

# GAP: Differentially Private Graph Neural Networks with Aggregation Perturbation

Sina Sajadmanesh<sup>1,2</sup> Ali Shahin Shamsabadi<sup>3</sup> Aurélien Bellet<sup>4</sup> Daniel Gatica-Perez<sup>1,2</sup>  
<sup>1</sup>*Idiap Research Institute* <sup>2</sup>*EPFL* <sup>3</sup>*The Alan Turing Institute* <sup>4</sup>*Inria*

## Abstract

Graph Neural Networks (GNNs) are powerful models designed for graph data that learn node representation by recursively aggregating information from each node’s local neighborhood. However, despite their state-of-the-art performance in predictive graph-based applications, recent studies have shown that GNNs can raise significant privacy concerns when graph data contain sensitive information. As a result, in this paper, we study the problem of learning GNNs with Differential Privacy (DP). We propose GAP, a novel differentially private GNN that safeguards the privacy of nodes and edges using aggregation perturbation, i.e., adding calibrated stochastic noise to the output of the GNN’s aggregation function, which statistically obfuscates the presence of a single edge (edge-level privacy) or a single node and all its adjacent edges (node-level privacy). To circumvent the accumulation of privacy cost at every forward pass of the model, we tailor the GNN architecture to the specifics of private learning. In particular, we first precompute private aggregations by recursively applying neighborhood aggregation and perturbing the output of each aggregation step. Then, we privately train a deep neural network on the resulting perturbed aggregations for any node-wise classification task. A major advantage of GAP over previous approaches is that we guarantee edge-level and node-level DP not only for training, but also at inference time with no additional costs beyond the training’s privacy budget. We theoretically analyze the formal privacy guarantees of GAP using Rényi DP. Empirical experiments conducted over three real-world graph datasets demonstrate that GAP achieves a favorable privacy-accuracy trade-off and significantly outperforms existing approaches.

## 1 Introduction

Real-world datasets are often represented by graphs, such as social [45], financial [54], transportation [11], or biological [28] networks, modeling the relations (i.e., edges) between a collection of entities (i.e., nodes). Graph Neural Networks (GNNs)

have achieved state-of-the-art performance in learning over such relational data in various graph-based machine learning tasks, such as node classification, link prediction, and graph classification [9, 20, 29, 61, 68, 69]. Due to their superior performance, GNNs are now widely used in many applications such as recommendation systems, credit issuing, traffic forecasting, drug discovery, and medical diagnosis [3, 7, 17, 27, 37, 63].

**Privacy concerns.** However, real-world deployments of GNNs raise privacy concerns when graphs contain users’ data: for instance, social or financial networks involve sensitive information about individuals and their interactions. Recent works [24, 25, 41, 57] have extended the study of the privacy leakage of standard deep learning models to GNNs, showing that the risk of leaking information about the training data is even higher in GNNs, as they incorporate not only node features and labels but also the graph structure itself [12]. Consequently, GNNs are vulnerable to various privacy attacks, such as node membership inference [25, 41] and edge stealing [24, 57]. For example, a GNN trained over a social network for friendship recommendation could reveal the existing friendship links between the users through its predictions. As another example, a GNN trained over the social graph of COVID-19 patients to predict the spread of the disease could be used as a service by government authorities, but an adversary might be able to recover the attributes and social interactions of individual patients that were used for training.

**Problem and motivation.** Motivated by these privacy concerns, in this paper, we investigate the problem of designing algorithms for training GNNs on private, sensitive graphs. Our goal is to protect the sensitive graph structure and other accompanying data using the framework of Differential Privacy (DP) [13], the gold standard for formalizing privacy guarantees. In the context of graph data, two different variants of DP have been defined: *edge-level* and *node-level* DP [46]. Informally, an edge-level  $\epsilon$ -DP algorithm should have roughly the same output (as measured by  $\epsilon$ ) if one edge is removed from the input graph. This ensures that the algorithm’s output does not reveal the existence of a particular edge in the graph.

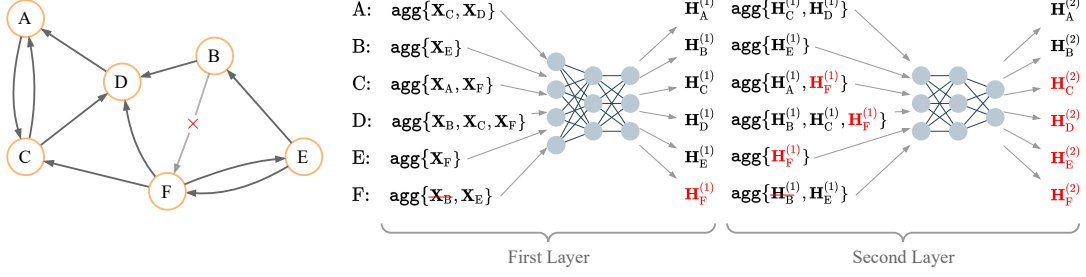


Figure 1: Simplified representation of an unfolded 2-layer GNN taking an example graph as input. At each layer, every node aggregates its neighbors’ embedding vectors (initially node features, e.g.  $\mathbf{X}_A$  for node A). This aggregated vector is then transformed using a neural network into a new embedding vector (e.g.,  $\mathbf{H}_A$ ). Removing an arbitrary edge (here, the edge from node B to F) excludes the source node (B) from the aggregation set of the destination node (F). At the first layer, this will only alter the destination node’s embedding, but due to the existing graph connectivities, this change is propagated to the other nodes in the next layer. Node embeddings that are affected by the removal of edge (B,F) are indicated in red.

Correspondingly, node-level private algorithms conceal the existence of a particular node in the graph together with all its associated edges and attributes. Clearly, node-level DP is a stronger privacy definition, but it is much harder to attain because it requires the output distribution of the algorithm to hide much larger differences in the input graph.

**Challenges.** The problem of learning GNNs with differential privacy has unique challenges, mainly due to the relational nature of graph data and the unique way it is processed through GNNs, which is quite different from standard deep learning models like Convolutional Neural Networks (CNNs). In fact, GNNs generalize classical CNNs to graphs by aggregating node features from the local neighborhood of every node in the graph. Specifically, a  $K$ -layer GNN, comprising  $K$  stacked graph convolution layers, iteratively learns lower-dimensional node representations by aggregating and utilizing information from every node’s  $K$ -hop neighborhood (i.e., from nodes that are at distance at most  $K$  in the graph). As a result, the representation of each node is influenced by all the nodes in its  $K$ -hop neighborhood (Figure 1). This fact voids the privacy guarantee of standard differentially private learning algorithms, such as DP-SGD [2], as the training loss of GNNs can no longer be decomposed into individual samples.

Another major challenge in designing differentially private GNNs is to guarantee *inference privacy*, i.e., preserving the privacy of graph data not only during training, but also after training and during the inference time, when the GNN model is queried to make predictions for test nodes. When learning standard deep learning models with DP, it is sufficient to make the training algorithm private as training data is not re-used at inference time, thus the inference mechanism is independent of the training dataset. On the contrary, in the case of GNNs, training and test nodes are not always disjoint. Specifically, GNNs can be trained in the *transductive setting*, where the model is trained on a subset of labeled nodes and then tested on the rest of the nodes within *the same graph*.

Considering a social network as an example, a GNN could be trained on a graph where a small number of users are known to be either humans or bots (training nodes). After training, the GNN is used to predict the label (human/bot) for the unlabeled users (test nodes) in the same social network. Since the same graph is used for prediction, even with a privately trained GNN model, the private node features and adjacency information can still be leaked. Note that even in the *inductive learning setting*, where the training and test graphs are completely disjoint, the interactions between the test nodes can still be private, and without a privacy-preserving inference mechanism, this private information could be leaked through the GNN’s predictions. As a result, it is critical to ensure that both the training and inference stages of a GNN satisfy differential privacy. This is illustrated in Figure 2.

**Prior work.** A few recent works have attempted to learn GNNs with DP [10, 40, 47, 57], but they either consider restrictive assumptions or their methods incur high privacy costs to achieve reasonable accuracy, which is often unacceptable for end users. Wu *et al.* [57] propose an edge-level DP GNN learning algorithm by directly perturbing the adjacency matrix of the graph. However, due to the huge dimensions of the adjacency matrix and the excessive noise injected, their method only work well for relatively large privacy budgets ( $\epsilon \geq 10$ ). Olatunji *et al.* [40] extend the framework of PATE [42] to GNNs and propose a node-level private GNN learning algorithm, but their method requires additional public graph data, which is not very realistic as public graphs are hardly available for many tasks. Daigavane *et al.* [10] also propose a node-level private approach by extending the standard DP-SGD algorithm to bounded-degree graph data. However, their method fails to provide inference privacy and is restricted to 1-layer GNNs and thus cannot benefit from higher-order aggregations. Both of these node-level private algorithms [10, 40] require a very high privacy budget ( $\epsilon \geq 30$ ) to perform reasonably well.

**Our contributions.** In this paper, we propose GAP, a GNN

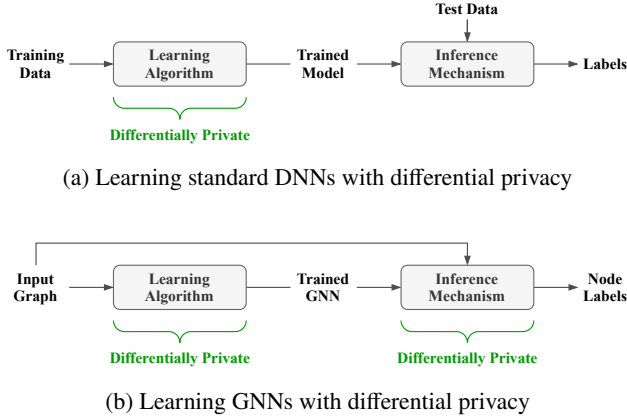


Figure 2: Comparison of differentially private deep learning with (a) conventional deep neural networks, and (b) graph neural networks. With DNNs, the training and test data are not related, so it is sufficient to only make the training procedure differentially private. With GNNs however, the training and test nodes can both belong to the same graph, so both training and inference steps need to be differentially private.

learning framework that satisfies edge-level privacy for both training and inference. Motivated by the observation that perturbing an edge in the input graph can practically be viewed as changing a sample in a node’s neighborhood aggregation function, we propose *aggregation perturbation*, where we use the Gaussian mechanism to add calibrated noise to the output of the GNN’s aggregation function. As this process happens inside the forward propagation of the model, both training and inference stages satisfy edge-level DP. To avoid the accumulation of privacy cost at every forward pass of the model, instead of stacking multiple graph convolution layers as usual, we decouple the GNN’s graph convolutions and learnable transformations into an aggregation module and a classification module. This decomposition allows us to benefit from higher-order, multi-hop aggregations of graph data while significantly reducing the privacy cost as the perturbed aggregations are computed once and reused throughout training (and inference). Moreover, our method does not incur additional privacy cost at inference time beyond the privacy cost of training, neither in the transductive nor the inductive learning setting. More importantly, our proposed method has the capability to be extended to node-level privacy setting for bounded-degree graphs by combining it with standard private learning algorithms, such as DP-SGD.

We evaluate GAP on three large-scale graph datasets and compare its accuracy-privacy performance with existing approaches. We demonstrate that GAP can achieve near-optimal accuracy with much lower privacy budgets under both edge-level DP ( $\epsilon \geq 5$ ) and node-level DP ( $\epsilon \geq 15$ ). Furthermore, it always performs better than a (private) neural network model which does not utilize the graph’s adjacency information. In

summary, our main contributions are as follows:

- We formalize the problem of learning and deploying GNNs with differential privacy in [Section 3](#) and discuss the key differences between standard deep learning models and graph neural networks when it comes to learning with DP.
- We propose a novel GNN learning method with DP guarantees in [Section 4](#), which is based on aggregation perturbation and relies on a custom architecture tailored to the requirements of private learning. Unlike existing methods, our approach ensures both edge-level and node-level DP, does not rely on public data, is not limited to 1-layer aggregations, and guarantees inference privacy for both transductive and inductive learning settings without additional cost.
- To empirically demonstrate the superior privacy-accuracy performance of GAP, we conduct extensive experiments in [Section 5](#). Our method allows us to establish new state-of-the-art trade-offs between privacy and accuracy on three graph datasets, namely Facebook, Reddit and Amazon.

## 2 Related Work

**Graph neural networks.** Deep learning on graphs has emerged in the past few years to tackle different kinds of graph-based learning tasks. A variety of GNN models and various architectures have been proposed, including Graph Convolutional Networks [29], Graph Attention Networks [51], GraphSAGE [20], Graph Isomorphism Networks [61], Jumping Knowledge Networks [62], and Gated Graph Neural Networks [32]. For the latest advances and trends in GNNs, we refer the reader to the available surveys [1, 21, 59, 71, 74].

**Privacy attacks on GNNs.** Several recent works have investigated the possibility of performing privacy attacks against GNNs and quantified the privacy leakage of publicly released GNN models or node embeddings trained on private graph datasets. These attacks mainly fall in three categories: attribute inference, membership inference, and link stealing/graph reconstruction. Duddu *et al.* [12] conduct a comprehensive study on quantifying the privacy leakage of graph embedding algorithms trained on sensitive graph data and propose various privacy attacks on GNNs under practical threat models and adversary assumptions. Zhang *et al.* [73] also study the information leakage in graph embeddings and propose three different inference attacks against GNNs: inferring graph properties (such as number of nodes and edges), inferring whether a given subgraph is contained in the target graph, and graph reconstruction with similar statistics to the target graph. He *et al.* [24] propose a series of link stealing attacks on GNN models, to which the adversary has black-box access. They show that an adversary can accurately infer a link between any pair of nodes in a graph used to train the GNN. Zhang *et al.* [72] study the connection between model inversion risk and edge

influence, and show that edges with greater influence are more likely to be inferred. Wu *et al.* [57] also study the link stealing attack via influence analysis and propose an effective attack against GNNs based on the node influence information. The feasibility of the membership inference attack against GNNs has also been studied and several attacks with different threat models have been proposed in the literature [4, 12, 25, 41]. Overall, these works underline the privacy risks of GNNs trained on sensitive graph data and confirm the vulnerability of these models to various privacy attacks.

**Privacy-preserving GNNs.** There have been attempts to mitigate the above-mentioned attacks by proposing privacy-preserving GNN models. Li *et al.* [31] presents a graph adversarial training framework that integrates disentangling and purging mechanisms to remove users’ private attributes from learned node representations. Liao *et al.* [33] also follow an adversarial learning approach to address the attribute inference attack on GNNs, where they introduce a mini-max game between the desired graph feature encoder and the worst-case attacker. Wang *et al.* [52] propose a privacy-preserving GNN learning framework from the mutual information perspective, where they learn node representations, such that the primary learning task achieves high performance, while the privacy protection task performs arbitrarily bad. Hsieh *et al.* [26] also propose an adversarial defense mechanism against attribute inference attacks on GNNs by maintaining the accuracy of target label classification and reducing the accuracy of private label classification. There is also increasing trend in using federated and split learning to address the privacy issues of GNNs in distributed learning settings [5, 6, 23, 34, 35, 39, 44, 53, 56, 60, 66, 67]. However, none of the aforementioned works employ the notion of DP, and thus they do not provide provable privacy guarantees.

**Differentially private GNNs.** Very recently, differential privacy has also been used to provide formal privacy guarantees in various GNN learning settings. Sajadmanesh and Gatica-Perez [47] propose a locally private GNN model by considering a distributed learning setting, where node features and labels are private but training the GNN is federated by a central server with access to graph edges. They use local DP to protect node data and design a learning algorithm utilizing the public graph structure for feature and label denoising. However, their method cannot be used in applications where the graph edges are private. Wu *et al.* [57] propose an edge-level DP learning algorithm for GNNs by following the input graph perturbation approach. Their method directly injects Laplace noise into the adjacency matrix of the graph, and the GNN is trained over the resulting noisy graph. However, their method cannot be extended trivially to the node-level privacy setting. Olatunji *et al.* [40] consider a centralized learning setting and propose a node-level private GNN learning algorithm by adapting the framework of PATE [42] to GNNs. They train the student GNN model using public graph data, which

is privately labeled using the teacher GNN models trained exclusively for each query node. However, their dependence on public graph data restricts the applicability of their method. Daigavane *et al.* [10] also propose a node-level private approach for training 1-layer GNNs by extending the standard DP-SGD algorithm and privacy amplification by subsampling results to bounded-degree graph data. However, their approach fails to provide inference privacy. Moreover, it is limited to 1-layer models and thus cannot leverage the information from higher-order neighborhoods, which is crucial for GNNs to achieve high accuracy in many applications.

**Comparison with existing methods.** To our best knowledge, the proposed GAP method is the first approach that can be used under either edge-level or node-level privacy settings depending on the specifics of the application. Unlike existing methods, our approach does not rely on public data and can leverage multi-hop aggregations beyond first-order neighbors. Furthermore, GAP guarantees inference privacy at no additional cost. In Section 5, we also show that GAP outperforms other baselines in terms of accuracy-privacy performance.

### 3 Background and Problem Formulation

#### 3.1 Graph Neural Networks

GNNs aim to learn a representation for every node in the input graph by incorporating the initial node features and the graph structure (edges). The learned node representations, or embeddings, can then be used for the downstream machine learning task. In this paper, we focus on node classification, where the embeddings are used to predict the label of the graph nodes. Node-wise prediction problems can be tackled in either a transductive or an inductive setting. In the transductive setting, both training and testing are performed on the same graph, but different nodes are used for training and testing. Conversely, in the inductive setting, training and testing are performed on different graphs.

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{Y})$  be an unweighted directed graph dataset consisting of sets of nodes  $\mathcal{V}$  and edges  $\mathcal{E}$  represented by a binary adjacency matrix  $\mathbf{A} \in \{0, 1\}^{N \times N}$ , where  $N = |\mathcal{V}|$  denotes the number of nodes, and  $\mathbf{A}_{i,j} = 1$  if there is a directed edge  $(i, j) \in \mathcal{E}$  from node  $i$  to node  $j$ . Nodes are characterized by  $d$ -dimensional feature vectors stacked up in an  $N \times d$  matrix  $\mathbf{X}$ , where  $\mathbf{X}_v$  denotes the feature vector of the  $v$ -th node.  $\mathbf{Y} \in \{0, 1\}^{N \times C}$  represents the labels of the nodes, where  $\mathbf{Y}_v$  is a  $C$ -dimensional one-hot vector denoting the label of the  $v$ -th node, and  $C$  is the number of classes. Note that in the transductive learning setting, only a subset  $\mathcal{V}_l \subset \mathcal{V}$  of the nodes are labeled, and thus  $\mathbf{Y}_v$  is a zero vector for all  $v \notin \mathcal{V}_l$ .

A typical  $K$ -layer GNN consists of  $K$  sequential graph convolution layers. The  $i$ -th layer receives node embeddings from  $(i - 1)$ -th layer and outputs a new embedding for each node by aggregating the previous layer’s embeddings of its adjacent



neighbors followed by applying a learnable transformation, as defined below:

$$\mathbf{M}_v^{(i)} = \text{agg} \left( \{ \mathbf{H}_u^{(i-1)} : \forall u \in \mathcal{N}_v \} \right),$$

$$\mathbf{H}_v^{(i)} = \text{upd} \left( \mathbf{M}_v^{(i)} ; \boldsymbol{\Theta}^{(i)} \right),$$

where  $\mathcal{N}_v = \{u : \mathbf{A}_{u,v} \neq 0\}$  denotes the set of adjacent nodes to node  $v$  (i.e., nodes with outbound edges toward  $v$ ), and  $\mathbf{H}_u^{(i-1)}$  is the embedding of an adjacent node  $u$  at layer  $i-1$ . The aggregation function,  $\text{agg}(\cdot)$ , is a (sub)differentiable, permutation invariant aggregator function, such as SUM, MEAN, or MAX, which outputs an aggregated vector  $\mathbf{M}_v^{(i)}$  for each node  $v$ . Finally,  $\text{upd}(\cdot)$  is a learnable function, such as a multi-layer perceptron (MLP), parameterized by  $\boldsymbol{\Theta}^{(i)}$  that takes the aggregated vector and outputs the new embedding  $\mathbf{H}_v^{(i)}$ . For convenience, we use the matrix notation by stacking the vectors of all the nodes into a matrix, and accordingly, we define the matrix-based version of  $\text{agg}(\cdot)$  and  $\text{upd}(\cdot)$  as:

$$\text{AGG}(\mathbf{H}, \mathbf{A}) = [\text{agg}(\{ \mathbf{H}_u : \forall u \in \mathcal{N}_v \}) : \forall v \in \mathcal{V}]^T,$$

$$\text{UPD}(\mathbf{M}; \boldsymbol{\Theta}) = [\text{upd}(\mathbf{M}_v; \boldsymbol{\Theta}) : \forall v \in \mathcal{V}]^T,$$

where we omitted the layer indicator superscripts for simplicity. Initially, we have  $\mathbf{H}^{(0)} = \mathbf{X}$  as the input to the GNN's first layer, i.e., the input embedding of a node  $v$  is its feature vector  $\mathbf{X}_v$ . The last layer generates an output embedding vector for each node, which can be used in different ways depending on the downstream task. For node classification, a softmax layer is applied to the final embeddings  $\mathbf{H}^{(K)}$  to obtain the posterior class probabilities  $\hat{\mathbf{Y}}$ . The illustration of a typical 3-layer GNN is depicted in Figure 3.

### 3.2 Differential Privacy

Differential privacy (DP) [14] is the de facto standard for formalizing the privacy guarantees of (randomized) algorithms that process sensitive data. Informally, DP requires that, for every individual, the output distribution of the algorithm be roughly the same whether or not this individual's data is present in the dataset. Therefore, an adversary having access to the data of all but the target individual cannot distinguish (with high probability) whether the target's record is among the input data. The formal definition of DP is as follows.

**Definition 1** (Differential Privacy [14]). *Given  $\epsilon > 0$  and  $\delta > 0$ , a randomized algorithm  $\mathcal{A}$  satisfies  $(\epsilon, \delta)$ -differential privacy, if for all possible pairs of adjacent datasets  $X$  and  $X'$  differing by at most one record, denoted as  $X \sim X'$ , and for any possible set of outputs  $S \subseteq \text{Range}(\mathcal{A})$ , we have:*

$$\Pr[\mathcal{A}(X) \in S] \leq e^\epsilon \Pr[\mathcal{A}(X') \in S] + \delta.$$

Here, the parameter  $\epsilon$  is called the privacy budget (or privacy cost) and is used to tune the privacy-utility trade-off of the

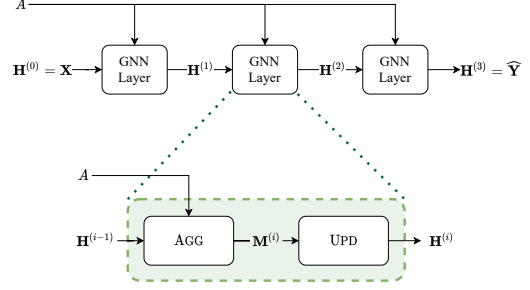


Figure 3: Typical 3-layer GNN for node classification. Each layer  $i$  takes the adjacency matrix  $\mathbf{A}$  along with the node embedding matrix  $\mathbf{H}^{(i-1)}$  from the previous layer, and outputs a new embedding matrix  $\mathbf{H}^{(i)}$ . Internally, the input embeddings  $\mathbf{H}^{(i-1)}$  are aggregated based on the adjacency information in  $\mathbf{A}$ , and the resulting aggregations,  $\mathbf{M}^{(i)}$ , are fed into a neural network (UPD) to generate new embeddings  $\mathbf{H}^{(i)}$ . The initial embedding is equal to the node features  $\mathbf{X}$  and the final embedding is the predicted class labels  $\hat{\mathbf{Y}}$ .

algorithm: a lower privacy budget leads to stronger privacy guarantees but reduced utility. The parameter  $\delta$  is treated as a failure probability, and is usually chosen to be very small with respect to the inverse of the number of data records. DP has the following important properties that help us design complex algorithms from simpler ones [13]:

- *Robustness to post-processing:* Any post-processing of the output of an  $(\epsilon, \delta)$ -DP algorithm remains  $(\epsilon, \delta)$ -DP.
- *Sequential composition:* If an  $(\epsilon, \delta)$ -DP algorithm is applied  $k$  times on the same data, the result is at least  $(k\epsilon, k\delta)$ -DP.
- *Parallel composition:* Executing an  $(\epsilon, \delta)$ -DP algorithm on disjoint chunks of data yields an  $(\epsilon, \delta)$ -DP algorithm.

In this paper, we use an alternative definition of DP, called *Rényi Differential Privacy* (RDP) [36], as it allows to obtain tighter sequential composition results.

**Definition 2** (Rényi Differential Privacy [36]). *A randomized algorithm  $\mathcal{A}$  is  $(\alpha, \epsilon)$ -RDP for  $\alpha > 1, \epsilon > 0$  if for every adjacent datasets  $X \sim X'$ , we have:*

$$D_\alpha(\mathcal{A}(X) \parallel \mathcal{A}(X')) \leq \epsilon,$$

where  $D_\alpha(P \parallel Q)$  is the Rényi divergence of order  $\alpha$  between probability distributions  $P$  and  $Q$  defined as:

$$D_\alpha(P \parallel Q) = \frac{1}{\alpha - 1} \log \mathbb{E}_{x \sim Q} \left[ \frac{P(x)}{Q(x)} \right]^\alpha.$$

As RDP is a generalization of DP, it can be easily converted back to standard  $(\epsilon, \delta)$ -DP using the following proposition.

**Proposition 1** (From RDP to  $(\epsilon, \delta)$ -DP [36]). *If  $\mathcal{A}$  is an  $(\alpha, \epsilon)$ -RDP algorithm, then it also satisfies  $(\epsilon + \log(1/\delta)/\alpha - 1, \delta)$ -DP for any  $\delta \in (0, 1)$ .*

A basic mechanism to achieve RDP is the *Gaussian mechanism*, where we add Gaussian noise to the output of the algorithm we want to make private. Specifically, let  $f : \mathcal{X} \rightarrow \mathbb{R}^d$  be the non-private algorithm taking a dataset as input and outputting a  $d$ -dimensional vector. We first define the *sensitivity* of  $f$  as the maximum  $L_2$  distance achievable when applying  $f(\cdot)$  to adjacent datasets  $X$  and  $X'$ :

$$\Delta_f = \max_{X \sim X'} \|f(X) - f(X')\|_2.$$

Then, adding Gaussian noise with variance  $\sigma^2$  to  $f$  as:

$$\mathcal{A}(X) = f(X) + \mathcal{N}(\sigma^2 \mathbb{I}_d),$$

where  $\mathbb{I}_d$  is  $d \times d$  identity matrix, yields an  $(\alpha, \epsilon)$ -RDP algorithm for all  $\alpha > 1$  with  $\epsilon = \frac{\Delta_f^2 \alpha}{2\sigma^2}$  [36].

### 3.3 Problem Definition

Let  $\hat{\mathbf{Y}} = \mathcal{F}(\mathbf{X}, \mathbf{A}; \Theta)$  be a GNN-based node classification model with parameter set  $\Theta$  that takes node features  $\mathbf{X}$  and the graph's adjacency matrix  $\mathbf{A}$  as input, and outputs the corresponding predicted labels  $\hat{\mathbf{Y}}$ . To learn the model parameters  $\Theta$ , we minimize a standard classification loss function (e.g., cross-entropy) with respect to  $\Theta$  as follows:

$$\Theta^* = \arg \min_{\Theta} \sum_{v \in \mathcal{V}_l} \ell(\hat{\mathbf{Y}}_v, \mathbf{Y}_v), \quad (1)$$

where  $\ell(\cdot, \cdot)$  is the loss function,  $\mathbf{Y}$  is the ground-truth labels, and  $\mathcal{V}_l \subseteq \mathcal{V}$  is the set of labeled nodes. After training, in the transductive setting, the learned GNN is used to infer the labels of the unlabeled nodes in  $\mathcal{G}$ :

$$\hat{\mathbf{Y}} = \mathcal{F}(\mathbf{X}, \mathbf{A}; \Theta^*), \quad (2)$$

Otherwise, in the inductive setting, a new graph dataset  $\mathcal{G}_{test}$  is used as input to the learned GNN for label inference.

The goal of this paper is to preserve the privacy of graph datasets for both the training step (Eq. 1) and the inference step (Eq. 2) using differential privacy. Note that preserving privacy in the inference step is critical as the adjacency information is still used in this step for obtaining the predicted labels.

However, as graph datasets are different from standard tabular datasets due to the existence of links between data records, one needs to adapt the definition of DP to graphs. As the semantic interpretation of DP relies on the definition of adjacent datasets, we first define two different notions of adjacency in graphs, namely edge-level and node-level adjacent graph datasets [22].

**Definition 3** (Edge-level adjacent graphs). *Two graphs  $\mathcal{G}$  and  $\mathcal{G}'$  are edge-level adjacent if one can be obtained by removing a single edge from the other. Therefore,  $\mathcal{G}$  and  $\mathcal{G}'$  differ by at most one edge.*

**Definition 4** (Node-level adjacent graphs). *Two graphs  $\mathcal{G}$  and  $\mathcal{G}'$  are node-level adjacent if one can be obtained by removing a single node (with its features, labels, and all attached edges) from the other. Therefore,  $\mathcal{G}$  and  $\mathcal{G}'$  differ by at most one node.*

Accordingly, the definition of edge-level and node-level DP follows from the above definitions: an algorithm  $\mathcal{A}$  is edge-level (respectively, node-level)  $(\epsilon, \delta)$ -DP if for every two edge-level (respectively, node-level) adjacent graph datasets  $\mathcal{G}$  and  $\mathcal{G}'$  and any set of outputs  $S \subseteq \text{Range}(\mathcal{A})$ , we have:

$$\Pr[\mathcal{A}(\mathcal{G}) \in S] \leq e^\epsilon \Pr[\mathcal{A}(\mathcal{G}') \in S] + \delta.$$

Intuitively, edge-level DP protects edges (which could represent connections between people), while node-level DP protects nodes together with their adjacent edges (i.e., all information pertaining to an individual, including features, labels, and connections).

## 4 Proposed Method: GAP

In this section, we present our proposed differentially private GNN with Aggregation Perturbation (GAP). First, we provide an overview of GAP and then describe each of its individual components in details. Next, we analyze GAP's privacy guarantees at training and inference under both edge-level and node-level DP. We conclude this section with a discussion of GAP's architectural choices and practical limitations.

### 4.1 Overview

Figure 4 exhibits an overview of our proposed GAP method, which is composed of the following three components: (i) a graph-agnostic *Encoder Module* that is pre-trained on initial node features and labels to extract low-dimensional node representation without using the private graph structure; (ii) a parameter-independent *Aggregation Module* that privately aggregates extracted low-dimensional node embeddings at multiple hops; and (iii) a learnable *Classification Module* that is trained over private aggregations to predict node labels.

The fundamental mechanism of GAP for preserving the privacy of graph data, which is part of the aggregation module, is the *aggregation perturbation* approach: we use the Gaussian mechanism to add stochastic perturbations to the output of the GNN's aggregation function proportional to its sensitivity. This approach is motivated by the fact that adding or removing an edge to/from the input graph can practically be viewed as adding or removing a sample to/from the neighborhood aggregation function of the edge's destination node. Therefore, we can effectively hide the presence of a single edge (corresponding to edge-level privacy) or a group of edges (corresponding to node-level privacy) by adding an appropriate amount of noise to the aggregation function.

Even though aggregation perturbation can be applied on standard sequential GNN models, such as the one shown in

Figure 3, the total privacy cost of such a solution would be too high to be practical. The reason is that with sequential GNNs, the aggregation function of all layers are constantly queried at every training iteration as their inputs depend on model parameters (except the aggregation of the very first layer, which only depends on initial node features). Therefore, all the aggregations need to be perturbed after every update to ensure DP. With a  $K$ -layer GNN and  $T$  training iterations, this results in an excessive accumulated privacy cost of  $\mathcal{O}(\sqrt{KT})$  as the Gaussian mechanism is composed  $KT$  times (at each forward pass and each layer).

To circumvent this issue, we propose to decouple graph-based aggregations (the aggregation module) from learnable transformations (the classification module). This decomposition, which is similar in spirit to the Inception architecture and scalable networks [16, 48, 49, 58, 70], significantly reduces the privacy cost of the model as the aggregations no longer depend on the model parameters and thus can be precomputed and perturbed only once and reused throughout training and inference. More specifically, in GAP, the Gaussian mechanism is composed only  $K$  times instead of  $KT$ . Since  $K$  is usually small ( $1 \leq K \leq 5$ ) compared to  $T$  (in the order of hundreds), this leads to a substantial reduction in the privacy budget compared to a sequential GNN model.

As shown in Figure 4, the graph’s adjacency matrix  $\mathbf{A}$  is only queried inside the aggregation module, and the subsequent classification module only post-processes the aggregated data without further accessing the graph structure. Therefore, due to the robustness of DP to post-processing, the aggregation perturbation approach ensures that the forward propagation (and by extension, backward propagation) of our GAP model preserves the privacy of graph’s adjacency information, which is sufficient to guarantee edge-level privacy at both training and inference stages. To ensure node-level privacy, however, in addition to aggregation perturbation, the encoder and classification modules will also need to be trained privately as they access private node features and labels (we will discuss this in Section 4.6).

Finally, as most graph datasets usually come with high dimensional and often sparse node features, it is essential to transform them into a low-dimensional space before computing the private aggregations, since the utility of the perturbed data typically degrades as the dimensionality of the node features increases. To this end, prior to learning the classification module, we propose to perform a graph-independent pre-training step to learn the encoder module using only node features and labels. The learned encoder is then used to extract low-dimensional features for the aggregation module. Altogether, the summary of GAP’s end-to-end pipeline is as follows:

- **Step 1 (Encoder module):** Pre-train the encoder over node features and labels and then extract low-dimensional node embeddings using the encoder.

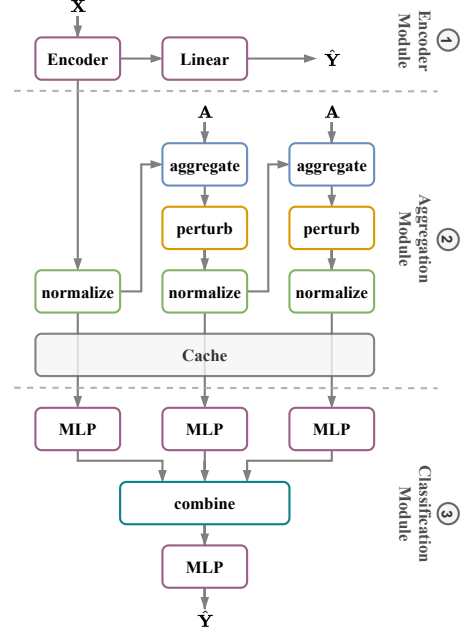


Figure 4: Overview of our GAP model. First, the encoder is pre-trained using only node features ( $\mathbf{X}$ ) and labels ( $\mathbf{Y}$ ). Then, the encoder’s extracted features are given to the aggregation module to compute private  $K$ -hop aggregations ( $K = 2$  in this case) using the graph’s adjacency matrix ( $\mathbf{A}$ ). This step is done only once and the result is cached. Finally, the aggregations are fed into the classification module for label prediction.

- **Step 2 (Aggregation module):** Compute private multi-hop aggregations from low-dimensional node embeddings.
- **Step 3 (Classification module):** Train the classification module over private aggregated features and labels.

## 4.2 Encoder Module

GAP uses a multi-layer perceptron (MLP) model as an encoder to transform the original node features into a low-dimensional representation suitable for the aggregation module. As we need the aggregation module to remain independent of the model parameters, we cannot train the encoder end-to-end with the classification module. Instead, we pre-train a simple MLP on the node features and labels, and then take the penultimate layer as the encoder to transform node features into a low-dimensional representation. More specifically, the model is composed of an encoder MLP followed by a linear layer with softmax activation function as follows:

$$\hat{\mathbf{Y}} = \text{softmax}(\text{MLP}_{\text{enc}}(\mathbf{X}; \Theta_{\text{enc}}) \cdot \mathbf{W}), \quad (3)$$

where  $\text{MLP}_{\text{enc}}$  is an encoder MLP with parameter set  $\Theta_{\text{enc}}$ ,  $\mathbf{W}$  is the transformation matrix of the final linear layer,  $\mathbf{X}$  is the original node features, and  $\hat{\mathbf{Y}}$  is the corresponding posterior class probabilities. In order to train this model, we minimize the

cross-entropy (or any other classification-related) loss function  $\ell(\cdot, \cdot)$  with respect to the model parameters  $\Theta = \{\Theta_{\text{enc}}, \mathbf{W}\}$ :

$$\Theta^* = \arg \min_{\Theta} \sum_{v \in \mathcal{V}_l} \ell(\hat{\mathbf{Y}}_v, \mathbf{Y}_v), \quad (4)$$

where  $\mathbf{Y}$  is the ground-truth labels and  $\mathcal{V}_l \subseteq \mathcal{V}$  is the set of labeled nodes. Note that the graph's adjacency matrix is *not* used at this step. After pre-training, we detach the encoder MLP and use it to extract low-dimensional node features,  $\mathbf{X}^{(0)}$ , for the aggregation module:

$$\mathbf{X}^{(0)} = \text{MLP}_{\text{enc}}(\mathbf{X}; \Theta_{\text{enc}}^*). \quad (5)$$

### 4.3 Aggregation Module

The goal of aggregation module is to privately release node feature aggregations under different hops in the graph using the aggregation perturbation method. [Algorithm 1](#) presents our approach, the Private Multi-hop Aggregation (PMA) mechanism. We build our mechanism upon the SUM aggregation function, which is simply equivalent to the multiplication of the adjacency matrix  $\mathbf{A}$  by the input feature matrix  $\mathbf{X}$ , as:

$$\text{AGG}(\mathbf{X}, \mathbf{A}) = \mathbf{A}^T \cdot \mathbf{X}. \quad (6)$$

The PMA mechanism takes the initial node embeddings  $\tilde{\mathbf{X}}^{(0)}$ , which is the row-normalized version of the encoder's extracted features  $\mathbf{X}^{(0)}$  as:

$$\tilde{\mathbf{X}}_v^{(0)} = \mathbf{X}_v^{(0)} / \|\mathbf{X}_v^{(0)}\|_2, \quad \forall v \in \mathcal{V}. \quad (7)$$

It then outputs a set of  $K$  private aggregated node features  $\tilde{\mathbf{X}}^{(1)}$  to  $\tilde{\mathbf{X}}^{(K)}$  corresponding to different hops from 1 to  $K$ . Specifically, given  $\sigma > 0$ , the PMA mechanism performs the following steps to recursively compute and perturb the aggregations in  $k$ -th hop from  $(k-1)$ -th:

1. **Aggregation:** First, we compute  $k$ -th non-private aggregations using the normalized aggregations at step  $k-1$ :

$$\mathbf{X}^{(k)} = \mathbf{A}^T \cdot \tilde{\mathbf{X}}^{(k-1)}. \quad (8)$$

2. **Perturbation:** Next, we perturb the aggregations by adding noise with variance  $\sigma^2$  using the Gaussian mechanism:

$$\tilde{\mathbf{X}}^{(k)} = \mathbf{X}^{(k)} + \mathcal{N}(\sigma^2 \mathbb{I}). \quad (9)$$

3. **Normalization:** Finally, as it is essential to bound the effect of each feature vector on the subsequent aggregations, we again row-normalize the private aggregated node features, such that the L2-norm of each row is 1:

$$\tilde{\mathbf{X}}_v^{(k)} = \tilde{\mathbf{X}}_v^{(k)} / \|\tilde{\mathbf{X}}_v^{(k)}\|_2, \quad \forall v \in \mathcal{V}. \quad (10)$$

---

#### Algorithm 1: Private Multi-hop Aggregation

---

**Input** : Graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with adjacency matrix  $\mathbf{A}$ ; initial row-normalized feature matrix  $\tilde{\mathbf{X}}^{(0)}$ ; max hop  $K$ ; noise standard deviation  $\sigma$ ;  
**Output** : Private aggregated node feature matrices  $\tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(K)}$

```

1 for  $k \in \{1, \dots, K\}$  do
2    $\mathbf{X}^{(k)} \leftarrow \mathbf{A}^T \cdot \tilde{\mathbf{X}}^{(k-1)}$  // aggregate
3    $\tilde{\mathbf{X}}^{(k)} \leftarrow \mathbf{X}^{(k)} + \mathcal{N}(\sigma^2 \mathbb{I})$  // perturb
4   for  $v \in \mathcal{V}$  do
5      $\tilde{\mathbf{X}}_v^{(k)} \leftarrow \tilde{\mathbf{X}}_v^{(k)} / \|\tilde{\mathbf{X}}_v^{(k)}\|_2$  // normalize
6   end
7 end
8 return  $\tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(K)}$ 

```

---

### 4.4 Classification Module

The classification module takes as input the list of aggregated features  $\tilde{\mathbf{X}}^{(0)}, \tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(K)}$  under different hops provided by the aggregation module. To obtain the  $k$ -th layer representation  $\mathbf{H}^{(k)}$ , we apply a transformer MLP, denoted as  $\text{MLP}_{\text{pre}}^{(k)}$ :

$$\mathbf{H}^{(k)} = \text{MLP}_{\text{pre}}^{(k)}(\tilde{\mathbf{X}}^{(k)}; \Theta_{\text{pre}}^{(k)}), \quad \forall k \in \{0, 1, \dots, K\}, \quad (11)$$

where  $\Theta_{\text{pre}}^{(k)}$  is the learnable parameters of  $\text{MLP}_{\text{pre}}^{(k)}$ . Next, we combine these multi-hop representations to get an integrated node embedding matrix  $\mathbf{H}$ :

$$\mathbf{H} = \text{COMBINE}(\{\mathbf{H}^{(0)}, \mathbf{H}^{(1)}, \dots, \mathbf{H}^{(K)}\}; \Theta_{\text{comb}}), \quad (12)$$

where COMBINE is any differentiable combination strategy, with common choices being summation, concatenation, or attention, potentially with parameter set  $\Theta_{\text{comb}}$ . Finally, we feed the integrated representation into a classifier MLP, denoted as  $\text{MLP}_{\text{post}}$ , to get posterior class probabilities for the nodes:

$$\hat{\mathbf{Y}} = \text{MLP}_{\text{post}}(\mathbf{H}; \Theta_{\text{post}}), \quad (13)$$

where  $\Theta_{\text{post}}$  represents the parameters of  $\text{MLP}_{\text{post}}$ . To train the classification module, we minimize a similar loss function as [Eq. 4](#) but with respect to the parameters of the classification module:  $\Theta = \{\Theta_{\text{pre}}^{(0)}, \dots, \Theta_{\text{pre}}^{(K)}, \Theta_{\text{comb}}, \Theta_{\text{post}}\}$ . The overall training procedure of the GAP model is presented in [Algorithm 2](#).

### 4.5 Edge-Level Privacy Analysis

In the following, we provide a formal analysis of GAP's edge-level privacy guarantees at training and inference stages.

**Training privacy.** The following arguments establish the DP guarantees of the PMA mechanism and the GAP training algorithm. The detailed proofs can be found in [Appendix A](#).

**Theorem 1.** *Given the maximum hop  $K \geq 1$  and noise standard deviation  $\sigma > 0$ , the PMA mechanism presented in [Algorithm 1](#) satisfies edge-level  $(\alpha, K\alpha/2\sigma^2)$ -RDP for any  $\alpha > 1$ .*



---

**Algorithm 2:** Training the GAP Framework

---

**Input** : Graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  with adjacency matrix  $\mathbf{A}$ ; node feature matrix  $\mathbf{X}$ ; node labels  $\mathbf{Y}$ ; max hop  $K$ ; noise standard deviation  $\sigma$ ;

**Output** : Trained model parameters  
 $\{\Theta_{\text{enc}}^*, \Theta_{\text{pre}}^{*(0)}, \dots, \Theta_{\text{pre}}^{*(K)}, \Theta_{\text{comb}}^*, \Theta_{\text{post}}^*\}$

- 1 Pre-train the encoder module (Eq. 3) to obtain  $\Theta_{\text{enc}}^*$ .
  - 2 Use the pretrained encoder (Eq. 5) to obtain initial node embedding matrix  $\mathbf{X}^{(0)}$ .
  - 3 Row-normalize the initial node embedding matrix (Eq. 7) to obtain  $\tilde{\mathbf{X}}^{(0)}$ .
  - 4 Apply the PMA mechanism (Algorithm 1) to obtain private aggregations  $\tilde{\mathbf{X}}^{(1)}, \dots, \tilde{\mathbf{X}}^{(K)}$ .
  - 5 Train the classification module (Eq. 11 to Eq. 13) to obtain  $\Theta_{\text{pre}}^{*(0)}, \dots, \Theta_{\text{pre}}^{*(K)}, \Theta_{\text{comb}}^*$ , and  $\Theta_{\text{post}}^*$ .
  - 6 **return**  $\Theta_{\text{enc}}^*, \Theta_{\text{pre}}^{*(0)}, \dots, \Theta_{\text{pre}}^{*(K)}, \Theta_{\text{comb}}^*, \Theta_{\text{post}}^*$
- 

**Proposition 2.** For any  $\delta \in (0, 1)$ , maximum hop  $K \geq 1$ , and noise standard deviation  $\sigma > 0$ , Algorithm 2 satisfies edge-level  $(\epsilon, \delta)$ -DP with  $\epsilon = \frac{K}{2\sigma^2} + \sqrt{2K \log(1/\delta)}/\sigma$ .

The privacy guarantees in Proposition 2 show that the privacy cost grows with the number of layers ( $K$ ), but is independent from the number of training steps thanks to our GAP architecture.

**Inference privacy.** A major advantage of our GAP framework is that querying the model at inference time does not consume additional privacy budget. This is true for both the transductive and the inductive learning settings, as discussed below:

- *Transductive setting:* In this setting, both training and inference are conducted on the same graph, but using different nodes for training and inference steps. As the entire graph is available at training time, the PMA mechanism gives us the private aggregations of both training and test nodes. Therefore, there is no need to invoke again the PMA mechanism, as the aggregations for the test nodes are already cached at training time. As a result, no additional privacy cost is incurred at inference time.
- *Inductive setting:* In this setting, the training and inference steps are conducted on different graphs, and thus the PMA mechanism needs to be applied to the test graph as well as the training graph. However, since the training and test graphs are disjoint, the application of the PMA mechanism is subject to the parallel composition of differentially private mechanisms, and therefore, the total privacy budget does not increase beyond that of training's.

## 4.6 Extension to Node-Level Privacy

As defined in Section 3, the node-level differential privacy not only protects the graph edges, but also node features and labels. Our GAP architecture is readily extensible to provide node-level privacy guarantees as well. The only additional

assumption we need is that we have bounded-degree graphs, i.e., the degree (number of inbound plus outbound edges) of each node should be bounded above by a constant  $D$ . This allows to bound the sensitivity of the aggregation function in the PMA mechanism when adding/removing a node, as in this case each node can influence and be influenced by at most  $D$  other nodes. If the input graph has nodes with very high degrees, we can use randomized neighbor sampling to reduce the degree of each node to be at most  $D$ , as proposed in [10].

For bounded-degree graphs, adding or removing a node corresponds (in the worst case) to adding or removing  $D$  edges. Therefore, our PMA mechanism also ensures node-level privacy, albeit with increased privacy costs compared to the edge-level setting (see Theorem 2 below).

However, since the node features and labels are also private under node-level DP, the encoder and classification modules need to be trained privately as they access node features/labels. To this end, we can simply use standard DP-SGD [2] or any other differential private learning algorithm for pre-training the encoder as well as training the classification module with DP. In other words, steps 1 and 5 of Algorithm 2 must be done with DP instead of regular non-private training. This way, since each of the three GAP modules become node-level private, the entire GAP model, as an adaptive composition of several node-level private mechanisms, satisfies node-level DP. The formal node-level privacy analysis of GAP's training and inference is provided below.

**Training privacy.** The node-level privacy guarantees of the PMA mechanism and the GAP training algorithm are as follows. Detailed proofs are deferred to Appendix A.

**Theorem 2.** Given the maximum degree  $D \geq 1$ , maximum hop  $K \geq 1$ , and noise standard deviation  $\sigma > 0$ , the PMA mechanism presented in Algorithm 1 satisfies node-level  $(\alpha, DK\alpha/2\sigma^2)$ -RDP for any  $\alpha > 1$ .

**Proposition 3.** For any  $\alpha > 1$ , let encoder pre-training (Step 1 of Algorithm 2) classification module training (Step 5 of Algorithm 2) satisfy  $(\alpha, \epsilon_1(\alpha))$ -RDP and  $(\alpha, \epsilon_5(\alpha))$ -RDP, respectively. Then, for any  $0 < \delta < 1$ , maximum hop  $K \geq 1$ , maximum degree  $D \geq 1$ , and noise standard deviation  $\sigma > 0$ , Algorithm 2 satisfies node-level  $(\epsilon, \delta)$ -DP with:

$$\epsilon = \epsilon_1(\alpha) + \epsilon_5(\alpha) + \frac{DK\alpha}{2\sigma^2} + \frac{\log(1/\delta)}{\alpha - 1}.$$

Note that in Proposition 3, we cannot optimize  $\alpha$  in closed form as we do not know the precise form of  $\epsilon_1(\alpha)$  and  $\epsilon_5(\alpha)$ . However, in our experiments, we numerically optimize the choice of  $\alpha$  on a per-case basis.

**Inference privacy.** The arguments stated for edge-level inference privacy hold for node-level privacy as well. Note that in the inductive setting, the test graph should also have bounded degree for the node-level privacy guarantees to hold at inference time.

## 4.7 Discussion

**Choice of aggregation function.** In this paper, we used SUM as the default choice of aggregation function. Although other choices of aggregation functions are also possible, we empirically found that SUM is the most efficient choice to privatize, as its sensitivity does not depend on the size of the aggregation set (i.e., number of neighbors), which is itself a quantity that should be computed privately. For example, the calculation of both MEAN and GCN [29] aggregation functions depend on the node degrees, and thus requires additional privacy budget to be spent on perturbing node degrees. In any case, SUM is recognized as one of the most expressive aggregation functions in the GNN literature [9, 61].

**Normalization instead of clipping.** The PMA mechanism uses normalization instead of clipping to bound the effect of each individual feature on the SUM aggregation function. While clipping is more common in the private learning literature (see for instance gradient clipping in DP-SGD [2]), we empirically found that normalization is a better choice for aggregation perturbation: the classification module is then trained on normalized data, which tends to facilitate learning. Normalizing the node embeddings is actually commonly done in non-private GNNs to stabilize training [20, 64].

**Limitations.** As the aggregation perturbation approach used in the PMA mechanism adds random noise to the output of the aggregation function, the utility of the PMA mechanism naturally depends on the size of the node’s aggregation set, i.e., the node’s degree. Specifically, with a certain amount of noise, the more inbound neighbors a node has, the more accurate its noisy aggregated vector will be. This implies that graphs with higher average degree per node can tolerate larger noise in the aggregation function, and thus GAP can achieve a better privacy-accuracy trade-off on such graphs. Conversely, GAP’s performance will suffer if the average degree of the graph is too low, requiring higher privacy budgets to achieve acceptable accuracy. Note however that this is an expected behavior: nodes with fewer inbound neighbors are more easily influenced by a change in their neighborhood compared to nodes with higher degrees, and thus the privacy of low-degree nodes is harder to preserve than high-degree ones. This limitation is not specific to GAP: it is shared by all DP algorithms, whose performance generally suffer from lack of sufficient data.

## 5 Experiments

### 5.1 Datasets

We evaluate the proposed method on three publicly available node classification datasets, which are medium to large scale in terms of number of nodes and edges, and have been used by prior works to benchmark GNNs [8, 10, 16, 20, 40].

**Facebook** [50]. This dataset contains the anonymized Face-

Table 1: Overview of dataset statistics.

DATASET	NODES	EDGES	DEGREE	FEATURES	CLASSES
FACEBOOK	26,406	2,117,924	62	501	6
REDDIT	116,713	46,233,380	209	602	8
AMAZON	1,790,731	80,966,832	22	100	10

book social network between the students of the University of Illinois collected in September 2005. Nodes represent Facebook users and edges indicate friendship. Each node (user) has the following attributes: student/faculty status, gender, major, minor, and housing status, and the task is to predict the class year of users.

**Reddit** [20]. The Reddit dataset consist of a set of posts from the Reddit social network, where each node represents a post and an edge indicates if the same user commented on both posts. Node features are extracted based on the embedding of the post contents, and the task is to predict the community (subreddit) that a post belongs to.

**Amazon** [8]. The largest dataset used in this paper represents Amazon product co-purchasing network, where nodes represent products sold on Amazon and an edge indicates if two products are purchased together. Node features are bag-of-words vectors of the product description followed by PCA, and the task is to predict the category of the products.

We preprocess the datasets by limiting the classes to those having 1k, 10k, and 100k nodes on Facebook, Reddit, and Amazon, respectively. We then randomly split the remaining nodes into training, validation, and test sets with 75/10/15% ratios, respectively. Table 1 summarizes the statistics of the datasets after preprocessing.

### 5.2 Competing Methods

**Edge-level private methods.** The following methods are evaluated under edge-level privacy:

- *GAP-EDP*: Our proposed edge-level DP algorithm.
- *SAGE-EDP*: This is the method of Wu *et al.* [57] that uses the graph perturbation approach, with the popular GraphSAGE architecture [20] as its backbone GNN model. We perturb the graph’s adjacency matrix using the top-m-filter mechanism [38], which is a more memory-efficient version of the LAPGRAPH method [57] for perturbing the adjacency matrix without materializing the whole matrix.
- *MLP*: A simple MLP model that does not use the graph edges, and thus provides perfect edge-level privacy ( $\epsilon = 0$ ).

**Node-level private methods.** We compare the following node-level private algorithms:

- *GAP-NDP*: Our proposed node-level DP approach.
- *SAGE-NDP*: This is the method of Daigavane *et al.* [10] that adapts the standard DP-SGD method for 1-layer GNNs, with

Table 2: Test accuracy of different methods on the three datasets. The best performing method in each category — none-private, edge-level DP and node-level DP — is highlighted.

	METHOD	$\epsilon$	FACEBOOK	REDDIT	AMAZON
NONE	GAP- $\infty$	$\infty$	$80.2 \pm 0.24$	<b><math>99.4 \pm 0.01</math></b>	$91.2 \pm 0.07$
	SAGE- $\infty$	$\infty$	<b><math>83.3 \pm 0.42</math></b>	$99.1 \pm 0.02$	<b><math>92.8 \pm 0.09</math></b>
EDGE DP	GAP-EDP	5	<b><math>76.5 \pm 0.32</math></b>	<b><math>98.8 \pm 0.03</math></b>	<b><math>84.6 \pm 0.23</math></b>
	SAGE-EDP	5	$56.6 \pm 0.83$	$89.5 \pm 0.98$	$70.0 \pm 0.34$
	MLP	0	$50.3 \pm 0.39$	$82.7 \pm 0.06$	$70.7 \pm 0.30$
NODE DP	GAP-NDP	15	<b><math>64.3 \pm 0.41</math></b>	<b><math>95.1 \pm 0.08</math></b>	<b><math>78.4 \pm 0.12</math></b>
	SAGE-NDP	15	$43.9 \pm 0.91$	$74.7 \pm 1.33$	$33.8 \pm 1.07$
	MLP-DP	15	$50.2 \pm 0.71$	$81.3 \pm 0.18$	$70.8 \pm 0.24$

the same GraphSAGE architecture as its backbone model. Since this method does not inherently ensure inference privacy, as suggested by its authors, we add noise to the aggregation function based on its node-level sensitivity at test time and account for the additional privacy cost.

- *MLP-DP*: Similar to MLP, but trained with DP-SGD so as to provide node-level DP without using the graph edges.

We do not consider the approach of [40] as it requires public graph data and is thus not directly comparable to the others.

**Non-private methods.** To quantify the accuracy loss of our methods, we use the following non-private baselines ( $\epsilon = \infty$ ):

- *GAP- $\infty$* : a non-private counterpart of the GAP method, where we do not perturb the aggregations.
- *SAGE- $\infty$* : a non-private GNN model based on the GraphSAGE architecture.

### 5.3 Experimental Setup

**Model implementation details.** For our GAP models (GAP-EDP, GAP-NDP, and GAP- $\infty$ ), we set the number of  $\text{MLP}_{\text{enc}}$ ,  $\text{MLP}_{\text{pre}}$ , and  $\text{MLP}_{\text{post}}$  layers to be 2, 1, and 1, respectively. We use concatenation as the COMBINE function (Eq. 12) and tune the number of hops  $K$  in  $\{1, 2, \dots, 5\}$ . For the GraphSAGE models (SAGE-EDP, SAGE-NDP, and SAGE- $\infty$ ), we use the SUM aggregation function and tune the number of message-passing layers in  $\{1, 2, \dots, 5\}$ , except for SAGE-NDP that only supports one message-passing layer. We use a 2-layer and a 1-layer MLP as pre-processing and post-processing before and after the message-passing layers, respectively. For the MLP baselines (MLP and MLP-DP), we set the number of layers to 3. In addition, for both the GAP-NDP and SAGE-NDP methods, we use randomized neighbor sampling to bound the maximum degree  $D$  and search for the best  $D$  within  $\{10, 20, 50, 100, 200\}$  on Facebook and Amazon, and  $\{50, 100, 200, 300, 400\}$  on Reddit (ranges are chosen based on datasets’ average degree). For all methods, we set the number of hidden units to 16 (including the dimension of GAP’s encoded representation) and use the SeLU activation function [30] at every layer. Batch-normalization is used for all methods except the node-

level private ones (GAP-NDP, SAGE-NDP, and MLP-DP), for which batch-normalization is not supported. We initialize the weights using Glorot initialization [18].

**Training and evaluation details.** We train the non-private and edge-level private methods using the Adam optimizer over 100 epochs with full-sized batches. For the node-level private algorithms (GAP-NDP, SAGE-NDP, MLP-DP), we use DP-Adam [19] with maximum gradient norm set to 1, and train each model for 10 epochs with a batch size of 256, 2048, 4096 on Facebook, Reddit, and Amazon, respectively. For our GAP models (GAP- $\infty$ , GAP-EDP, and GAP-NDP), we use the same parameter setting for training both the encoder and classification modules. We train all the methods with a learning rate of 0.01 and repeat each combination of possible hyper-parameter values 10 times. We pick the best performing model based on validation accuracy, and report the average test accuracy with 95% confidence interval calculated by bootstrapping with 1000 samples.

**Privacy accounting and calibration.** Privacy budget accounting is done via the Analytical Moments Accountant [55]. We numerically calibrate the noise scale (i.e., the noise standard deviation  $\sigma$  divided by the sensitivity) of the PMA mechanism (for GAP-EDP and GAP-NDP), top-m-filter mechanism (for SAGE-EDP), DP-SGD (for GAP-NDP, SAGE-NDP, and MLP-DP) and the Gaussian mechanism (for inference privacy in SAGE-NDP) to achieve the desired  $(\epsilon, \delta)$ -DP. We report results for several values of  $\epsilon$ , while  $\delta$  is fixed to  $10^{-7}$  for Facebook (the smallest dataset) and  $10^{-8}$  for Reddit and Amazon. For both GAP-NDP and SAGE-NDP, we use the same noise scale for perturbing the gradients (in DP-SGD) and the aggregations (in the PMA and Gaussian mechanisms).

**Software and hardware.** All the models are implemented in PyTorch [43] using PyTorch-Geometric (PyG) [15]. We use the autodp library<sup>1</sup> which implements analytical moments accountant, and utilize Opacus [65] for training the node-level private models with differential privacy. Experiments are conducted on Sun Grid Engine with NVIDIA GeForce RTX 3090 and NVIDIA Tesla V100 GPUs, Intel Xeon 6238 CPUs, and 32 GB RAM. Our code is publicly available on GitHub<sup>2</sup>.

## 5.4 Experimental Results

### 5.4.1 Trade-offs between Privacy and Accuracy

We first compare the accuracy of our proposed methods against the non-private, edge-level private, and node-level private competitors. We fix the privacy budget to  $\epsilon = 15$  for the node-level private methods and  $\epsilon = 5$  for the edge-level private ones (except for MLP, which does not use the graph structure and thus achieves  $\epsilon = 0$ ). The results are presented in Table 2. We observe that in the non-private setting, the proposed GAP

<sup>1</sup><https://github.com/yuxiangw/autodp>

<sup>2</sup>Will be released upon acceptance

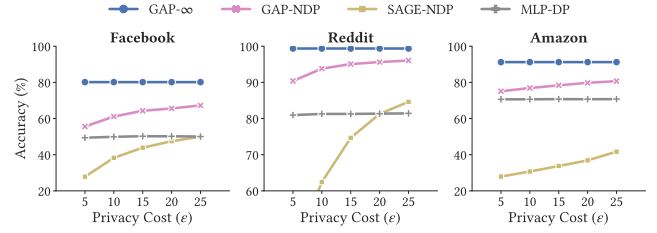
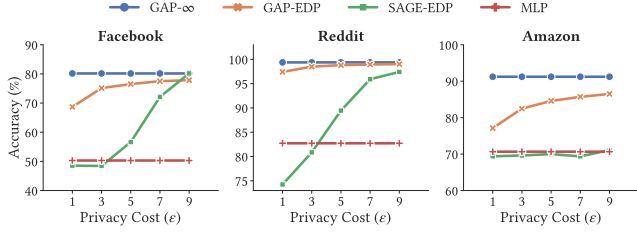


Figure 5: Accuracy vs. privacy cost ( $\epsilon$ ) of edge-level private algorithms (left) and node-level private methods (right).

architecture is competitive with SAGE, with only a slightly decrease in accuracy on Facebook and Amazon. Under both edge-level and node-level privacy settings, however, our proposed methods GAP-EDP and GAP-NDP significantly outperform their competitors. Specifically, under edge-level privacy, GAP-EDP’s accuracy is roughly 20, 10, and 14 points higher than the best competitor over Facebook, Reddit, and Amazon, respectively. Under node-level privacy, our proposed GAP-NDP method outperforms the best performing competitor by approximately 14, 14, and 8 accuracy points, respectively.

Next, to investigate how different methods perform under different privacy budgets, we vary  $\epsilon$  from 1 to 9 for edge-level private methods and from 5 to 25 for node-level private algorithms and report the accuracy of the methods under each privacy budget. The result for both edge-level and node-level privacy settings is depicted in Figure 5.

Under edge-level privacy (Figure 5, left side), we observe that our approach GAP-EDP consistently outperforms its direct competitor, SAGE-EDP, especially at lower privacy costs, and always performs significantly better than a vanilla MLP. The relative gap between GAP-EDP and SAGE-EDP is influenced by the average degree of the dataset. SAGE-EDP requires a privacy budget of  $\epsilon \geq 9$  to achieve its best accuracy and reach a similar accuracy as GAP-EDP on Facebook and Reddit (which have high and moderate average degrees respectively), and does even beat the MLP baseline on Amazon (which has low average degree). In comparison, the accuracy of GAP-EDP converges to its best at much lower privacy budgets. Specifically, GAP-EDP’s accuracy almost saturates at  $\epsilon \geq 3$  on Reddit,  $\epsilon \geq 5$  on Facebook, and  $\epsilon \geq 7$  on Amazon. The amount of accuracy loss with respect to the non-private method, GAP- $\infty$ , also depends on the average degree of the graph. For example, on Reddit at  $\epsilon = 3$ , GAP- $\infty$ ’s accuracy is less than 1 point higher than GAP-EDP’s, while on Amazon at  $\epsilon = 7$ , GAP-EDP’s accuracy fall behind GAP- $\infty$  by around 5 points. These observations are in line with our discussion of Section 4.7.

We can observe similar trends under node-level privacy (Figure 5, right side). We see that for  $\epsilon \geq 5$ , our GAP-NDP method is always much more accurate than its competitors SAGE-NDP and MLP-DP. While the accuracy of the MLP-DP method is already saturated at  $\epsilon \geq 5$ , we see that GAP-NDP’s accuracy keeps increasing with larger privacy budgets, reaching its peak at  $\epsilon \approx 15$  on Facebook,  $\epsilon \approx 10$  on Reddit, and

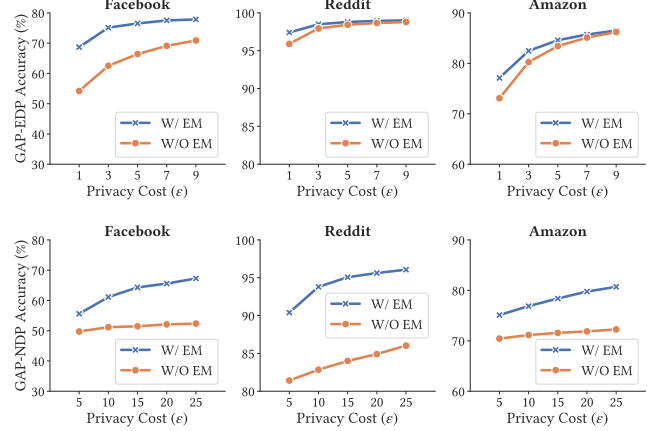


Figure 6: Effect of the encoder module (EM) on the accuracy/privacy performance of the edge-level private GAP-EDP (top) and the node-level private GAP-NDP (bottom).

$\epsilon \approx 20$  on Amazon. In the meantime, SAGE-NDP requires very high privacy budgets ( $\epsilon > 25$ ) to even beat the simple MLP-DP baseline. As expected, since the node-level private GAP-NDP hides more information (e.g., node features, labels, and all the adjacent edges to a node) than the edge-level private GAP-EDP, it requires larger privacy budgets to achieve a reasonable accuracy. Still, the accuracy loss with respect to the non-private method is higher in the node-level private method as we have further information loss due to neighborhood sampling (to bound the maximum degree of the nodes) and gradient clipping (to bound the sensitivity in DP-SGD/Adam).

#### 5.4.2 Ablation Studies

**Effectiveness of the encoder module.** In this experiment, we investigate the effect of the encoder module (EM) on the accuracy/privacy performance of the proposed methods, GAP-EDP and GAP-NDP. We compare the case in which the encoder module is used to reduce the dimensionality of the node features with the case where we remove the encoder module and just input the original node features to the aggregation module. The results under different privacy budgets are given in Figure 6. We can observe that in all cases, the accuracy of GAP-EDP and GAP-NDP is higher



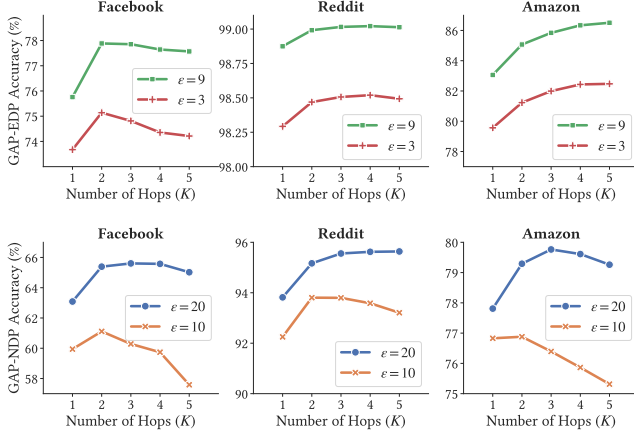


Figure 7: Effect of the number of hops  $K$  on the accuracy/privacy performance of the edge-level private GAP-EDP (top) and the node-level private GAP-NDP (bottom).

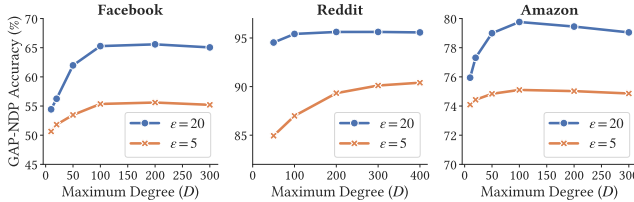


Figure 8: Effect of the maximum degree  $D$  on the accuracy/privacy performance of the node-level private GAP-NDP method with  $\epsilon = 5$  (top) and  $\epsilon = 20$  (bottom).

with EM than without it. For example, leveraging EM results in a gain of around 10, 1, and 2 accuracy points for GAP-EDP with  $\epsilon = 5$  on Facebook, Reddit, and Amazon datasets, respectively. GAP-NDP with EM also benefits from a gain of more than 15, 10, and 10 points with  $\epsilon = 15$  on Facebook, Reedit, and Amazon datasets, respectively. The effect of EM is more significant on GAP-NDP, as the amount of noise injected into the aggregations is generally larger for node-level privacy, hence dimensionality reduction becomes more critical to mitigate the impact of noise.

**Effect of the number of hops.** In this experiment, we investigate how changing the number of hops  $K$  affects the accuracy/privacy performance of our proposed methods, GAP-EDP and GAP-NDP. We vary  $K$  within  $\{1, 2, 3, 4, 5\}$  and report the accuracy under a low and a high privacy budget:  $\epsilon \in \{3, 9\}$  for GAP-EDP and  $\epsilon \in \{10, 20\}$  for GAP-NDP. The result is depicted in Figure 7. We observe that both of our methods can effectively benefit from allowing multiple hops, but there is a trade-off in increasing the number of hops. As we increment  $K$ , the accuracy of both GAP-EDP and GAP-NDP method increase up to a point and then decrease in almost all cases. The reason is that with a larger  $K$  the model is able to aggregate and utilize information from more distant nodes (all the nodes

within the  $K$ -hop neighborhood of a node) for prediction, which can increase the final accuracy. However, as more hops are involved, the amount of noise in the aggregations is also increased, which adversely affects the model’s accuracy. We can see that with the lower privacy budgets where the noise is more severe, both GAP-EDP and GAP-NDP achieve their peak accuracy at smaller  $K$  values. But as the privacy budget increases, the magnitude of the noise is reduced, enabling the models to benefit from larger  $K$  values.

**Effect of the maximum degree.** We now analyze the effect of  $D$  on the performance of GAP-NDP. We vary  $D$  in  $\{10, 20, 500, \dots, 300\}$  on Facebook and Amazon, and in  $\{50, 100, \dots, 400\}$  on Reddit, and report the accuracy under two different privacy budgets  $\epsilon \in \{5, 20\}$ . Figure 8 shows that in all cases, GAP-NDP’s accuracy increases with  $D$  up to a peak point, after which the accuracy starts to decrease slightly. This is due to the trade-off between having more samples for aggregation and the amount of noise injected: the larger  $D$ , the fewer neighbors are excluded from the aggregations (i.e., less information loss), but on the other hand, the larger the sensitivity of the aggregation function, leading to more noise injection. We also observe that the accuracy gain as a result of increasing  $D$  gets bigger as the privacy budget is increased from 5 to 20, since a higher privacy budget compensates for the higher sensitivity by reducing the amount of noise.

## 6 Conclusion

In this paper, we presented GAP, a privacy-preserving GNN framework that ensures both edge-level and node-level differential privacy at training and inference. We used aggregation perturbation, where the Gaussian mechanism is applied to the output of the GNN’s aggregation function, as a fundamental technique to achieve DP in our approach. To avoid spending privacy budget at every training iteration, we decoupled the neighborhood aggregation steps and the learnable transformations into separate aggregation and classification modules, respectively. To further reduce the privacy cost, we proposed to train an encoder module to transform node features into a low-dimensional space without relying on graph adjacency information. Experimental results over real-world graph datasets showed that our approach achieves favorable privacy/accuracy trade-offs and significantly outperforms existing methods. Promising future directions include: (i) investigating robust aggregation functions that provide specific benefits for private learning; (ii) exploiting the redundancy of information in recursive aggregations to achieve tighter composition when the number of hops  $K$  gets large, which might prove useful for specific applications; (iii) extending the framework to other tasks and scenarios, such as link-wise prediction or learning over dynamic graphs; and (iv) conducting an extended theoretical analysis of differentially private GNNs, such as proving utility bounds and characterizing their expressiveness.

## Acknowledgments

This work was supported by the European Commission’s Horizon 2020 Program ICT-48-2020, under grant number 951911, AI4Media project. It was also supported by the French National Research Agency (ANR) through grant ANR-20-CE23-0015 (Project PRIDE). Ali Shahin Shamsabadi acknowledges support from The Alan Turing Institute.

## References

- [1] Sergi Abadal, Akshay Jain, Robert Guirado, Jorge López-Alonso, and Eduard Alarcón. Computing graph neural networks: A survey from algorithms to accelerators. *ACM Computing Surveys (CSUR)*, 54(9):1–38, 2021.
- [2] Martin Abadi, Andy Chu, Ian Goodfellow, H Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, pages 308–318, 2016.
- [3] David Ahmedt-Aristizabal, Mohammad Ali Armin, Simon Denman, Clinton Fookes, and Lars Petersson. Graph-based deep learning for medical diagnosis and analysis: Past, present and future. *arXiv preprint arXiv:2105.13137*, 2021.
- [4] Wu Bang, Yang Xiangwen, Pan Shirui, and Yuan Xingliang. Adapting membership inference attacks to gnn for graph classification: Approaches and implications. In *2021 IEEE International Conference on Data Mining (ICDM)*. IEEE, 2021.
- [5] Chuan Chen, Weibo Hu, Ziyue Xu, and Zibin Zheng. Fedgl: Federated graph learning framework with global self-supervision. *arXiv preprint arXiv:2105.03170*, 2021.
- [6] Fahao Chen, Peng Li, Toshiaki Miyazaki, and Celimuge Wu. Fedgraph: Federated graph learning with intelligent sampling. *IEEE Transactions on Parallel and Distributed Systems*, 2021.
- [7] Mark Cheung and José M. F. Moura. Graph neural networks for covid-19 drug discovery. In *2020 IEEE International Conference on Big Data (Big Data)*, pages 5646–5648, 2020.
- [8] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 257–266, 2019.
- [9] Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. In *Advances in Neural Information Processing Systems*, 2020.
- [10] Ameya Daigavane, Gagan Madan, Aditya Sinha, Abhradeep Guha Thakurta, Gaurav Aggarwal, and Prateek Jain. Node-level differentially private graph neural networks. *arXiv preprint arXiv:2111.15521*, 2021.
- [11] Zulong Diao, Xin Wang, Dafang Zhang, Yingru Liu, Kun Xie, and Shaoyao He. Dynamic spatial-temporal graph convolutional neural networks for traffic forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, pages 890–897, 2019.
- [12] Vasisht Duddu, Antoine Boutet, and Virat Shejwalkar. Quantifying privacy leakage in graph embedding. In *MobiQuitous 2020 - 17th EAI International Conference on Mobile and Ubiquitous Systems: Computing, Networking and Services*, MobiQuitous ’20, page 76–85, New York, NY, USA, 2020. Association for Computing Machinery.
- [13] Cynthia Dwork. Differential privacy: A survey of results. In *International conference on theory and applications of models of computation*, pages 1–19. Springer, 2008.
- [14] Cynthia Dwork, Frank McSherry, Kobbi Nissim, and Adam Smith. Calibrating noise to sensitivity in private data analysis. In *Theory of cryptography conference*, pages 265–284. Springer, 2006.
- [15] Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- [16] Fabrizio Frasca, Emanuele Rossi, Davide Eynard, Ben Chamberlain, Michael Bronstein, and Federico Monti. Sign: Scalable inception graph neural networks. *arXiv preprint arXiv:2004.11198*, 2020.
- [17] Nikolaos Gkalelis, Andreas Goulas, Damianos Galanopoulos, and Vasileios Mezaris. Objectgraphs: Using objects and a graph convolutional network for the bottom-up recognition and explanation of events in video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3375–3383, 2021.
- [18] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256. JMLR Workshop and Conference Proceedings, 2010.

- [19] Roan Gylberth, Risman Adnan, Setiadi Yazid, and T Basaruddin. Differentially private optimization algorithms for deep neural networks. In *2017 International Conference on Advanced Computer Science and Information Systems (ICACSIS)*, pages 387–394. IEEE, 2017.
- [20] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 1025–1035, 2017.
- [21] William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584*, 2017.
- [22] Michael Hay, Chao Li, Gerome Miklau, and David Jensen. Accurate estimation of the degree distribution of private networks. In *2009 Ninth IEEE International Conference on Data Mining*, pages 169–178. IEEE, 2009.
- [23] Chaoyang He, Keshav Balasubramanian, Emir Ceyani, Carl Yang, Han Xie, Lichao Sun, Lifang He, Liangwei Yang, Philip S Yu, Yu Rong, et al. Fedgraphnn: A federated learning system and benchmark for graph neural networks. *arXiv preprint arXiv:2104.07145*, 2021.
- [24] Xinlei He, Jinyuan Jia, Michael Backes, Neil Zhenqiang Gong, and Yang Zhang. Stealing links from graph neural networks. In *30th {USENIX} Security Symposium ({USENIX} Security 21)*, 2021.
- [25] Xinlei He, Rui Wen, Yixin Wu, Michael Backes, Yun Shen, and Yang Zhang. Node-level membership inference attacks against graph neural networks. *arXiv preprint arXiv:2102.05429*, 2021.
- [26] I-Chung Hsieh and Cheng-Te Li. Netfense: Adversarial defenses against privacy attacks on neural networks for graph data. *IEEE Transactions on Knowledge and Data Engineering*, 2021.
- [27] Weiwei Jiang and Jiayun Luo. Graph neural network for traffic forecasting: A survey. *arXiv preprint arXiv:2101.11174*, 2021.
- [28] Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. Molecular graph convolutions: moving beyond fingerprints. *Journal of computer-aided molecular design*, 30(8):595–608, 2016.
- [29] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [30] Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural networks. In *Proceedings of the 31st international conference on neural information processing systems*, pages 972–981, 2017.
- [31] Kaiyang Li, Guangchun Luo, Yang Ye, Wei Li, Shihao Ji, and Zhipeng Cai. Adversarial privacy preserving graph embedding against inference attack. *arXiv preprint arXiv:2008.13072*, 2020.
- [32] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493*, 2015.
- [33] Peiyuan Liao, Han Zhao, Keyulu Xu, Tommi Jaakkola, Geoffrey J Gordon, Stefanie Jegelka, and Ruslan Salakhutdinov. Information obfuscation of graph neural networks. In *International Conference on Machine Learning*, pages 6600–6610. PMLR, 2021.
- [34] Zhiwei Liu, Liangwei Yang, Ziwei Fan, Hao Peng, and Philip S Yu. Federated social recommendation with graph neural network. *arXiv preprint arXiv:2111.10778*, 2021.
- [35] Chuizheng Meng, Sirisha Rambhatla, and Yan Liu. Cross-node federated graph neural network for spatio-temporal data modeling. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining, KDD '21*, page 1202–1211, New York, NY, USA, 2021. Association for Computing Machinery.
- [36] Ilya Mironov. Rényi differential privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*, pages 263–275. IEEE, 2017.
- [37] Alireza Mohammadshahi and James Henderson. Graph-to-graph transformer for transition-based dependency parsing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: Findings*, pages 3278–3289, Online, November 2020. Association for Computational Linguistics.
- [38] Hiep Nguyen, Abdessamad Imine, and Michaël Rusinowitch. Network structure release under differential privacy. *Transactions on Data Privacy*, 9(3):26, 2016.
- [39] Xiang Ni, Xiaolong Xu, Lingjuan Lyu, Changhua Meng, and Weiqiang Wang. A vertical federated learning framework for graph convolutional network. *arXiv preprint arXiv:2106.11593*, 2021.
- [40] Iyiola E Olatunji, Thorben Funke, and Megha Khosla. Releasing graph neural networks with differential privacy guarantees. *arXiv preprint arXiv:2109.08907*, 2021.

- [41] Iyiola E Olatunji, Wolfgang Nejdl, and Megha Khosla. Membership inference attack on graph neural networks. *arXiv preprint arXiv:2101.06570*, 2021.
- [42] Nicolas Papernot, Martín Abadi, Ulfar Erlingsson, Ian Goodfellow, and Kunal Talwar. Semi-supervised knowledge transfer for deep learning from private training data. *arXiv preprint arXiv:1610.05755*, 2016.
- [43] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [44] Yang Pei, Renxin Mao, Yang Liu, Chaoran Chen, Shifeng Xu, Feng Qiang, and Blue Elephant Tech. Decentralized federated graph neural networks. In *International Workshop on Federated and Transfer Learning for Data Sparsity and Confidentiality in Conjunction with IJCAI*, 2021.
- [45] Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. Deepinf: Social influence prediction with deep learning. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 2110–2119, 2018.
- [46] Sofya Raskhodnikova and Adam Smith. Differentially private analysis of graphs. *Encyclopedia of Algorithms*, 2016.
- [47] Sina Sajadmanesh and Daniel Gatica-Perez. Locally private graph neural networks. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 2130–2145, 2021.
- [48] Chuxiong Sun and Guoshi Wu. Scalable and adaptive graph neural networks with self-label-enhanced training. *arXiv preprint arXiv:2104.09376*, 2021.
- [49] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1–9, 2015.
- [50] Amanda L Traud, Peter J Mucha, and Mason A Porter. Social structure of facebook networks. *Physica A: Statistical Mechanics and its Applications*, 391(16):4165–4180, 2012.
- [51] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [52] Binghui Wang, Jiayi Guo, Ang Li, Yiran Chen, and Hai Li. Privacy-preserving representation learning on graphs: A mutual information perspective. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1667–1676, 2021.
- [53] Binghui Wang, Ang Li, Hai Li, and Yiran Chen. Graphfl: A federated learning framework for semi-supervised node classification on graphs. *arXiv preprint arXiv:2012.04187*, 2020.
- [54] Jianian Wang, Sheng Zhang, Yanghua Xiao, and Rui Song. A review on graph neural network methods in financial applications. *arXiv preprint arXiv:2111.15367*, 2021.
- [55] Yu-Xiang Wang, Borja Balle, and Shiva Prasad Kasiviswanathan. Subsampled renyi differential privacy and analytical moments accountant. In Kamalika Chaudhuri and Masashi Sugiyama, editors, *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pages 1226–1235. PMLR, 16–18 Apr 2019.
- [56] Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. Fedgcn: Federated graph neural network for privacy-preserving recommendation. *arXiv preprint arXiv:2102.04925*, 2021.
- [57] Fan Wu, Yunhui Long, Ce Zhang, and Bo Li. Linkteller: Recovering private edges from graph neural networks via influence analysis. *arXiv preprint arXiv:2108.06504*, 2021.
- [58] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pages 6861–6871. PMLR, 2019.
- [59] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2020.
- [60] Han Xie, Jing Ma, Li Xiong, and Carl Yang. Federated graph classification over non-iid graphs. *Advances in Neural Information Processing Systems*, 34, 2021.



- [61] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- [62] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5453–5462, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR.
- [63] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 974–983, 2018.
- [64] Jiaxuan You, Zhitao Ying, and Jure Leskovec. Design space for graph neural networks. *Advances in Neural Information Processing Systems*, 33, 2020.
- [65] Ashkan Yousefpour, Igor Shilov, Alexandre Sablayrolles, Davide Testuggine, Karthik Prasad, Mani Malek, John Nguyen, Sayan Ghosh, Akash Bharadwaj, Jessica Zhao, Graham Cormode, and Ilya Mironov. Opacus: User-friendly differential privacy library in PyTorch. *arXiv preprint arXiv:2109.12298*, 2021.
- [66] Huanding Zhang, Tao Shen, Fei Wu, Mingyang Yin, Hongxia Yang, and Chao Wu. Federated graph learning—a position paper. *arXiv preprint arXiv:2105.11099*, 2021.
- [67] Ke Zhang, Carl Yang, Xiaoxiao Li, Lichao Sun, and Siu Ming Yiu. Subgraph federated learning with missing neighbor generation. *Advances in Neural Information Processing Systems*, 34, 2021.
- [68] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in Neural Information Processing Systems*, 31:5165–5175, 2018.
- [69] Muhan Zhang, Pan Li, Yinglong Xia, Kai Wang, and Long Jin. Labeling trick: A theory of using graph neural networks for multi-node representation learning. *Advances in Neural Information Processing Systems*, 34, 2021.
- [70] Wentao Zhang, Ziqi Yin, Zeang Sheng, Wen Ouyang, Xiaosen Li, Yangyu Tao, Zhi Yang, and Bin Cui. Graph attention multi-layer perceptron. *arXiv preprint arXiv:2108.10097*, 2021.
- [71] Z. Zhang, P. Cui, and W. Zhu. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge & Data Engineering*, 34(01):249–270, jan 2022.
- [72] Zaixi Zhang, Qi Liu, Zhenya Huang, Hao Wang, Chengqiang Lu, Chuanren Liu, and Enhong Chen. Graphmi: Extracting private graph data from graph neural networks. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 3749–3755. International Joint Conferences on Artificial Intelligence Organization, 8 2021.
- [73] Zhikun Zhang, Min Chen, Michael Backes, Yun Shen, and Yang Zhang. Inference attacks against graph neural networks. In *31st USENIX Security Symposium (USENIX Security 22)*, Boston, MA, August 2022. USENIX Association.
- [74] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.

## A Deferred Theoretical Arguments

### A.1 Proof of Theorem 1

To prove Theorem 1, we first establish the following lemma.

**Lemma 1.** *Let  $\text{AGG}(\mathbf{X}, \mathbf{A}) = \mathbf{A}^T \cdot \mathbf{X}$  be the summation aggregation function. Assume that the input feature matrix  $\mathbf{X}$  is row-normalized, such that  $\forall v \in \mathcal{V} : \|\mathbf{X}_v\|_2 = 1$ . Then, the edge-level sensitivity of the aggregation function is  $\Delta_{\text{AGG}} = 1$ .*

*Proof.* Let  $\mathbf{A}$  and  $\mathbf{A}'$  be the adjacency matrices of two arbitrary edge-level adjacent graphs. Therefore, there exist two nodes  $u$  and  $v$  such that:

$$\begin{cases} \mathbf{A}'_{i,j} \neq \mathbf{A}_{i,j}, & \text{if } i = u \text{ and } j = v, \\ \mathbf{A}'_{i,j} = \mathbf{A}_{i,j}, & \text{otherwise.} \end{cases} \quad (14)$$

Without loss of generality, we can assume that  $\mathbf{A}_{v,u} = 1$  and  $\mathbf{A}'_{v,u} = 0$ . The goal is to bound the following quantity:

$$\|\text{AGG}(\mathbf{X}, \mathbf{A}) - \text{AGG}(\mathbf{X}, \mathbf{A}')\|_F.$$

Let  $\mathbf{M} = \text{AGG}(\mathbf{X}, \mathbf{A})$  be the aggregation function output on  $\mathbf{A}$ , and

$$\mathbf{M}_i = \sum_{j=1}^N \mathbf{A}_{j,i} \mathbf{X}_j,$$

be the  $i$ -th row of  $\mathbf{M}$  corresponding to the aggregated vector

for the  $i$ -th node. Analogously, let  $\mathbf{M}' = \text{AGG}(\mathbf{X}, \mathbf{A}')$ . Then:

$$\begin{aligned}
\|\text{AGG}(\mathbf{X}, \mathbf{A}) - \text{AGG}(\mathbf{X}, \mathbf{A}')\|_F &= \|\mathbf{M} - \mathbf{M}'\|_F \\
&= \left( \sum_{i=1}^N \|\mathbf{M}_i - \mathbf{M}'_i\|_2^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^N \left\| \sum_{j=1}^N (\mathbf{A}_{j,i} \mathbf{X}_j - \mathbf{A}'_{j,i} \mathbf{X}_j) \right\|_2^2 \right)^{1/2} \\
&= \left( \|\mathbf{A}_{v,u} \mathbf{X}_v - \mathbf{A}'_{v,u} \mathbf{X}_v\|_2^2 \right)^{1/2} \\
&= \|(\mathbf{A}_{v,u} - \mathbf{A}'_{v,u}) \mathbf{X}_v\|_2 \\
&= \|\mathbf{X}_v\|_2 \\
&= 1,
\end{aligned}$$

which concludes the proof.  $\square$

We can now prove [Theorem 1](#).

*Proof.* The PMA mechanism applies the Gaussian mechanism on the output of the summation aggregation function  $\text{AGG}(\mathbf{X}, \mathbf{A}) = \mathbf{A}^T \cdot \mathbf{X}$ . Based on [Lemma 1](#), the edge-level sensitivity of  $\text{AGG}(\cdot)$  is 1. Therefore, according to Corollary 3 of [36], each individual application of the Gaussian mechanism is  $(\alpha, \alpha/2\sigma^2)$ -RDP. As PMA can be seen as an adaptive composition of  $K$  such mechanisms, based on Proposition 1 of [36], the total privacy cost is  $(\alpha, K\alpha/2\sigma^2)$ -RDP.  $\square$

## A.2 Proof of Proposition 2

*Proof.* Under edge-level DP, only the adjacency information is protected. In [Algorithm 2](#), the only step where the graph's adjacency is used is the application of the PMA mechanism (step 4), which according to [Theorem 1](#) is  $(\alpha, K\alpha/2\sigma^2)$ -RDP. Since the encoder module does not use the graph's edges and the classification module only post-process the private aggregated features without accessing the edges again, the total privacy cost remains  $(\alpha, K\alpha/2\sigma^2)$ -RDP. Therefore, according to [Proposition 1](#) it is equivalent to edge-level  $(\epsilon, \delta)$ -DP with  $\epsilon = \frac{K\alpha}{2\sigma^2} + \frac{\log(1/\delta)}{\alpha-1}$ . Minimizing this expression over  $\alpha > 1$  gives  $\epsilon = \frac{K}{2\sigma^2} + \sqrt{2K \log(1/\delta)}/\sigma$ .  $\square$

## A.3 Proof of Theorem 2

We first prove the following lemma and then prove [Theorem 2](#).

**Lemma 2.** Let  $\text{AGG}(\mathbf{X}, \mathbf{A}) = \mathbf{A}^T \cdot \mathbf{X}$  be the summation aggregation function. Assume that the input feature matrix  $\mathbf{X}$  is row-normalized, such that  $\forall v \in \mathcal{V} : \|\mathbf{X}_v\|_2 = 1$ . Then, the node-level sensitivity of the aggregation function on any degree-bounded graph with maximum degree  $D$  is  $\Delta_{\text{AGG}} = \sqrt{D}$ .

*Proof.* Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{X}, \mathbf{Y})$  be an arbitrary graph dataset with  $N \times N$  adjacency matrix  $\mathbf{A}$ , feature matrix  $\mathbf{X}$ , and maximum

degree bounded above by  $D$ , i.e., for all  $v \in \mathcal{V}$ :

$$\sum_{i=1}^N \mathbf{A}_{v,i} + \sum_{j=1}^N \mathbf{A}_{j,v} - \mathbf{A}_{v,v} \leq D. \quad (15)$$

Let  $\mathcal{G}' = (\mathcal{V}', \mathcal{E}', \mathbf{X}', \mathbf{Y}')$  be a node-level adjacent graph dataset to  $\mathcal{G}$  with adjacency matrix  $\mathbf{A}'$  and feature matrix  $\mathbf{X}'$ , also having maximum degree bounded above by  $D$ . The goal is to show the following inequality:

$$\|\text{AGG}(\mathbf{X}, \mathbf{A}) - \text{AGG}(\mathbf{X}', \mathbf{A}')\|_F \leq \sqrt{D}.$$

Let  $\mathbf{M} = \text{AGG}(\mathbf{X}, \mathbf{A})$  be the aggregation function output on  $\mathbf{A}$ , and

$$\mathbf{M}_i = \sum_{j=1}^N \mathbf{A}_{j,i} \mathbf{X}_j,$$

be the  $i$ -th row of  $\mathbf{M}$  corresponding to the aggregated vector for the  $i$ -th node. Analogously, let  $\mathbf{M}' = \text{AGG}(\mathbf{X}, \mathbf{A}')$ . Note that  $\mathcal{G}'$  has one node less than  $\mathcal{G}$ , thus both  $\mathbf{X}'$  and  $\mathbf{A}'$  have respectively one row less than  $\mathbf{X}$  and  $\mathbf{A}$ . Accordingly,  $\mathbf{M}'$  has also one row less than  $\mathbf{M}$ . To work with matrices  $\mathbf{M}$  and  $\mathbf{M}'$  of the same dimension, we model the removal of a node from  $\mathcal{G}$  by the addition of an isolated node to  $\mathcal{G}'$ , which is equivalent to filling the missing row in  $\mathbf{M}'$  with an all zero vector. Note that an isolated node does not have any outbound edges, so it does not contribute to the aggregation of any other node, and thus the aggregation output of all the other nodes remain the same regardless of the presence of the isolated node. Furthermore, as an isolated node does not have any inbound edges either, the size of its aggregation set is zero, and thus its aggregation output is a zero vector (regardless of the values of its features and label).

With the above notations, bounding the sensitivity of the aggregation function amounts to bounding  $\|\mathbf{M} - \mathbf{M}'\|_F$ . Note that there exists a node  $v$  (isolated in  $\mathcal{G}'$ ) such that:

$$\begin{cases} \mathbf{A}'_{i,j} = 0, & \text{if } i = v \text{ or } j = v, \\ \mathbf{A}'_{i,j} = \mathbf{A}_{i,j}, & \text{otherwise.} \end{cases} \quad (16)$$

Then:

$$\begin{aligned}
\|\mathbf{M} - \mathbf{M}'\|_F &= \left( \sum_{i=1}^N \|\mathbf{M}_i - \mathbf{M}'_i\|_2^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^N \left\| \sum_{j=1}^N (\mathbf{A}_{j,i} \mathbf{X}_j - \mathbf{A}'_{j,i} \mathbf{X}'_j) \right\|_2^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^N \|\mathbf{A}_{v,i} \mathbf{X}_v - \mathbf{A}'_{v,i} \mathbf{X}'_v\|_2^2 + \left\| \sum_{j=1}^N (\mathbf{A}_{j,v} \mathbf{X}_j - \mathbf{A}'_{j,v} \mathbf{X}'_j) \right\|_2^2 \right. \\
&\quad \left. - \|\mathbf{A}_{v,v} \mathbf{X}_v - \mathbf{A}'_{v,v} \mathbf{X}'_v\|_2^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^N \|\mathbf{A}_{v,i} \mathbf{X}_v\|_2^2 + \left\| \sum_{j=1}^N \mathbf{A}_{j,v} \mathbf{X}_j \right\|_2^2 - \|\mathbf{A}_{v,v} \mathbf{X}_v\|_2^2 \right)^{1/2} \\
&= \left( \sum_{i=1}^N \mathbf{A}_{v,i} \|\mathbf{X}_v\|_2^2 + \sum_{j=1}^N \mathbf{A}_{j,v} \|\mathbf{X}_j\|_2^2 - \mathbf{A}_{v,v} \|\mathbf{X}_v\|_2^2 \right)^{1/2} \\
&\leq \left( \sum_{i=1}^N \mathbf{A}_{v,i} + \sum_{j=1}^N \mathbf{A}_{j,v} - \mathbf{A}_{v,v} \right)^{1/2} \leq \sqrt{D},
\end{aligned}$$

which concludes the proof.  $\square$

We can now prove [Theorem 2](#).

*Proof.* The PMA mechanism applies the Gaussian mechanism on the output of the summation aggregation function  $\text{AGG}(\mathbf{X}, \mathbf{A}) = \mathbf{A}^T \cdot \mathbf{X}$ . Based on [Lemma 2](#), the node-level sensitivity of  $\text{AGG}(\cdot)$  for bounded graphs with maximum degree  $D$  is  $\sqrt{D}$ . Therefore, each individual application of the Gaussian mechanism is  $(\alpha, \alpha D/2\sigma^2)$ -RDP. As PMA can be seen as an adaptive composition of  $K$  such mechanisms, based on Proposition 1 of [\[36\]](#), the total privacy cost is  $(\alpha, \alpha DK/2\sigma^2)$ -RDP.  $\square$

#### A.4 Proof of [Proposition 3](#)

*Proof.* Under node-level DP, all the information pertaining to an individual node, including its features, label, and edges, are private. The first step of [Algorithm 2](#) privately processes the node features and labels so as to satisfy  $(\alpha, \epsilon_1(\alpha))$ -RDP. Steps 2 and 3 of the algorithm, however, expose the private node features, but then they are processed by steps 4 and 5, which are  $(\alpha, DK\alpha/2\sigma^2)$ -RDP (according to [Theorem 2](#)) and  $(\alpha, \epsilon_5(\alpha))$ -RDP, respectively. As a result, [Algorithm 2](#) can be seen as an adaptive composition of an  $(\alpha, \epsilon_1(\alpha))$ -RDP mechanism, an  $(\alpha, DK\alpha/2\sigma^2)$ -RDP mechanism, and an  $(\alpha, \epsilon_5(\alpha))$ -RDP mechanism. Therefore, based on Proposition 1 of [\[36\]](#), the total node-level privacy cost of [Algorithm 2](#) is  $(\alpha, \epsilon_1(\alpha) + DK\alpha/2\sigma^2 + \epsilon_5(\alpha))$ -RDP, which ensures  $(\epsilon_1(\alpha) + \epsilon_5(\alpha) + \frac{DK\alpha}{2\sigma^2} + \frac{\log(1/\delta)}{\alpha-1}, \delta)$ -DP based on [Proposition 1](#).  $\square$