# Graph Unlearning with Efficient Partial Retraining

Jiahao Zhang, Lin Wang, Shijie Wang, Wenqi Fan

The Hong Kong Polytechnic University

Kowloon, Hong Kong SAR

{tony-jiahao.zhang,comp-lin.wang,shijie.wang}@connect.polyu.hk,wenqifan03@gmail.com

## ABSTRACT

Graph Neural Networks (GNNs) have achieved remarkable success in various real-world applications. However, GNNs may be trained on undesirable graph data, which can degrade their performance and reliability. To enable trained GNNs to efficiently unlearn unwanted data, a desirable solution is retraining-based graph unlearning, which partitions the training graph into subgraphs and trains sub-models on them, allowing fast unlearning through partial retraining. However, the graph partition process causes information loss in the training graph, resulting in the low model utility of sub-GNN models. In this paper, we propose GraphRevoker, a novel graph unlearning framework that better maintains the model utility of unlearnable GNNs. Specifically, we preserve the graph property with graph property-aware sharding and effectively aggregate the sub-GNN models for prediction with graph contrastive sub-model aggregation. We conduct extensive experiments to demonstrate the superiority of our proposed approach.

## CCS CONCEPTS

• **Computing methodologies → Machine learning**; • **Mathematics of computing → Graph algorithms**.

## KEYWORDS

Machine Unlearning; Graph Neural Networks; Graph Machine Unlearning; Contrastive Learning

## 1 INTRODUCTION

As one of the crucial data representations, graphs are used to describe data with complex relations between objects. Recently, Graph Neural Networks (**GNNs**) have achieved remarkable success in learning from graph data in various real-world applications.

Despite the aforementioned success of GNNs, the concerns about undesirable data posing detrimental effects on GNNs are rising. For instance, *malicious data* [6, 11], *low-quality data* [9, 24], and *sensitive data* [5, 7] are threatening GNNs' safety, prediction performance, and compliance to "the right to be forgotten" (as shown in Fig. 1), which necessities the development of **graph unlearning**.

To efficiently remove the effect of undesirable data from trained GNN models, recent years have witnessed two separate lines of research: 1) *approximate unlearning* and 2) *retraining-based unlearning*. Approximate unlearning relies on directly manipulating the model parameters to erase the effects of unwanted data points [8, 27], which sacrifices the reliability of the removal, as discussed in previous literatures [2, 22]. On the other hand, retraining approaches partition the training graph into small subgraphs and train disjoint sub-models on them, allowing the model owner only to retrain a small sub-model to remove the effect of some bad data. Despite their desirable security properties [2, 22], these
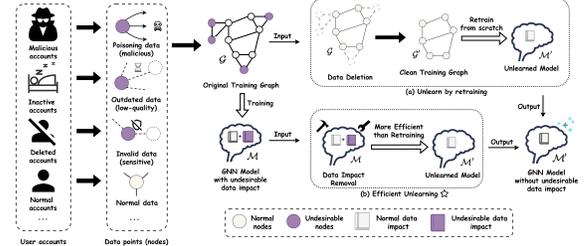


**Figure 1: An example of graph unlearning in a social network.**

unlearning approaches may destroy the graph structure and label semantic during the partitioning process, degrading the model utility of sub-model GNNs.

In this work, our research objective revolves around 1) preserving the desirable property of retraining-based unlearning and 2) significantly improving the model utility of sub-GNN models in this paradigm. To achieve this goal, we systematically review the limitations of prior works in graph unlearning (as discussed in Section 3), and then introduce a novel graph unlearning framework, namely **GraphRevoker**, equipped with *graph property-aware sharding* and *graph contrastive sub-model aggregation*. Our contributions can be summarized as follows: **1)** We propose the *graph property-aware sharding* module to preserve the sub-GNN models' prediction performance by keeping the structural and semantic properties in the training graph. **2)** To effectively leverage the disjoint sub-models for prediction, we propose the *graph contrastive sub-model aggregation* module, a lightweight GNN ensemble network, empowered by local-local structural reconstruction and local-global contrastive learning. **3)** Extensive experiments are conducted to illustrate the model utility and unlearning efficiency of the proposed method.

## 2 PROBLEM FORMULATION

**Notations.** In general, a graph can be denoted as $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{u_1, \cdots, u_N\}$ denotes the node set, and the edge set $\mathcal{E}$ is represented by adjacency matrix $\boldsymbol{A}$. We define the diagonal degree matrix $\boldsymbol{D}$ as $\boldsymbol{D}_{i,i} = deg(u_i) = \sum_{j=1}^{N} \boldsymbol{A}_{i,j}$. For node classification, the labels of nodes can be represented as $\mathcal{Y} = [y_1, \cdots, y_N]$, in which $y_i \in \{\ell_1, \cdots, \ell_C\}$ denotes the $C$ different categories.

**Problem Definition.** With the concerns on malicious and low-quality graph data, given a trained GNN model $\mathcal{F}_{\boldsymbol{\theta}}$, the goal of graph unlearning is to eliminate the impact of an undesirable subset of training data $\mathcal{D}^-$ from $\mathcal{F}_{\boldsymbol{\theta}}$ while preserving its model utility. Specifically, in the context of node classification, the undesirable subset with size $t$ can be represented as a set of undesirable nodes $\mathcal{D}^- = \{v_{j_1}, \cdots, v_{j_t}\} \subset \mathcal{V}$.

Since undesirable knowledge from $\mathcal{D}^-$ has been encoded into the parameters of $\mathcal{F}_{\boldsymbol{\theta}}$, *retraining-based graph unlearning* aims to obtain an unlearned GNN model $\mathcal{F}_{\boldsymbol{\theta}'_u}$, which is equivalent to the GNN model $\mathcal{F}_{\boldsymbol{\theta}'_r}$ that *has never been trained* on $\mathcal{D}^-$. This definition

Jiahao Zhang, Wenqi Fan, Lin Wang, Shijie Wang.

**Table 1: Characteristics of different unlearning frameworks.**

| Method / Aspect | Approximate[8, 27] | Retraining-based | | |
|---|---|---|---|---|
| | | SISA[2] | GraphEraser[7] | Ours |
| Accurate Removal | ✗ | ✓ | ✓ | ✓ |
| Structural Preservation | ✓ | ✗ | ✓ | ✓ |
| Semantic Preservation | ✓ | ✓ | ✗ | ✓ |
| Effective Ensemble | N/A | ✗ | ✗ | ✓ |

of graph unlearning can be formalized as follows:

$$\text{Retraining: } \mathcal{G} \xrightarrow{\text{data removal}} \mathcal{G}/\mathcal{D}^- \xrightarrow{\text{retrain}} \mathcal{F}_{\theta'_r};$$

$$\text{Unlearning: } \mathcal{G} \xrightarrow{\text{train}} \mathcal{F}_\theta \xrightarrow{\text{unlearn}(\mathcal{D}^-)} \mathcal{F}_{\theta'_u},$$

where $\mathcal{F}_{\theta'_u}$ and $\mathcal{F}_{\theta'_r}$ are expected to be equivalent.

In graph unlearning, we also aim to achieve two goals: *1) Model Utility*: After unlearning data points $\mathcal{D}^-$, the latest model $\mathcal{F}_{\theta'_u}$ should have comparable performance in comparison with retrained model $\mathcal{F}_{\theta'_r}$; *2) Unlearning Efficiency*: Compared to retraining from scratch, the unlearning process is expected to be more efficient.

## 3 RELATED WORK

In this section, we briefly review the state-of-the-art approaches in graph unlearning, and further discuss their limitations compared with our proposed framework. Additional related works can be found in Section A.4.

**Approximate Unlearning**. To avoid the costly retraining process when unlearning undesirable data, approximate approaches directly manipulate the parameters of GNNs, and alleviate the impact of the unwanted data. For instance, GIF [27] develops a novel graph influence function tailored to GNNs and updates the parameters with the gradients from the influence function, while GNNDelete [8] inserts a trainable delete operator between each GNN layer and optimizing the parameters of the delete operator to achieve unlearning. However, recent literature [2, 22] finds these methods do not guarantee the accurate removal of the deleted data, which necessitates the existence of retraining in unlearning frameworks.

**Efficient Retraining**. Regarding both accurate removal and unlearning efficiency, SISA [2] has introduced the retraining-based paradigm, randomly partitioning the training graph into small subgraphs, thereby training sub-models on them. Thus, exact data impact removal can be achieved by only retraining a sub-model. Later on, GraphEraser [7] adapts the SISA [2] paradigm to the context of graphs by introducing a clustering-based partitioning strategy and learnable sub-model aggregation, enhancing the model prediction performance while still allowing efficient and exact unlearning.

Despite the desirable removal guarantee, retraining-based graph unlearning still faces challenges in sub-optimal model utility, due to the graph structural loss in the partition stage. It is true that clustering-based partitions [7] provide a remedy for this issue, but cluster-based partitions tend to put nodes with the same label into one subgraph, sacrificing the semantic balance between different submodels. In addition, the lack of expressive sub-GNN aggregators also deteriorates the model utility problem in retraining approaches.

To this end, we propose GraphRevoker, which achieves efficient and accurate data removal, while preserving both the structural and semantic information in the graph partition phase. We also introduce a powerful contrastive learning empowered aggregation module to ensemble the subgraph GNNs. The advantages of our proposed model can be seen in Table. 1.

## 4 PROPOSED METHOD

To endow GNNs with efficient and exact unlearning capability, GraphRevoker follows the widely adopted SISA framework, which mainly contains three stages: 1) subgraph partition, 2) isolated training, and 3) sub-model aggregation, as shown in Fig. 2. In both partitioning and aggregation stages, we introduce innovative designs to preserve the model utility of the framework. After the aforementioned steps, we can apply GraphRevoker to make predictions by aggregating all the predictions from subgraph GNN models, or efficiently unlearn data points by partially retraining the affected sub-models.

### 4.1 Graph Property-Aware Sharding

To preserve model utility and unlearning efficiency in the partition phase, we first formulate the unlearning goals in Section 2 into three reachable optimization objectives, and then solve them with an effective neural framework to give desirable graph partitions. In comparison to previous random [2] and clustering-based [7] partition methods, our framework preserves both graph structure and label semantics, resulting in stronger sub-GNN models.

**Unlearning Time.** The efficiency of retraining-based unlearning mainly depends on the time cost of retraining corresponding sub-models, which relies on two key factors: 1) the probability of retraining a specific sub-model and 2) the time cost of retraining that sub-model. Clearly, the first factor relates to the number of nodes in a subgraph, and the second factor lies behind the cost of message-passing, which is proportional to the number of edges [34]. Let $S$ denote the number of partitioned subgraphs, and $\mathcal{P} = \{\mathcal{V}_1, \cdots, \mathcal{V}_S\}$ denotes node partition results, satisfying $\forall \mathcal{V}_i \in \mathcal{V}$ and $\forall i \neq j, \mathcal{V}_i \cap \mathcal{V}_j = \emptyset$. We propose the unlearning time objective as the expectation of retraining time as follows:

$$\mathcal{L}_{time} = \sum_{i=1}^{S} Pr(\mathcal{V}_i) \cdot \text{Cost}(\mathcal{V}_i) \approx \sum_{i=1}^{S} \frac{|\mathcal{V}_i|}{|\mathcal{V}|} |\mathcal{E}_i|. \quad (1)$$

**Graph Structure Preservation.** Edges that connect nodes in different subgraphs are inevitably removed during the partition phase. However, to maintain the model utility of sub-GNNs, preserving the original graph's structure is crucial. Prior works in graph unlearning [5, 7] achieve structural preservation by using balanced K-Means [7] or balanced Label Propagation [5] in an unsupervised setting, which cannot directly protect the graph structure. To address this issue, we propose a principled and supervised objective to quantize the structural loss in the graph partitions, namely the normalized edge-cut (*Ncut*) objective, counting the number of dropped edges between different subgraphs as follows:

$$\mathcal{L}_{struct} = Ncut(\mathcal{G}_1, \mathcal{G}_2, \cdots, \mathcal{G}_S) = \sum_{k=1}^{S} \frac{cut(\mathcal{G}_k, \bar{\mathcal{G}}_k)}{\sum_{u_i \in \mathcal{G}_k} deg(u_i)}, \quad (2)$$

where $cut(\mathcal{G}_k, \bar{\mathcal{G}}_k)$ denotes the edges between subgraph $\mathcal{G}_k$ and the remaining part of the training graph.

**Label Semantic Preservation.** Besides the structural destruction in the partition phase, the label distribution of the training graph is also perturbed, which may result in biased and less generalizable subgraph models. This phenomenon is more severe in clustering-based partitions [5, 7], which tend to assign nodes with similar
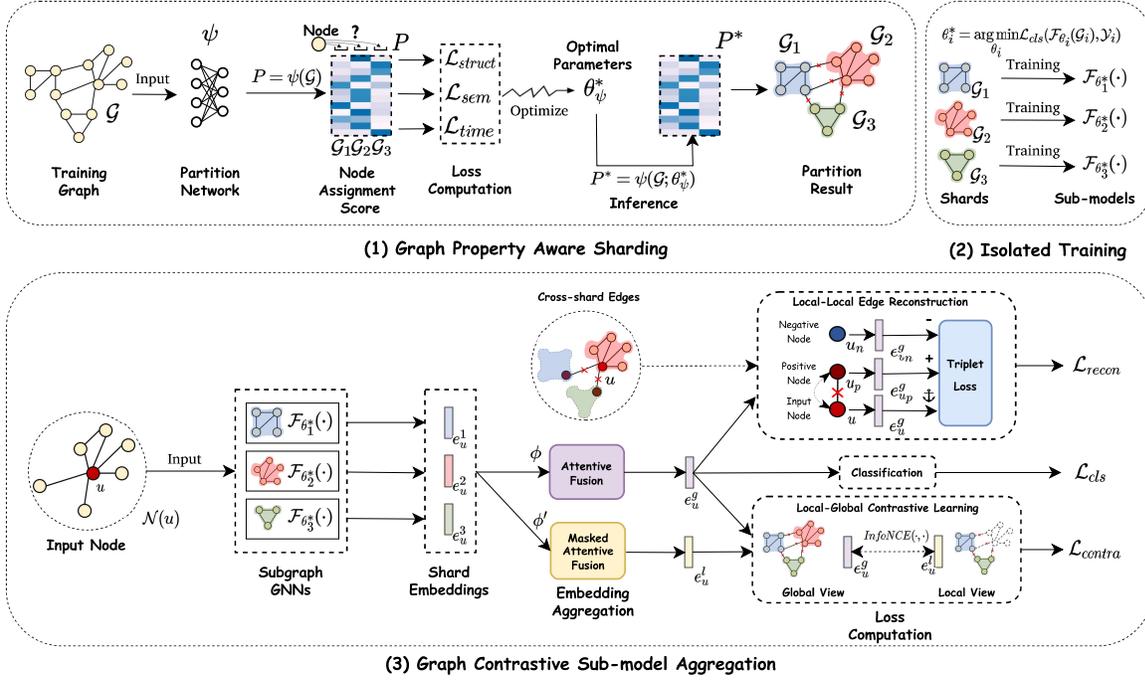
Figure 2: Illustration of the proposed framework.

labels into one subgraph. To overcome this problem, we propose an objective to capture the richness of label semantics in each subgraph with the entropy of its label distribution.

Let $c_{\ell_j}(\mathcal{V}_i) = |\{u_k | u_k \in \mathcal{V}_i, \mathcal{Y}_k = \ell_j\}|$ denote the number of nodes annotated with label $\ell_j$ in the $i$-th shard. We have the discrete label distribution $\boldsymbol{d}(\mathcal{V}_i)$ for the $i$-th subgraph, where $Pr(\ell_i) = \frac{c_{\ell_i}(\mathcal{V}_i)}{|\mathcal{V}_i|}$. Therefore, we propose the entropy-based semantic preservation objective as follows:

$$\mathcal{L}_{sem} = \frac{1}{S} \sum_{i=1}^{S} \text{Entropy}\left[\boldsymbol{d}(\mathcal{V}_i)\right] = \frac{1}{S} \sum_{i=1}^{S} \sum_{j=1}^{C} -\log\left[\frac{c_{\ell_j}(\mathcal{V}_i)}{|\mathcal{V}_i|}\right]\frac{c_{\ell_j}(\mathcal{V}_i)}{|\mathcal{V}_i|}. \quad (3)$$

**Differentiable Graph Partition Framework.** To partition the training graph $\mathcal{G}$ into disjoint subgraphs while minimizing objectives Eqs. (1) - (3), the main challenge lies behind optimization. As graph partition is a combinatorial optimization problem and cannot be solved in polynomial time, we relax the graph partition into a continuous form and solve it with a neural network.

To make graph partitioning continuous, we represent the partition result with a soft node assignment matrix $\boldsymbol{P} \in \mathbb{R}^{N \times S}$, where $P_{i,j} \in [0, 1]$ and $\sum P_{i,:} = 1$. The node assignment matrix can be computed with a graph partition network $\psi$ parameterized with $\boldsymbol{\theta}_\psi$, which contains GNN layers and a softmax output layer, transforming the node representations into the assignment matrix. Therefore, we can learn how to give effective graph partitions by training network $\psi$ with the following loss function:

$$\mathcal{L}_{part} = \mathcal{L}_{time} + \mathcal{L}_{struct} + \mathcal{L}_{sem} + \frac{1}{2}\gamma\|\boldsymbol{\theta}_\psi\|_2^2, \quad (4)$$

$$\text{s.t.} \quad \mathcal{L}_{time} = \frac{1}{|\mathcal{V}|} \sum_{i=1}^{S} (\mathbf{1}^T \boldsymbol{P}_{:,i}) \sum_{reduce} [(\boldsymbol{P}_{:,i}\boldsymbol{P}_{:,i}^T) \odot \boldsymbol{A}], \quad (5)$$

$$\mathcal{L}_{struct} = \sum_{reduce} [\boldsymbol{P} \cdot diag^{-1}(\mathbf{1}^T \boldsymbol{D}\boldsymbol{P})](1 - \boldsymbol{P}^T) \odot \boldsymbol{A}, \quad (6)$$

$$\mathcal{L}_{sem} = \frac{1}{S} \sum_{i=1}^{S} \sum_{j=1}^{C} -log\left(\frac{\boldsymbol{P}_{:,i}^T \boldsymbol{Y}_{:,j}}{\boldsymbol{P}_{:,i}^T \mathbf{1}}\right)\frac{\boldsymbol{P}_{:,i}^T \boldsymbol{Y}_{:,j}}{\boldsymbol{P}_{:,i}^T \mathbf{1}}, \quad (7)$$

where $\mathbf{1}$ denotes 1-valued column vector, and $\boldsymbol{Y} \in \{0, 1\}^{|\mathcal{V}| \times C}$ denotes the one-hot label matrix. After training $\psi$ on the training graph with $\mathcal{L}_{part}$, we can use this network to infer desirable partitions. Please refer to Section A.3 in our supplementary material for further explanations of the partition loss function.

## 4.2 Graph Contrastive Sub-model Aggregation

After acquiring the subgraph partition $\mathcal{P} = \{\mathcal{V}_1, \cdots, \mathcal{V}_S\}$, we can train $S$ isolated submodels $\mathcal{F}_{\boldsymbol{\theta}_1}, \cdots, \mathcal{F}_{\boldsymbol{\theta}_S}$ on these disjoint subgraphs. Thus, unlearning can be efficiently achieved by only retraining a single sub-model affected by unwanted data. However, unlearning efficiency is not the only goal of retraining-based unlearning, and another challenge lies behind the difficulty of leveraging the weak sub-models to obtain an accurate prediction. Though previous works have explored some straightforward solutions, including mean averaging [2], weighted averaging [7], and attention mechanism [5], it is still an open question on how to utilize the graph structures in the sub-model aggregation phase.

In this work, we develop a contrastive learning framework to learn an effective aggregator to ensemble the weak sub-GNN models. The aggregator only contains a few parameters and can be trained efficiently with only a small subset of training nodes, namely $\mathcal{U} = \{u_{j_1}, \cdots, u_{j_M}\} \subset \mathcal{V}, |\mathcal{U}| \ll |\mathcal{V}|$.

**Attentive Fusion.** When making predictions with the trained submodels, every sub-model $\mathcal{F}_{\boldsymbol{\theta}_i}(\cdot)$ generates a node embedding matrix $\boldsymbol{E}^i \in \mathbb{R}^{|\mathcal{U}| \times d}$, where row vector $\boldsymbol{e}_u^i$ represents the embedding of

node $u \in \mathcal{U}$. Towards a better global prediction, we first align embeddings from different feature spaces with a learnable linear projection, and then fuse all the embeddings with an attention mechanism. Let $\bar{E} \in \mathbb{R}^{|\mathcal{V}| \times d}$ denote the fused embedding matrix, and $\bar{e}_u$ denote a row vector in $\bar{E}$, we define the attentive fusion as follows:

$$\alpha_u^i = \frac{\exp(\boldsymbol{w}^T \text{ReLU}(\boldsymbol{W}_i \boldsymbol{e}_u^i + \boldsymbol{b}_i))}{\sum_{j=1}^{S} \exp(\boldsymbol{w}^T \text{ReLU}(\boldsymbol{W}_j \boldsymbol{e}_u^j + \boldsymbol{b}_j))}; \quad \bar{e}_u = \frac{1}{S} \sum_{i=1}^{S} \alpha_u^i \boldsymbol{e}_u^i, \quad (8)$$

where $\{\boldsymbol{W}_1, \cdots, \boldsymbol{W}_S\}$, $\{\boldsymbol{b}_1, \cdots, \boldsymbol{b}_S\}$ and $\boldsymbol{w}$ are trainable parameters, and $\alpha_u^i$ is the attention score for sub-model $i$ and node $u$.

**Local-global Contrastive Loss.** In unlearning frameworks based on partial retraining, the output from sub-models contains local knowledge of subgraphs in $\mathcal{G}$, while the aggregation result is expected to incorporate a full knowledge of $\mathcal{G}$. Inspired by previous Graph Contrastive Learning (GCL) approaches [26, 35], we find it natural to leverage the local views of a node to enhance its global aggregation result $\bar{e}^u$.

Specifically, we regard the fully aggregated embedding $\bar{e}_u$ as a global view of node $u$, while we generate a local view $\tilde{e}_u$ by randomly deactivating some sub-model attention scores in Eq. (8) (i.e., $\tilde{e}_u = \frac{S}{\|\boldsymbol{m}\|_1} \sum_{i=1}^{S} \boldsymbol{m}_i \alpha_u^i \boldsymbol{e}_u^i$, where $\boldsymbol{m}$ is a random 0-1 mask). Thus, two different views of the same node $(\bar{e}_u, \tilde{e}_u)$ are expected to be close, and the representation from another random node $v \neq u$ should be pulled apart. The local-global contrastive loss is formulated with the InfoNCE objective as follows:

$$\mathcal{L}_{contra} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \mathcal{L}_{InfoNCE}(\bar{e}_u, \tilde{e}_u), \quad (9)$$

$$\mathcal{L}_{InfoNCE}(\bar{e}_u, \tilde{e}_u) = -\log \frac{e^{\phi(\bar{e}_u, \tilde{e}_u)/\tau} + e^{\phi(\bar{e}_u, \bar{e}_v)/\tau} + e^{\phi(\bar{e}_u, \tilde{e}_v)/\tau}}{e^{\phi(\bar{e}_u, \tilde{e}_u)/\tau}}, \quad (10)$$

where $\phi(\cdot, \cdot)$ denotes the cosine similarity function, and $\tau$ denotes the softmax temperature. This formulation not only considers both inter-view negative pairs $(\bar{e}_u, \bar{e}_v)$ and intra-view negative pairs $(\bar{e}_u, \tilde{e}_v)$, but also requires zero external data augmentation, resulting in expressive representations and efficient computations.

**Local-local Reconstruction Loss.** In graph partitioning and isolated training, the links between subgraphs are dropped to ensure each node's impact only exists in one sub-model, which allows unlearning by only retraining one specific sub-model. Nevertheless, the dropped edges could include useful structural information of training graph $\mathcal{G}$. To address this problem, we propose the local-local reconstruction loss to restore the knowledge of the previously ignored edges as follows:

$$\mathcal{L}_{recon} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \max\{\phi(\bar{e}_u, \bar{e}_{v^+}) - \phi(\bar{e}_u, \bar{e}_{v^-}) + 1, 0\}, \quad (11)$$

where the positive sample $\bar{e}_{v^+}$ is sampled from the neighbors of $u$ in other subgraphs, and the negative sample $\bar{e}_{v^-}$ does not have any connection with $u$.

**Optimization.** Our aggregation module learns to fuse sub-GNN models by minimizing the following training objective:

$$\mathcal{L}_{aggr} = \mathcal{L}_{cls} + \mathcal{L}_{contra} + \mathcal{L}_{recon} + \frac{1}{2} \gamma' \|\Theta\|_2^2, \quad (12)$$

**Table 2: The comparison results of model utility in F1-score.**

| Datasets | Methods | GAT | GCN | SAGE | APPNP | JKNet |
|---|---|---|---|---|---|---|
| Cora | Retrain | 83.95±0.67 | 84.46±0.51 | 82.84±0.73 | 83.66±0.71 | 84.52±0.54 |
| | SISA | 67.47±4.43 | 55.02±1.06 | 58.82±4.87 | 66.64±4.61 | 68.21±4.30 |
| | GraphEraser | 75.70±0.33 | 71.27±2.99 | 72.10±0.58 | 68.84±2.56 | 69.94±3.12 |
| | **GraphRevoker** | **76.94±1.17** | **78.75±0.47** | **75.94±1.39** | **72.95±0.61** | **70.90±1.79** |
| Citeseer | Retrain | 73.37±0.83 | 72.94±0.63 | 73.17±0.75 | 71.04±0.83 | 73.24±1.01 |
| | SISA | 55.90±5.28 | 69.01±3.33 | 63.32±4.72 | 62.99±5.31 | 65.32±4.58 |
| | GraphEraser | 60.53±5.10 | 69.59±1.39 | 67.16±1.85 | 68.36±3.21 | 68.71±2.20 |
| | **GraphRevoker** | **70.27±0.80** | **72.79±1.04** | **69.35±1.42** | **70.09±1.67** | **70.62±1.45** |
| LastFM Asia | Retrain | 85.29±0.53 | 84.40±0.73 | 82.18±0.34 | 83.66±0.71 | 85.36±0.46 |
| | SISA | 75.77±0.44 | 75.46±0.52 | 74.42±0.61 | 76.06±0.35 | 75.67±0.38 |
| | GraphEraser | 68.12±2.10 | 64.77±0.96 | 64.52±1.91 | 71.55±0.90 | 69.87±0.77 |
| | **GraphRevoker** | **76.43±0.93** | **76.60±0.63** | **75.14±0.62** | **77.56±0.71** | **76.20±0.77** |
| Flickr | Retrain | 49.36±0.64 | 49.93±0.51 | 49.41±1.39 | 48.50±0.80 | 49.61±1.67 |
| | SISA | 42.81±0.86 | 46.01±1.24 | 46.92±1.24 | 46.72±1.02 | 44.38±0.95 |
| | GraphEraser | 43.82±1.68 | 46.42±1.10 | 47.52±0.62 | 45.96±1.66 | 44.32±1.55 |
| | **GraphRevoker** | **45.19±1.42** | **48.35±0.63** | **48.09±0.11** | **47.36±0.82** | **46.16±1.33** |

**Table 3: The comparison results of the time-cost of unlearning the undesirable data points ($\mathcal{D}^-$) (unit: second[s]).**

| Unlearned Nodes | Datastets | Methods | GAT | GCN | SAGE | APPNP | JKNet |
|---|---|---|---|---|---|---|---|
| 0.5% | Cora | Retrain | 30.71 | 24.93 | 17.55 | 26.26 | 24.42 |
| | | SISA | 7.79 | 5.04 | 5.43 | 5.74 | 6.41 |
| | | GraphEraser | 10.37 | 7.77 | 6.40 | 8.52 | 8.55 |
| | | **GraphRevoker** | **4.72** | **4.65** | **3.10** | **4.23** | **4.88** |
| | Citeseer | Retrain | 36.56 | 30.73 | 18.80 | 29.20 | 32.10 |
| | | SISA | 9.52 | 7.30 | 4.80 | 7.57 | 6.16 |
| | | GraphEraser | 16.39 | 10.51 | 9.89 | 12.90 | 12.70 |
| | | **GraphRevoker** | **6.25** | **6.05** | **3.85** | **5.32** | **5.87** |
| | LastFM Asia | Retrain | 106.86 | 88.57 | 73.96 | 96.63 | 100.31 |
| | | SISA | 31.59 | 46.33 | 30.58 | 55.04 | 51.81 |
| | | GraphEraser | 35.51 | 46.95 | 36.50 | 50.18 | 52.84 |
| | | **GraphRevoker** | **12.64** | **29.49** | **16.46** | **28.57** | **27.99** |
| | Flickr | Retrain | 177.11 | 138.67 | 108.06 | 155.04 | 151.35 |
| | | SISA | 139.68 | 110.88 | 72.29 | 123.16 | 115.53 |
| | | GraphEraser | 152.04 | 120.76 | 87.86 | 130.22 | 127.81 |
| | | **GraphRevoker** | **57.37** | **53.69** | **22.75** | **41.74** | **55.11** |

where $\Theta$ denotes the trainable parameters of the linear projection and attention mechanism. After unlearning each data point, the aggregation module must be retrained to ensure an accurate removal. Fortunately, we find this module needs 10 to 20 epochs to be trained, only introducing a small computational overhead.

## 5 EXPERIMENT

**Datasets.** We evaluate GraphRevoker on four real-world graph datasets, including *Cora* [31], *Citeseer* [31], *LastFM-Asia* [20], and *Flickr* [32]. For the first three datasets, we randomly divide nodes into train/val/test sets with a 0.7/0.2/0.1 ratio. For Flickr, we use the pre-defined dataset splits in PyG. More implementation details and additional experiments are presented in Section A.1 and Section A.2.

**Model Performance.** We have evaluated all the unlearning frameworks on the inductive node classification task. The model utility and unlearning efficiency results can be found in Tables. 2 - 3, and we draw the following conclusions: 1) GraphRevoker witnessed the state-of-the-art model performance among all the three efficient unlearning frameworks in both model utility and unlearning efficiency; 2) Despite the strong model utility of Retrain, its efficiency is far worse than efficient unlearning approaches, which is unacceptable in large-scale settings; 3) The propose GraphRevoker even outperforms SISA in efficiency evaluations, which can be attributed to our retraining-time-aware design in the partition module and our parallelized sub-model retraining implementations.

## 6 CONCLUSION

Graph unlearning, which involves removing the impact of undesirable data points from trained GNNs, is significant in various real-world scenarios. In this work, we propose a novel graph unlearning framework, which fully unleashes the potential of retraining-based

unlearning in graphs and GNNs to achieve accurate, model utility preserving, and efficient unlearning. To alleviate the model degradation attributed to the partition and isolated training process, we introduce two main contributions: a neural graph partition network and a graph contrastive sub-model ensemble module. In our future works, we plan to explore the effectiveness of GraphRevoker on more settings and downstream tasks in graph mining.

## REFERENCES

[1] Zainab Abbas, Vasiliki Kalavri, Paris Carbone, and Vladimir Vlassov. 2018. Streaming graph partitioning: an experimental study. In *VLDB*.
[2] Lucas Bourtoule, Varun Chandrasekaran, Christopher A Choquette-Choo, Hengrui Jia, Adelin Travers, Baiwu Zhang, David Lie, and Nicolas Papernot. 2021. Machine unlearning. In *IEEE S&P*.
[3] Jonathan Brophy and Daniel Lowd. 2021. Machine unlearning for random forests. In *ICML*.
[4] Yinzhi Cao and Junfeng Yang. 2015. Towards making systems forget with machine unlearning. In *IEEE S&P*.
[5] Chong Chen, Fei Sun, Min Zhang, and Bolin Ding. 2022. Recommendation unlearning. In *WWW*.
[6] Jingfan Chen, Wenqi Fan, Guanghui Zhu, Xiangyu Zhao, Chunfeng Yuan, Qing Li, and Yihua Huang. 2022. Knowledge-enhanced black-box attacks for recommendations. In *KDD*.
[7] Min Chen, Zhikun Zhang, Tianhao Wang, Michael Backes, Mathias Humbert, and Yang Zhang. 2022. Graph unlearning. In *ACM CCS*.
[8] Jiali Cheng, George Dasoulas, Huan He, Chirag Agarwal, and Marinka Zitnik. 2023. GNNDelete: A General Unlearning Strategy for Graph Neural Networks. In *ICLR*.
[9] Wenqi Fan, Xiaorui Liu, Wei Jin, Xiangyu Zhao, Jiliang Tang, and Qing Li. 2022. Graph trend filtering networks for recommendation. In *SIGIR*.
[10] Wenqi Fan, Shijie Wang, Xiao-yong Wei, Xiaowei Mei, and Qing Li. 2023. Untargeted Black-box Attacks for Social Recommendations. *arXiv preprint* (2023).
[11] Wenqi Fan, Han Xu, Wei Jin, Xiaorui Liu, Xianfeng Tang, Suhang Wang, Qing Li, Jiliang Tang, Jianping Wang, and Charu Aggarwal. 2023. Jointly attacking graph neural network and its explanations. In *ICDE*.
[12] Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. 2018. Predict then Propagate: Graph Neural Networks meet Personalized PageRank. In *ICLR*.
[13] Aditya Golatkar, Alessandro Achille, and Stefano Soatto. 2020. Eternal sunshine of the spotless net: Selective forgetting in deep networks. In *CVPR*.
[14] Chuan Guo, Tom Goldstein, Awni Hannun, and Laurens Van Der Maaten. 2020. Certified data removal from machine learning models. In *ICML*.
[15] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive representation learning on large graphs. In *NeurIPS*.
[16] George Karypis and Vipin Kumar. 1997. METIS: A software package for partitioning unstructured graphs, partitioning meshes, and computing fill-reducing orderings of sparse matrices. *Technical Report* (1997).
[17] Thomas N Kipf and Max Welling. 2016. Semi-Supervised Classification with Graph Convolutional Networks. In *ICLR*.
[18] Wenqing Lin. 2021. Large-scale network embedding in apache spark. In *KDD*.
[19] Yao Ma, Suhang Wang, Tyler Derr, Lingfei Wu, and Jiliang Tang. 2021. Graph adversarial attack via rewiring. In *KDD*.
[20] Benedek Rozemberczki and Rik Sarkar. 2020. Characteristic functions on graphs: Birds of a feather, from statistical descriptors to parametric models. In *CIKM*.
[21] Isabelle Stanton and Gabriel Kliot. 2012. Streaming graph partitioning for large distributed graphs. In *KDD*.
[22] Anvith Thudi, Hengrui Jia, Ilia Shumailov, and Nicolas Papernot. 2022. On the necessity of auditable algorithmic definitions for machine unlearning. In *USENIX Security*.
[23] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *ICLR*.
[24] Lin Wang, Wenqi Fan, Jiatong Li, Yao Ma, and Qing Li. 2024. Fast graph condensation with structure-based neural tangent kernel. In *WWW*.
[25] Alexander Warnecke, Lukas Pirch, Christian Wressnegger, and Konrad Rieck. 2023. Machine unlearning of features and labels. In *NDSS*.
[26] Jiahao Wu, Wenqi Fan, Jingfan Chen, Shengcai Liu, Qing Li, and Ke Tang. 2022. Disentangled contrastive learning for social recommendation. In *CIKM*.
[27] Jiancan Wu, Yi Yang, Yuchun Qian, Yongduo Sui, Xiang Wang, and Xiangnan He. 2023. GIF: A General Graph Unlearning Strategy via Influence Function. In *WWW*.
[28] Cong Xie, Ling Yan, Wu-Jun Li, and Zhihua Zhang. 2014. Distributed power-law graph computing: Theoretical and empirical analysis. *NeurIPS*.
[29] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2018. Representation learning on graphs with jumping knowledge networks. In *ICML*.
[30] Haonan Yan, Xiaoguang Li, Ziyao Guo, Hui Li, Fenghua Li, and Xiaodong Lin. 2022. ARCANE: An Efficient Architecture for Exact Machine Unlearning.. In *IJCAI*.
[31] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. 2016. Revisiting semi-supervised learning with graph embeddings. In *ICML*.
[32] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. 2019. GraphSAINT: Graph Sampling Based Inductive Learning Method. In *ICLR*.
[33] Chenzi Zhang, Fan Wei, Qin Liu, Zhihao Gavin Tang, and Zhenguo Li. 2017. Graph edge partitioning via neighborhood heuristic. In *KDD*.
[34] Jiahao Zhang, Rui Xue, Wenqi Fan, Xin Xu, Qing Li, Jian Pei, and Xiaorui Liu. 2024. Linear-Time Graph Neural Networks for Scalable Recommendations. In *WWW*.
[35] Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2021. Graph contrastive learning with adaptive augmentation. In *WWW*.

# A SUPPLEMENTARY MATERIAL

## A.1 Evaluation Settings

**Table 4: Dataset statistics.**

| Dataset | # Nodes | # Edges | # Classes | Type |
|---|---|---|---|---|
| **Cora** [31] | 2,708 | 5,429 | 7 | Citation |
| **Citeseer** [31] | 3,327 | 4,732 | 6 | Citation |
| **LastFM-Asia** [20] | 7,264 | 55,612 | 18 | Social Network |
| **Flickr** [32] | 89,250 | 899,756 | 7 | Web Image |

*A.1.1 **Dataset Statistics.*** We evaluate GraphRevoker on four real-world graph datasets from various origins, including *Cora*, *Citeseer*, *LastFM-Asia*, and *Flickr*, which are widely used to evaluate the performance of GNN models. The statistic of the datasets is shown in Table. 4.

*A.1.2 **General Settings.*** We select *GAT* [23], *GCN* [17], *SAGE* [15], *APPNP* [12], and *JKNet* [29] as our base models to evaluate the unlearning frameworks. All the GNN models except *JKNet* include two message-passing layers and an MLP classifier, and all the *JKNets* contain three message-passing layers. The embedding size for GNN models is set to 64. The partition module is optimized with an AdamW optimizer with lr=1e-3 and weight_decay=1e-5. Specifically, the partition network $\psi$ is trained in 10 to 30 epochs until convergence, and weights for $\mathcal{L}_{time}$ and $\mathcal{L}_{sem}$ are both set to 1e-3. The aggregation module is optimized with an AdamW optimizer with lr=0.01 and weight_decay=1e-5. The weights for both auxiliary loss functions in aggregator training ($\mathcal{L}_{contra}$ and $\mathcal{L}_{recon}$) are set to 1e-4. For compared baselines (*SISA* and *GraphEraser*), we followed their official implementations and settings, and carefully tuned their hyperparameters based on their suggested hyperparameter search space (e.g., the learning rate of the aggregator).

*A.1.3 **Model Utility Settings.*** For all unlearning methods except *Retrain*, the number of shards is fixed to 20. Each experiment is conducted 10 times, and we report both the average and standard deviation of the results. All the GNN models are trained for 100 epochs except on the *Flickr* dataset, where we train the GNN models for 20 epochs in *SISA*, *GraphEraser* and our *GraphRevoker* to avoid over-fitting. For the number of samples to train the learnable aggregator ($|\mathcal{U}|$), we follow the settings of *GraphEraser*, which selects 1000 samples in datasets with small or medium sizes (*Cora*, *Citerseer* and *LastFM Asia*), and selected 10% of the training nodes in larger datasets (*Flickr*).

*A.1.4 **Unlearning Efficiency Settings.*** We randomly select 0.5% nodes from the node-set $\mathcal{V}$ as the undesirable data points $\mathcal{D}^-$, and report the total time cost, including retraining sub-models affected by undesirable data points and reconstructing the aggregator. The settings of GNN models and the aggregator are the same as we mentioned in Section A.1.2. With a single RTX 3090 GPU, 112× Intel Xeon Gold 6238R CPU cores, and 10 sub-processes, we run the experiments and present the unlearning efficiency in Table. 3.

## A.2 Additional Experiments

*A.2.1 **Unlearning Power.*** We verify the unlearning power of GraphRevoker by enabling models affected by low-quality data points to regain their utility through unlearning undesirable data



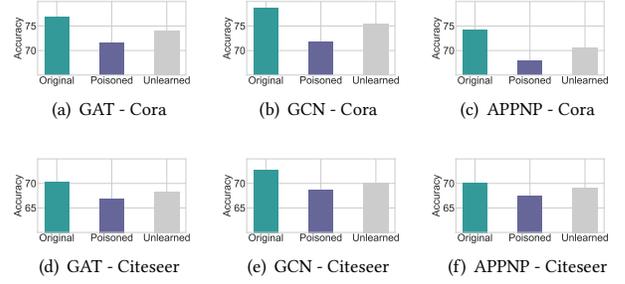| (a) GAT - Cora | (b) GCN - Cora | (c) APPNP - Cora |
|---|---|---|
| (d) GAT - Citeseer | (e) GCN - Citeseer | (f) APPNP - Citeseer |

**Figure 3: The unlearning power of our proposed framework. This figure shows the F1-score (%) on Cora and Citeseer of the original model, the model affected by low-quality data, and the model unlearned the low-quality data.**

points. In a real-world setting, low-quality data (e.g., outdated user-item interactions in recommendations) can be regarded as data points with noisy labels, so here we adopt the random attack [10, 19] to inject 100 nodes with wrong labels on Cora and Citeseer. We follow their settings and add 10 random edges for each injected node. In Figure. 3, we show the F1-score of the original GNN models and models affected by noisy data in (a), (b), (d), and (e), and then we unlearn the injected negative data points with GraphRevoker and record the model utility in (c) and (f). As illustrated in Figure. 3, we can find that the GNN models poisoned by low-quality data suffer a dramatic performance decrease, and after unlearning the injected noisy data points, the model utility could be partly regained. Though the model utility cannot be fully recovered as the noisy data also degrades the graph sharding module, our proposed graph unlearning framework could efficiently unlearn undesirable data, and the recovered model utility is comparable.

**Table 5: Ablation study for the partition module.**

| Datasets | Design | GAT | | GCN | | APPNP | |
|---|---|---|---|---|---|---|---|
| | | F1-Score | Time | F1-Score | Time | F1-Score | Time |
| Cora | **GraphRevoker** | **76.94±1.17** | **4.72** | **77.66±1.60** | **4.65** | **75.35±1.31** | **4.23** |
| | w/o $\mathcal{L}_{sem}$ | 71.99±1.60 | 4.77 | 73.82±0.68 | 4.6 | 69.54±1.26 | 4.16 |
| | w/o $\mathcal{L}_{time}$ | 75.89±1.21 | 4.92 | 75.70±1.73 | 4.73 | 71.85±1.63 | 4.3 |
| | w/o $\mathcal{L}_{time} + \mathcal{L}_{sem}$ | 70.94±1.50 | 5.07 | 71.66±0.82 | 4.93 | 72.64±1.75 | 4.47 |
| Flickr | **GraphRevoker** | **45.19±1.42** | **57.37** | **48.35±0.63** | **53.69** | **47.36±0.82** | **41.74** |
| | w/o $\mathcal{L}_{sem}$ | 44.46±0.33 | 57.94 | 47.64±1.11 | 52.23 | 46.68±1.01 | 40.48 |
| | w/o $\mathcal{L}_{time}$ | 44.54±0.89 | 58.44 | 48.00±1.20 | 54.91 | 47.19±0.88 | 42.54 |
| | w/o $\mathcal{L}_{time} + \mathcal{L}_{sem}$ | 44.09±0.36 | 59.98 | 47.88±0.91 | 55.22 | 46.09±1.39 | 44.41 |

**Table 6: Ablation study for the aggregation module.**

| Datasets | Design | GAT | GCN | APPNP |
|---|---|---|---|---|
| Cora | **GraphRevoker** | **76.94±1.17** | **78.75±0.47** | **74.27±1.16** |
| | w/o $\mathcal{L}_{contra}$ | 76.03±1.02 | 78.56±0.90 | 73.75±1.41 |
| | w/o $\mathcal{L}_{recon}$ | 75.70±1.50 | 78.34±1.21 | 74.13±0.31 |
| | w/o $\mathcal{L}_{contra} + \mathcal{L}_{recon}$ | 74.91±2.19 | 78.15±1.06 | 73.39±1.19 |
| Flickr | **GraphRevoker** | **45.19±1.42** | **48.35±0.63** | **47.36±0.82** |
| | w/o $\mathcal{L}_{contra}$ | 44.21±0.72 | 47.97±0.49 | 47.12±0.83 |
| | w/o $\mathcal{L}_{recon}$ | 44.48±1.55 | 48.20±0.82 | 47.00±1.27 |
| | w/o $\mathcal{L}_{contra} + \mathcal{L}_{recon}$ | 43.86±1.20 | 47.47±1.07 | 46.30±1.15 |

*A.2.2 **Ablation Study.*** As shown in Tables 5-6, we present the performance of GraphRevoker with different designs of training objectives in both partition and aggregation modules. We find that ablating any part of our design will result in a model performance drop. This illustrates the effectiveness of different components in our proposed GraphRevoker framework.

## A.3 Explanations of the Partition Loss Function

In Section 4.1, we directly present the loss function $\mathcal{L}_{part}$ of the differentiable graph partition framework in Eq. (4) to Eq. (7) w.r.t. the soft node assignment matrix $P$. However, it is still unclear why Eq. (5) to Eq. (7) based on $P$ are sufficient to represent the concepts (e.g., edge cuts, node counts, etc.) in Eq. (1) to Eq. (3). In this section, we supplement the derivation details of the partition training objective, illustrating the intuition behind the computations in Eq. (5) to Eq. (7).

The general idea of understanding the soft node assignment matrix $P$ is to interpret it from a probabilistic perspective. Specifically, $P_{i,j}$ denotes the probability of assigning node $u_i$ to subgraph $\mathcal{V}_j$, so we can estimate the expectation of the number of nodes, number of edges, and edge cuts in each subgraph. When inferring actual graph partitions with trained partition network $\psi$, we directly assign nodes to the subgraph with the largest probability (i.e., assigning node $u_i$ to subgraph indexed $\arg\max_j (P_{i,j})$ ).

### A.3.1 The Unlearning Time Objective.
To compute Eq. (1) with the soft assignment matrix $P \in \mathbb{R}^{N \times S}$, the key problem is representing the number of nodes $|\mathcal{V}_i|$ and edges $|\mathcal{E}_i|$ in each partitioned subgraph $\mathcal{V}_i$. It is straightforward to count the nodes in each subgraph, which is

$$\mathbb{E}|\mathcal{V}_i| = \sum_{j=1}^{N} P_{j,i} \times 1 = \mathbf{1}^T P_{:,i}, \tag{13}$$

where $\mathbf{1}$ denotes the 1-valued column vector. To count the edges in the partitioned subgraph, we have to traverse all the edges $(u_j, u_k)$ in $\mathcal{E}$ and compute the expectation as follows:

$$\mathbb{E}|\mathcal{E}_i| = \sum_{(u_j, u_k) \in \mathcal{E}} Pr(u_j, u_k \in \mathcal{V}_i) \times 1 = \sum_{A_{j,k}=1} P_{j,i} P_{k,i}. \tag{14}$$

However, computing $|\mathcal{E}_i|$ with loop structures is computationally inefficient, so we further give a vectorized form, which is

$$\mathbb{E}|\mathcal{E}_i| = \sum_{reduce} [(P_{:,i} P_{:,i}^T) \odot A], \tag{15}$$

where $\odot$ denotes point-wise multiplication, and $\sum_{reduce}$ denotes the reduce-sum operator. Thus, with Eq. (13) and Eq. (15), we can recover the proposed unlearning time loss in Eq. (5) as follows:

$$\mathcal{L}_{time} = \sum_{i=1}^{S} \frac{|\mathcal{V}_i|}{|\mathcal{V}|} |\mathcal{E}_i|$$
$$\approx \frac{1}{|\mathcal{V}|} \sum_{i=1}^{S} \mathbb{E}|\mathcal{V}_i| \mathbb{E}|\mathcal{E}_i|$$
$$= \frac{1}{|\mathcal{V}|} \sum_{i=1}^{S} (\mathbf{1}^T P_{:,i}) \sum_{reduce} [(P_{:,i} P_{:,i}^T) \odot A].$$

### A.3.2 The Structure Preservation Objective.
The normalized edge cut objective in Eq. (2) includes the edge cut in the numerator and the subgraph degree summation in the denominator. First, we show the computation of the edge cut based on the node assignment

probabilities:

$$\mathbb{E}[cut(\mathcal{G}_i, \bar{\mathcal{G}}_i)] = \sum_{(u_j, u_k) \in \mathcal{E}} Pr(u_j \in \mathcal{V}_i) Pr(u_k \notin \mathcal{V}_i) \times 1$$
$$= \sum_{A_{j,k}=1} P_{j,i}(1 - P_{k,i}) \tag{16}$$
$$= \sum_{reduce} [P_{:,i}(1 - P_{:,i})^T] \odot A.$$

The expected node degree in a subgraph can be computed as follows:

$$\mathbb{E}\left[\sum_{u_j \in \mathcal{V}_i} deg(u_j)\right] = \sum_{j=1}^{N} P_{j,i} deg(u_j) = \mathbf{1}^T DP_{:,i}. \tag{17}$$

Therefore, we can combine previous Eq. (16) and Eq. (17) to recover the *Ncut* objective in the structure preservation loss as follows:

$$\mathcal{L}_{sem} = \sum_{k=1}^{S} \frac{cut(\mathcal{G}_k, \bar{\mathcal{G}}_k)}{\sum_{u_i \in \mathcal{G}_k} deg(u_i)}$$
$$\approx \sum_{k=1}^{S} \frac{\mathbb{E}\left[cut(\mathcal{G}_k, \bar{\mathcal{G}}_k)\right]}{\mathbb{E}\left[\sum_{u_i \in \mathcal{G}_k} deg(u_i)\right]}$$
$$= \sum_{k=1}^{S} \frac{\sum_{reduce} [P_{:,k}(1 - P_{:,k})^T] \odot A}{\mathbf{1}^T DP_{:,k}}$$
$$= \sum_{reduce} \left[\sum_{k=1}^{S} \frac{P_{:,k}(1 - P_{:,k})^T}{\mathbf{1}^T DP_{:,k}}\right] \odot A$$
$$= \sum_{reduce} [P \cdot diag^{-1}(\mathbf{1}^T DP)](1 - P^T) \odot A.$$

### A.3.3 The Semantic Preservation Objective.
To recover the semantic preservation loss in Eq. (7), the computation of the number of nodes $|\mathcal{V}_i|$ is shown in Eq. (13). The only problem is counting the number of nodes annotated with label $\ell_j$ (i.e., $c_{\ell_j}(\mathcal{V}_i)$), which is computed as follows:

$$\mathbb{E}[c_{\ell_j}(\mathcal{V}_i)] = \sum_{\mathcal{Y}_k = \ell_j} P_{k,i} \times 1 = \sum_{i=1}^{N} P_{k,i} \times Y_{k,j} = P_{:,k}^T Y_{:,j}. \tag{18}$$

Therefore, with Eq. (13) and Eq. (18), we have the following derivations for the semantic preservation loss in Eq. (7):

$$\mathcal{L}_{sem} = \frac{1}{S} \sum_{i=1}^{S} \sum_{j=1}^{C} -\log[\frac{c_{\ell_j}(\mathcal{V}_i)}{|\mathcal{V}_i|}] \frac{c_{\ell_j}(\mathcal{V}_i)}{|\mathcal{V}_i|}$$
$$\approx \frac{1}{S} \sum_{i=1}^{S} \sum_{j=1}^{C} -\log[\frac{\mathbb{E}[c_{\ell_j}(\mathcal{V}_i)]}{\mathbb{E}|\mathcal{V}_i|}] \frac{\mathbb{E}[c_{\ell_j}(\mathcal{V}_i)]}{\mathbb{E}|\mathcal{V}_i|}$$
$$= \frac{1}{S} \sum_{i=1}^{S} \sum_{j=1}^{C} -log\left(\frac{P_{:,i}^T Y_{:,j}}{P_{:,i}^T \mathbf{1}}\right) \frac{P_{:,i}^T Y_{:,j}}{P_{:,i}^T \mathbf{1}}.$$

## A.4 Extended Related Works

### A.4.1 Machine Unlearning.
Machine Unlearning is a novel concept that denotes removing the impact of some undesirable data points from trained ML models. Pioneering efforts in this domain, including statistical query unlearning [4] and certified removal [14],

have provided effective and provable data impact removal solutions for simple ML models (e.g., Logistic Regression). Nevertheless, these methods cannot be simply adapted to more complex models, such as Deep Neural Networks (DNNs). To mitigate this, researchers have introduced approximate unlearning strategies aimed at removing the effect of adverse data in DNNs, which estimates the influence of these data on parameters and manipulates the corresponding parameters [13, 25]. In retrospect, approximate unlearning methods only ensure that the impact of data is alleviated, falling short of achieving exact data impact removal. To address this problem, Bourtoule et al. proposed the SISA [2] framework to achieve model-agnostic and exact unlearning via data partitioning and partial retraining. Owing to its flexibility and reliable removal guarantee, SISA and its variants have made great progress in machine unlearning tailored to image recognition [30] and decision trees [3].

Despite the aforementioned success of machine unlearning approaches, graph unlearning for graph neural networks (GNNs) is still a less explored problem. It is also non-trivial to adapt the existing solutions to the context of graphs and GNNs, since they cannot address the non-Euclidean and non-i.i.d. nature of graph learning. In this paper, we mainly focus on the graph unlearning problem, and provide innovative designs to overcome the limitations in existing solutions, as discussed in previous Section 3.

*A.4.2* **Graph Partitioning.** Graph partitioning, the process of splitting a graph into a fixed number of shards, is a prevalent technique in graph data management. Generally, previous partitioning techniques strive to balance storage load and inter-server communication overhead, thereby facilitating the efficient processing of large-scale graphs stored in distributed systems [1]. These methods partition the graph into subgraphs with balanced numbers of nodes and edges to achieve load-balancing, and minimize the edge cut between subgraphs to lower the communication cost between different server machines [16, 18, 21, 28, 33]. However, in line with our discussions in Section 3, such methods may lead to imbalanced sub-datasets and suboptimal sub-models due to their tendency to cluster interconnected nodes into the same subgraph. Furthermore, the conventional metrics of balance in graph partitioning, which focus solely on nodes or edges, do not account for the multifaceted nature of balancing the unlearning time of each subgraph GNN model. As discussed in Section 4.1, the unlearning time is simultaneously related to the number of both elements (i.e., nodes and edges) in each subgraph, highlighting the inadequacy of existing node-only or edge-only balance metrics. These limitations strongly underscore the need for a novel and systematic graph partitioning method, specifically designed for the unique demands of retraining-dependent graph unlearning.

## A.5 Limitations

While this paper represents a significant leap forward in enhancing model utility within retraining-centric graph unlearning methods, it is still not without its limitations. Notably, there remains a substantial utility loss compared with unlearning by retraining from scratch (as shown in Table. 2), which may impede practical applications of the proposed GraphRevoker framework. Furthermore, a more comprehensive evaluation is needed, encompassing both

retraining-based and approximate graph unlearning techniques, to fully show the effectiveness of the proposed method.

Besides, another concern arises from the potential for information leakage during the graph partitioning phase of retraining-based unlearning methods, such as GraphEraser and our own method. Since the graph partition results remain unchanged during sub-model retraining, some residual effects of the undesirable data could persist in the partitions even after unlearning, potentially leading to unforeseen consequences. Therefore, a more thorough theoretical analysis is necessary to ensure the robustness of data impact removal provided by these methods.