 $\pm$ 4.37
MAS	13.57 $\pm$ 1.83	33.62 $\pm$ 2.02	-21.95 $\pm$ 2.53	28.36 $\pm$ 1.24	49.34 $\pm$ 1.51	-11.68 $\pm$ 1.47
TWP	12.62 $\pm$ 2.58	21.91 $\pm$ 0.73	-58.41 $\pm$ 3.34	7.86 $\pm$ 3.14	18.87 $\pm$ 1.25	-17.47 $\pm$ 4.27

Table 15: Performance comparison on CoraFull with MLP (no graph structure).

METHOD	BATCH SIZE 10			BATCH SIZE 50		
	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$	AA% $\uparrow$	AAA <sub>val</sub> % $\uparrow$	AF% $\uparrow$
BARE	2.90 $\pm$ 1.07	16.14 $\pm$ 1.17	-36.09 $\pm$ 2.86	2.62 $\pm$ 0.58	13.22 $\pm$ 0.24	-20.90 $\pm$ 0.53
ER	25.25 $\pm$ 0.63	33.18 $\pm$ 0.34	-67.41 $\pm$ 0.51	13.61 $\pm$ 1.15	31.74 $\pm$ 0.12	-76.10 $\pm$ 1.10
EWC	18.28 $\pm$ 3.78	30.83 $\pm$ 2.21	-38.09 $\pm$ 4.47	2.22 $\pm$ 0.91	18.50 $\pm$ 0.63	-45.75 $\pm$ 2.93
GEM	26.56 $\pm$ 1.70	38.95 $\pm$ 0.11	-45.83 $\pm$ 1.85	3.10 $\pm$ 0.22	16.80 $\pm$ 0.17	-20.85 $\pm$ 0.49
LWF	5.65 $\pm$ 1.09	18.34 $\pm$ 0.68	-35.84 $\pm$ 2.24	3.53 $\pm$ 0.48	14.69 $\pm$ 0.33	-22.45 $\pm$ 1.18
MAS	2.69 $\pm$ 0.20	20.41 $\pm$ 0.36	-46.12 $\pm$ 0.73	10.78 $\pm$ 2.29	26.87 $\pm$ 0.58	-42.59 $\pm$ 2.54
TWP	4.98 $\pm$ 1.28	17.42 $\pm$ 0.92	-37.57 $\pm$ 4.05	2.57 $\pm$ 0.25	13.03 $\pm$ 0.39	-21.42 $\pm$ 0.97

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

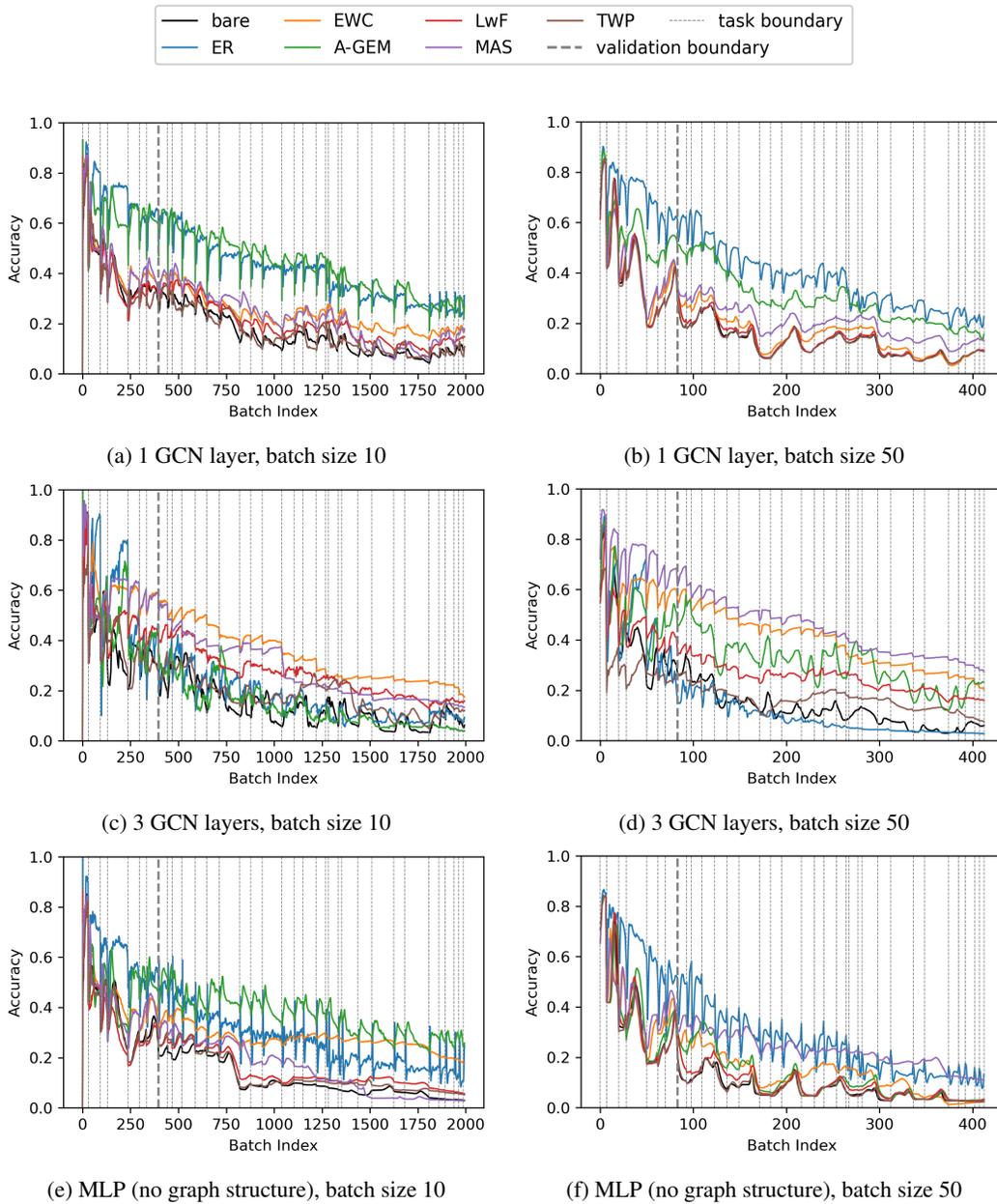


Figure 10: Anytime evaluation on CoraFull, showing AA on validation nodes after training on each mini-batch. We highlight the boundaries between tasks and the threshold up to which hyperparameter selection is performed.