# GLASU: A Communication-Efficient Algorithm for Federated Learning with Vertically Distributed Graph Data

Xinwei Zhang\*† zhan6234@umn.edu

 $\label{lem:computer} Department\ of\ \bar{E}lectrical\ and\ Computer\ Engineering \\ University\ of\ Minnesota$ 

Mingyi Hong<sup>†</sup> mhong@umn.edu

 $\label{lem:computer_engineering} Department\ of\ Electrical\ and\ Computer\ Engineering\ University\ of\ Minnesota$ 

 ${\bf Jie\ Chen} \\ chenjie@us.ibm.com$ 

 $MIT ext{-}IBM$  Watson AI Lab IBM Research

Reviewed on OpenReview: https://openreview.net/forum?id=LHl2I2rWZa

#### **Abstract**

Vertical federated learning (VFL) is a distributed learning paradigm, where computing clients collectively train a model based on the partial features of the same set of samples they possess. Current research on VFL focuses on the case when samples are independent, but it rarely addresses an emerging scenario when samples are interrelated through a graph. In this work, we train a graph neural network (GNN) through VFL, where each client owns a part of the node features and a different edge set. This data scenario incurs a significant communication overhead, not only because of the handling of distributed features but also due to neighborhood aggregation in a GNN. Moreover, the training analysis is faced with a challenge caused by the biased stochastic gradients. We propose a model-splitting method that splits a backbone GNN across the clients and the server and a communicationefficient algorithm, GLASU, to train such a model. GLASU adopts lazy aggregation and stale updates to skip communication in neighborhood aggregation and in model updates, respectively, greatly reducing communication while enjoying convergence guarantees. We conduct extensive numerical experiments on real-world datasets, showing that GLASU effectively trains a GNN that matches the accuracy of centralized training, while using only a fraction of the time due to communication saving.

## 1 Introduction

Vertical federated learning (VFL) is a newly developed machine learning scenario in distributed optimization, where clients share data with the same sample identity but each client possesses only a subset of the features for each sample. The goal is for the clients to collaboratively learn a model based on all features. Such a scenario appears in many applications, including healthcare, finance, and recommendation systems.

Most of the current VFL solutions (Chen et al., 2020b; Liu et al., 2022) treat the case where samples are independent, but omit their relational structure. However, the pairwise relationship between samples

<sup>\*</sup>This work was done while X. Zhang was an intern at MIT-IBM Watson AI Lab, IBM Research.

<sup>&</sup>lt;sup>†</sup>M. Hong and X. Zhang are partially supported by NSF grant EPCN-2311007 and CNS-2003033. This work is also part of AI-CLIMATE: "AI Institute for Climate-Land Interactions, Mitigation, Adaptation, Tradeoffs and Economy," and is supported by USDA National Institute of Food and Agriculture (NIFA) and the National Science Foundation (NSF) National AI Research Institutes Competitive Award no. 2023-67021-39829.

emerges on many occasions, and it can be crucial in several learning scenarios, including the low-labeling-rate scenario in semi-supervised learning and the no-labeling scenario in self-supervised learning.

Consider, for example, a large technology company with multiple services under different subsidiaries, which offers news recommendations to its subscribed users. For each service (video, music, workspace, search tools, etc.) held by one subsidiary, the users could form different connections to each other and generate different historical information with the same user account. It naturally creates multiple vertically split graphs, each held by one subsidiary: a professional network where users are connected through occupational ties; a personal network where users are connected through personal life interactions; a follower network where a user is a follower of another on social media, etc. Further, the user data in each graph may contain different features (e.g., occupation-related, life-related, and interest-related, respectively). In this scenario, the portion of overlapping nodes should be large, as the users could access different products and services with the same account. To offer personal recommendations, the company sets up a server that communicates with each client (i.e., each subsidiary's data center), to train a model that predicts multiple labels for each user without revealing each client's raw local data. See Figure 1 for an illustration.

One of the most effective machine learning models for such a prediction task is graph neural networks (GNNs) (Kipf & Welling, 2016; Hamilton et al., 2017; Chen et al., 2018; Velickovic et al., 2018; Chen et al., 2020a). This model performs neighborhood aggregation in every feature transformation layer, such that the prediction of a graph node is based on not only the information of this node but also that of its neighbors.

VFL on graph-structured data is not as well studied as that on other data, in part because of the challenges incurred by an enormous amount of communication. The communication overhead comes not only from the aggregation of the partial features/representations of a datum in federated learning, but also from the neighborhood aggregation

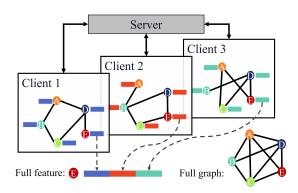


Figure 1: Data isolation of vertically distributed graph-structured data over three clients.

unique to GNNs. That is, communication occurs in each layer of the neural network, so that the latest representation of a neighboring node can be used to update the representation of the center node. One solution to reduce communication is that each client uses a local GNN to extract node representations from its own graph and the server aggregates these representations to make predictions (Zhou et al., 2020). The drawback of this method is that the partial features of a node outside one client's neighborhood are not used, even if this node appears in another client's neighborhood. Another solution is to simulate centralized training: intermediate representations of each node are aggregated by the server, from where neighborhood aggregation is performed (Ni et al., 2021). This method suffers the communication overhead incurred in each layer's computation.

In this work, we propose GLASU for communication-efficient VFL on graph data. The GNN model is split across the clients and the server, such that the clients can use a majority of existing GNNs as the backbone, while the server contains no model parameters. The server only aggregates and disseminates processed data (e.g., node embeddings) with the clients. The communication frequency between the clients and the server is mitigated through lazy aggregation and stale updates (hence the name of the method). For an L-layer GNN, GLASU communicates partial node representations only in K layers and in every other Q iterations, enjoying the reduction of communication by a factor of QL/K. GLASU can be considered as a framework that encompasses several well-known models and algorithms as special cases, including Liu et al. (2022) when the graphs are absent, Zhou et al. (2020) when all aggregations but the final one are skipped (K = 1), Ni et al. (2021) when no aggregations are skipped (K = L), and centralized training when only a single client exists.

With the enjoyable reduction in communication, another difficulty is the convergence analysis, which admits two challenges: the biased gradient caused by neighborhood sampling in training GNNs and the correlated updates due to the use of stale node representations. We conduct an analysis based on the error decomposition of the gradient, showing that the training admits a convergence rate of  $\mathcal{O}((TQ)^{-1})$ , where T is the number of training rounds, each of which contains Q iterations.

We summarize the main contributions of this work below:

- 1. Model design: We propose a flexible, federated GNN architecture that is compatible with a majority of existing GNN backbones.
- 2. Algorithm design: We propose the communication-efficient GLASU algorithm to train the model. Therein, lazy aggregation saves communication for each joint inference round, through skipping some aggregation layers in the GNN; while stale updates further save communication by allowing the clients to use stale global information for multiple local model updates.
- 3. Theoretical analysis: We provide theoretical convergence analysis for GLASU by addressing the challenges of biased stochastic gradient estimation caused by neighborhood sampling and correlated update steps caused by using stale global information. To the best of our knowledge, this is the first convergence analysis for federated learning with graph data.
- 4. Numerical results: We conduct extensive experiments on seven datasets, together with ablation studies, to demonstrate that GLASU can achieve a comparable performance as the centralized model on multiple datasets and multiple GNN backbones, and that GLASU effectively saves communication and reduces training time.

## 2 Problem, background, and related works

**Problem setup:** Consider M clients, indexed by  $m=1,\ldots,M$ , each of which holds a part of a graph with the node feature matrix  $\mathbf{X} \in \mathbb{R}^{N \times d}$  and the edge set  $\mathcal{E}$ . Here, N is the number of nodes in the graph and d is the feature dimension. The number of clients is restricted by the feature dimension and is typically small. We assume that each client has the same node set and the same set of training labels,  $\mathbf{y}$ , but a different edge set  $\mathcal{E}_m$  and a non-overlapping node feature matrix  $\mathbf{X}_m \in \mathbb{R}^{N \times d_m}$ , such that  $\mathcal{E} = \bigcup_{m=1}^M \mathcal{E}_m$ ,  $\mathbf{X} = [\mathbf{X}_1, \ldots, \mathbf{X}_M]$ , and  $d = \sum_{m=1}^M d_m$ . We denote the client dataset as  $\mathcal{D}_m = \{\mathbf{X}_m, \mathcal{E}_m, \mathbf{y}\}$  and the full dataset as  $\mathcal{D} = \{\mathbf{X}, \mathcal{E}, \mathbf{y}\}$ . The task is for the clients to collaboratively infer the labels of nodes in the test set.

#### 2.1 Graph convolutional network

The graph convolution network (GCN) (Kipf & Welling, 2016) is a typical example of the family of GNNs. Inside GCN, a graph convolution layer reads

$$\mathbf{H}[l+1] = \sigma \Big( \mathbf{A}(\mathcal{E}) \cdot \mathbf{H}[l] \cdot \mathbf{W}[l] \Big), \tag{1}$$

where  $\sigma(\cdot)$  denotes the point-wise nonlinear activation function,  $\mathbf{A}(\mathcal{E}) \in \mathbb{R}^{N \times N}$  denotes the adjacency matrix defined by the edge set  $\mathcal{E}$  with proper normalization,  $\mathbf{H}[l] \in \mathbb{R}^{N \times d[l]}$  denotes the node representation matrix at layer l, and  $\mathbf{W}[l] \in \mathbb{R}^{d[l] \times d[l+1]}$  denotes the weight matrix at the same layer. The initial node representation matrix  $\mathbf{H}[0] = \mathbf{X}$ . The classifier is denoted as  $\hat{\mathbf{y}} = f(\mathbf{H}[L], \mathbf{W}[L])$  with weight matrix  $\mathbf{W}[L]$  and the loss function is denoted as  $\ell(\mathbf{y}, \hat{\mathbf{y}})$ . Therefore, the overall model parameter is  $\mathbf{W} = {\mathbf{W}[0], \dots, \mathbf{W}[L-1], \mathbf{W}[L]}$ .

Mini-batch training of GCN (and GNNs in general) faces a scalability challenge, because computing one or a few rows of  $\mathbf{H}[L]$  (i.e., the representations of a mini-batch) requires more and more rows of  $\mathbf{H}[L-1]$ ,  $\mathbf{H}[L-2]$ , ..., recursively, in light of the multiplication with  $\mathbf{A}(\mathcal{E})$  in (1). This is known as the *explosive* neighborhood problem unique to graph-structured data. Several sampling strategies were proposed in the past to mitigate the explosion; in this work, we adopt the layer-wise sampling proposed by FastGCN (Chen et al., 2018). Starting from the output layer L, which is associated with a mini-batch of training nodes,  $\mathcal{E}[L]$ , we iterate over the layers backward such that at layer l, we sample a subset of neighbors for  $\mathcal{E}[l+1]$ , namely  $\mathcal{E}[l]$ . In doing so, at each layer, we form a bipartite graph with edge set  $\mathcal{E}[l] = \{(i,j) | i \in \mathcal{E}[l+1], j \in \mathcal{E}[l]\}$ .

Then, each graph convolution layer becomes

$$\mathbf{H}[l+1][\mathcal{S}[l+1]] = \sigma \Big( \mathbf{A}(\mathcal{E}[l]) \cdot \mathbf{H}[l][\mathcal{S}[l]] \cdot \mathbf{W}[l] \Big), \tag{2}$$

where  $\mathbf{A}(\mathcal{E}[l]) \in \mathbb{R}^{|\mathcal{S}[l+1]| \times |\mathcal{S}[l]|}$  is a properly scaled submatrix of  $\mathbf{A}(\mathcal{E})$  and  $\mathbf{H}[l][\mathcal{S}[l]]$  denotes the rows of  $\mathbf{H}[l]$  corresponding to the sampled neighbor set  $\mathcal{S}[l]$ .

#### 2.2 Related works

Vertical federated learning is a learning paradigm where the features of the data are distributed across clients, who collaborate to train a model that incorporate all features (Liu et al., 2022; Chen et al., 2020b; Romanini et al., 2021; Yang et al., 2019b; Gu et al., 2021; Yang et al., 2019a; Xu et al., 2021). Thus, the global model is split among clients and the key challenge is the heavy communication costs on exchanging partial sample information for computing the losses and the gradients for each sample. Most works consider simple models (e.g., linear) because complex models incur multiple rounds of communication for prediction.

**Federated learning with graphs** includes four scenarios. The *graph-level* scenario is *horizontal*, where each client possesses a collection of graphs and all clients collaborate to train a unified model (Zhang et al., 2021a; He et al., 2021; Bayram & Rekik, 2021; Xie et al., 2021). The task is to predict graph properties (such as molecular properties).

The subgraph-level scenario could be either vertical or horizontal. In the vertical scenario, each client holds a part of the node features, a part of the whole model, and additionally a subgraph of the global graph (Zhou et al., 2020; Ni et al., 2021). The clients aim to collaboratively train a global model (combined from those of each client) to predict node properties (such as the category of a paper in a citation network). Our work addresses this scenario.

The subgraph-level, horizontal scenario, on the other hand, considers training a GNN for node property prediction in a distributed manner: a graph is partitioned and each client holds one partition (Zhang et al., 2021b; Wu et al., 2021; Chen et al., 2022; Yao & Joe-Wong, 2022). Moreover, each client holds all features of the nodes that it possesses. A challenge to address is the aggregation of information along edges crossing different clients. This scenario differs from the vertical scenario in that features are not partitioned among clients and the graph partitions do not overlap.

The fourth scenario is *node-level*: the clients are connected by a graph and thus each of them is treated as a node. In other words, the clients, rather than the data, are graph-structured. It is akin to *decentralized learning*, where clients communicate to each other via the graph to train a unified model (Lalitha et al., 2019; Meng et al., 2021; Caldarola et al., 2021; Rizk & Sayed, 2021).

Please see Appendix A for in-depth discussions of the related works.

## 3 Proposed approach

In this section, we present the proposed model and the training algorithm GLASU for federated learning on vertically distributed graph data. The neighborhood aggregation in GNNs poses communication challenges distinct from conventional VFL. To mitigate this challenge, we propose lazy aggregation and stale updates to effectively reduce the communication between the clients and the server, while maintaining comparable prediction performance as centralized models. For notational simplicity, we present the approach by using the full-graph notation (1) but note that the implementation involves neighborhood sampling, where a more precise notation should follow (2), and that one can easily change the backbone from GCN to other GNNs.

## 3.1 GNN model splitting

We split the GNN model among the clients and the server, approximating a centralized model. Specifically, each GNN layer contains two sub-layers: the client GNN sub-layer and the server aggregation sub-layer. At

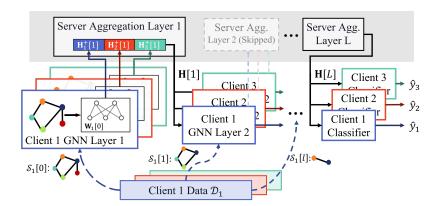


Figure 2: Illustration of the split model on M=3 clients with lazy aggregation. In the model, the second server aggregation layer is skipped and the graph size used by each layer gradually decreases, due to neighborhood aggregation (inverse of neighborhood sampling).

the l-th layer, each client computes the local feature matrix

$$\mathbf{H}_{m}^{+}[l] = \sigma \Big( \mathbf{A}(\mathcal{E}_{m}) \cdot \mathbf{H}_{m}[l] \cdot \mathbf{W}_{m}[l] \Big)$$

with the local weight matrix  $\mathbf{W}_m[l]$  and the local graph  $\mathcal{E}_m$ , where we use the superscript <sup>+</sup> to denote local representations before aggregation. Then, the server aggregates the clients' representations and outputs  $\mathbf{H}[l+1]$  as

$$\mathbf{H}[l+1] = \mathrm{Agg}(\mathbf{H}_1^+[l], \dots, \mathbf{H}_M^+[l]),$$

where  $Agg(\cdot)$  is an aggregation function. In this paper, we only consider parameter-free aggregations, including averaging and concatenation. The server broadcasts the aggregated  $\mathbf{H}[l+1]$  to the clients so that computation proceeds to the next layer. In the final layer, each client computes a prediction. This layer is the same among clients because they receive the same  $\mathbf{H}[L]$ .

The two aggregation operations of our choice render a rather simple implementation of the server. They bring in two advantages: parameter-free and memory-less. Since the operations do not contain any learnable parameters, the server does not need to perform gradient computations. Moreover, in the backward pass, these operations do not require data from the forward pass to back-propagate the gradients (memory-less). Specifically, for averaging, the server back-propagates  $\frac{1}{M}\nabla_{\mathbf{H}[l+1]}\mathcal{L}$  to each client, where  $\mathcal{L}$  denotes the loss; while for concatenation, the server back-propagates the corresponding block of  $\nabla_{\mathbf{H}[l+1]}\mathcal{L}$ .

Concatenation is more flexible as it does not require the output of the GNN layers at different clients to have the same length. With less restriction on the output feature length, the clients can have more flexibility in the layer width in model design so that GLASU can accommodate client hardware heterogeneity. In terms of communication and computation complexity, concatenation requires more communication and computation resources than averaging. Instead of broadcasting the averaged feature, the server needs to broadcast the longer concatenated feature to the clients. Moreover, the clients need to have GNN layers with larger input to process the concatenated features. If the output feature sizes of all client models are the same, then averaging can be viewed as a special case of concatenation followed by an averaging operation. In our experiments, we fix the clients' layer output dimension to be the same and use averaging as it is more communication efficient and computation/memory-friendly.

We illustrate in Figure 2 the split of each GNN layer among the clients and the server. Note the difference of our approach from existing approaches. Our model splitting resembles federated split learning (SplitFed) (Thapa et al., 2022). In SplitFed, each client can collaborate with the server to perform inference or model updates without accessing information from other clients; whereas in our case, all clients collectively perform the job. Our approach also differs from conventional VFL that splits the local feature processing and the final classifier among the clients and the server, respectively, such that each model update requires a single U-shape communication (Chen et al., 2020b). In our case, due to the graph structure, each GNN

layer contains one client-server interaction and the number of interactions is equal to the number of GNN layers (we will relax this in the following subsection).

## 3.2 Lazy aggregation

The development in the preceding subsection approximates a centralized model, but it is not communication friendly because each layer requires one round of client-server communication. We propose two communication-saving strategies in this subsection and the next. We first consider *lazy aggregation*, which skips aggregation in certain layers.

Instead of performing server aggregation at each layer, we specify a subset of K indices,  $\mathcal{I} = \{l_1, \ldots, l_K\} \subset [L]$ , such that aggregation is performed only at these layers. That is, at a layer  $l \in \mathcal{I}$ , the server performs aggregation and broadcasts the aggregated representations to the clients, serving as the input to the next layer:  $\mathbf{H}_m[l+1] = \mathbf{H}[l+1]$ ; while at a layer  $l \notin \mathcal{I}$ , each client uses the local representations as the input to the next layer:  $\mathbf{H}_m[l+1] = \mathbf{H}_m^+[l]$ . By doing so, the amount of communication is reduced from  $\mathcal{O}(L)$  to  $\mathcal{O}(K)$ .

There are subtleties caused by neighborhood sampling, similar to those faced by FastGCN (see Section 2.1). First, it requires additional rounds of communication to synchronize the sample indices, because whenever server aggregation is performed, it must be done on the same set of sampled nodes across clients. Hence, in the additional communication rounds, the server takes the union of the clients' index sets  $S_m[l_k]$  and broadcasts  $S[l_k] = \bigcup_{m=1}^M S_m[l_k]$  to the clients. Second, when server aggregation is skipped at a layer  $l \notin \mathcal{I}$ , each client can use its own set of sampled nodes,  $S_m[l]$ , which may differ from each other. Such a procedure is more flexible than conventional VFL where sample features are generally processed synchronously. The sampling procedure is summarized in Algorithm 2.

#### 3.3 Stale updates

To further reduce communication, we consider  $stale\ updates$ , which skip aggregation in certain iterations and use stale node representations to perform model updates. The key idea is to use the same minibatch, including the sampled neighbors at each layer, for training Q iterations. In every other Q iterations, the clients store the aggregated representations at the server aggregation layers. Then, in the subsequent iterations, every server aggregation is replaced by a local aggregation between a client's up-to-date node representations and other clients' stale node representations. By doing so, the clients and the server only need to communicate once in every Q iterations.

Specifically, let a round of training contain Q iterations and use t to index the rounds. At the beginning of each round, the clients and the server jointly decide the set of nodes used for training at each layer. Then, they perform a joint inference on the representations  $\mathbf{H}_m^{t,+}[l]$  at every layer  $l \in \mathcal{I}$ . Each client m will store the "all but m" representation  $\mathbf{H}_{-m}^t[l+1]$  through extracting such information from the aggregated representations  $\mathbf{H}_m^t[l+1]$ :

$$\mathbf{H}_{-m}^{t}[l+1] = \text{Extract}(\mathbf{H}_{m}^{t}[l+1], \mathbf{H}_{m}^{t,+}[l]).$$

For example, when the server aggregation is averaging, the extraction is

Extract(
$$\mathbf{H}_{m}^{t}[l+1], \mathbf{H}_{m}^{t,+}[l]$$
) =  $\mathbf{H}_{m}^{t}[l+1] - \frac{1}{M}\mathbf{H}_{m}^{t,+}[l]$ ,

Afterward, the clients perform Q iterations of model updates, indexed by  $q=0,\ldots,Q-1$ , on the local parameters  $\mathbf{W}_m^{t,q}$  in parallel, using the stored aggregated information  $\mathbf{H}_{-m}^t[l+1]$  to perform local computation, replacing server aggregation. The name "stale updates" comes from the fact that  $\mathbf{H}_{-m}^t[l+1]$  is computed by using stale model parameters  $\{\mathbf{W}_{m'}^{t,0}\}_{m'\neq m}$  at all iterations  $q\neq 0$ . The extraction and the local updates are summarized in Algorithm 3 and Algorithm 4, respectively.

## Algorithm 1 Training Procedure. All referenced algorithms are detailed in Appendix ??.

```
\begin{aligned} & \textbf{for } t = 0, \dots, T \textbf{ do} \\ & \textbf{Server/Client} \text{ (Algorithm 2): Sample } \{\mathcal{S}_m^t[l]\}_{l=0}^L. \\ & \textbf{Client: } \mathbf{W}_m^{t,0} = \begin{cases} \mathbf{W}_m^{t-1,Q}, & t > 0 \\ \mathbf{W}_m^0, & t = 0 \end{cases}. \\ & \textbf{Server/Client} \text{ (Algorithm 3): } \{\mathbf{H}_{-m}^t[l+1]\}_{l \in \mathcal{I}} = \textbf{JointInference}(\mathbf{W}_m^{t,0}, \mathcal{D}_m, \{\mathcal{S}_m^t[l]\}_{l=0}^L). \\ & \textbf{for } q = 0, \dots, Q-1 \textbf{ do} \\ & \textbf{Client} \text{ (Algorithm 4): } \mathbf{W}_m^{t,q+1} = \textbf{LocalUpdate}(\mathbf{W}_m^{t,q}, \mathcal{D}_m, \{\mathcal{S}_m^t[l]\}_{l=0}^L, \{\mathbf{H}_{-m}^t[l+1]\}_{l \in \mathcal{I}}). \\ & \textbf{end for} \\ & \textbf{output: } \{\mathbf{W}_m^{T,Q}\}_{m=1}^M \end{aligned}
```

# Algorithm 2 Sampling Procedure

```
1: Client:
                                                                1: Server:
 2: for k = K, ..., 2 do
                                                                2: Uniformly and independently sample
      Receive(S[l_k]).
                                                                     indices S[L] from training set.
3:
      Set S_m[l_k] = S[l_k].
                                                                3: Broadcast(S[L]).
 4:
      for l = l_k - 1, \dots, l_{k-1} do
                                                                4: for k = K - 1, \dots, 2 do
 5:
         Uniformly randomly sample indices \mathcal{S}_m[l]
                                                                     Aggregate(S_m[l_k]).
                                                                     Compute \mathcal{S}[l_k] = \bigcup_{m=1}^M \mathcal{S}_m[l_k].
         from neighbors of S_m[l+1].
                                                                6:
      end for
 7:
                                                                     Broadcast(S[l_k]).
                                                                7:
      Send(S_m[l_{k-1}]) if k > 2.
                                                                8: end for
9: end for
10: Output: \{S_m[l]\}_{l=0}^L
```

# ${\bf Algorithm~3~ Joint Inference}$

```
1: Client:
                                                                                                             1: Server:
 2: Input: \mathbf{W}_m, \mathcal{D}_m, \{\mathcal{S}_m[l]\}_{l=0}^L
                                                                                                             2: for l \in \mathcal{I} do
 3: Set \mathbf{H}_{m}[0] = \mathbf{X}_{m}[S_{m}[0]]
                                                                                                                      \mathbf{H}[l+1] = \text{Agg}(\mathbf{H}_{1}^{+}[l], \dots, \mathbf{H}_{M}^{+}[l]).
 4: for l = 0, ..., L - 1 do
                                                                                                                      Broadcast \mathbf{H}[l+1].
          \mathbf{H}_{m}^{+}[l] = \sigma(\mathbf{A}(\mathcal{E}(\mathcal{S}_{m}[l+1], \mathcal{S}_{m}[l]))\mathbf{H}_{m}[l]\mathbf{W}_{m}[l])
                                                                                                             5: end for
          if l \in \mathcal{I} then
 6:
              Send \mathbf{H}_{m}^{+}[l] to server
 7:
 8:
              Receive \mathbf{H}_m[l+1]
              \mathbf{H}_{-m}[l+1] = \text{Extract}(\mathbf{H}_m[l+1], \mathbf{H}_m^+[l])
 9:
10:
              Set \mathbf{H}_m[l+1] = \mathbf{H}_m^+[l]
11:
          end if
12:
13: end for
14: Output: \{\mathbf{H}_{-m}[l+1]\}_{l\in\mathcal{I}}
```

## Algorithm 4 LocalUpdate

```
1: Input: \mathbf{W}_{m}^{t,q}, \mathcal{D}_{m}, \{\mathcal{S}_{m}^{t}[l]\}_{l=0}^{L}, \{\mathbf{H}_{-m}^{t}[l+1]\}_{l\in\mathcal{I}}
  2: Set \mathbf{H}_{m}^{t,q}[0] = \mathbf{X}_{m}[\mathcal{S}_{m}^{t}[0]]
  3: for l = 0, ..., L - 1 do
              \mathbf{H}_m^{t,q,+}[l] = \sigma(\mathbf{A}(\mathcal{E}(\mathcal{S}_m^t[l+1],\mathcal{S}_m^t[l]))\mathbf{H}_m^{t,q}[l]\mathbf{W}_m^{t,q}[l])
              if l \in \mathcal{I} then
  5:
                     Set \mathbf{H}_{m}^{t,q}[l+1] = \text{Agg}(\mathbf{H}_{-m}^{t}[l+1], \mathbf{H}_{m}^{t,q,+}[l])
  6:
              else
  7:
                     Set \mathbf{H}_{m}^{t,q}[l+1] = \mathbf{H}_{m}^{t,q,+}[l]
  8:
              end if
  9:
10: end for
11: Compute loss \mathcal{L}_{m}^{t,q} = \ell\left(\mathbf{y}[\mathcal{S}_{m}^{t}[L]], f_{m}(\mathbf{H}_{m}^{t,q}[L], \mathbf{W}_{m}[L])\right)

12: Output: \mathbf{W}_{m}^{t,q+1} = \mathbf{W}_{m}^{t,q} - \eta^{t,q} \nabla_{\mathbf{W}_{m}^{t,q}} \mathcal{L}_{m}^{t,q}
```

## 3.4 Summary

The overall training procedure is summarized in Algorithm 1. For communication savings, lazy aggregation brings in a factor of L/K and stale updates bring in a factor of Q. Therefore, the overall saving factor is QL/K. Note that the algorithm assumes that all clients have the training labels. If the labels can be held by only one client (say, A), a slight modification by broadcasting the gradient with respect to the final-layer output possessed by A, suffices. See Appendix B for details.

#### 3.5 Special cases

It is interesting to note that GLASU encompasses several well-known methods as special cases.

Conventional VFL. VFL algorithms can be viewed as a special case of GLASU, where  $\mathbf{A}(\mathcal{E}_m) = \mathbf{I}$  for all m. In this case, no neighborhood sampling is needed and GLASU reduces to Liu et al. (2022).

Existing VFL algorithms for graphs. The model of Zhou et al. (2020) is a special case of GLASU, with K=1; i.e., no communication is performed except the final prediction layer. In this case, the clients omit the connections absent in the self subgraph but present in other clients' subgraphs. The model of Ni et al. (2021) is also a special case of GLASU, with K=L. This case requires communication at all layers and is less efficient.

Centralized GNNs. When there is a single client (M=1), our setting is the same as centralized GNN training. Specifically, by letting K=L and properly choosing the server aggregation function  $\mathrm{Agg}(\cdot)$ , our split model can achieve the same performance as a centralized GNN model. Note that using lazy aggregation  $(K \neq L)$  and choosing the server aggregation function as concatenation or averaging will make the split model different from a centralized GNN.

## 3.6 Privacy protection

Although this paper focuses on the utility side of FL and addresses a unique but unresolved issue pertaining to graph neural networks: the enormous communication overhead due to neighborhood aggregation inside a GNN model (besides the usual model communication problem faced by vertical FL), in this section, we provide a brief discussion on how privacy protection mechanisms can be applied to our GLASU algorithm.

First, let us clarify that the aim of differentially private GNN training is to protect local graph data. Specifically, we want to protect the data from node inference attacks (i.e., identifying if one node feature appears in some client's local subgraph.) In this case, GLASU can combine a secure aggregation method with a differentially private mechanism to guard against node inference attacks.

Secure Aggregation (SA) (Bonawitz et al., 2017; Hardy et al., 2017) is a form of secure multi-party computation approach used for aggregating information from a group of clients, without revealing the information

of any individual. This can be achieved by homomorphic encryption (Li et al., 2010; Hardy et al., 2017). In our case, when the server aggregation is averaging, homomorphic encryption can be directly applied.

**Differential Privacy (DP)** (Wei et al., 2020) is a probabilistic protection approach. By injecting stochasticity into the local outputs, this approach guarantees that an attacker cannot distinguish the sample from the dataset up to a certain probability. DP can be applied either solely or in combination with SA to our algorithm in the server-client communication, to offer privacy protection for the client data.

Specifically, for the DP mechanism, perturbation needs to be added to the first node aggregation layer; that is, for the smallest l in the set of aggregation layers,  $\mathbf{H}[l+1] = \operatorname{SecAgg}(\mathbf{H}_1^+[l] + \mathbf{w}_1, \dots, \mathbf{H}_M^+[l] + \mathbf{w}_M)$  where the injected noise  $\mathbf{w}_m, m = 1, \dots, M$  follows a normal distribution  $\mathcal{N}(0, \sigma^2 \cdot I)$  and the aggregation Agg proposed in Sec. 3.1 is replaced by a secure aggregation SecAgg (such as homomorphic encryption). With this scheme, the client's raw node features can stay private by properly choosing the noise variance  $\sigma^2$ , which depends on the sensitivity of the output of the l-th layer, the maximum degree of the nodes, and the number of training rounds T.

## 4 Convergence analysis

In this section, we analyze the convergence behavior of GLASU under lazy aggregation and stale updates. To start the analysis, denote by  $\mathcal{S}^t = \{\mathcal{S}_m^t[l]\}_{l=1,m=1}^{L,M}$  the samples used at round t (which include all sampled nodes at different layers and clients); by  $S = |\mathcal{S}_m^t[L]|$  the batch size; and by  $\mathcal{L}(\mathbf{W}; \mathcal{S})$  the training objective, which is evaluated at the overall set of model parameters across clients,  $\mathbf{W} = \{\mathbf{W}_m\}_{m=1}^M$ , and a batch of samples,  $\mathcal{S}$ .

A few assumptions are needed (see Appendix C.1 for formal statements). **A1**: The loss function  $\ell$  is  $G_{\ell}$ -smooth with  $L_{\ell}$ -Lipschitz gradient; and a client's prediction function  $f_m$  is  $G_f$ -smooth with  $L_f$ -Lipschitz gradient. **A2**: The training objective  $\mathcal{L}(\mathbf{W}; \mathcal{D})$  is bounded below by a finite constant  $\mathcal{L}^*$ . **A3**: The samples  $\mathcal{S}^t$  are uniformly sampled from the neighbor set in each layer.

**Theorem 1.** Under assumptions A1-A3, by running Algorithm 1 with constant step size  $\eta \leq C_0^{-1} \cdot (1 + 2Q^2M)^{-1}$ , with probability at least  $p = 1 - \delta$ , the averaged squared gradient norm is bounded by:

$$\frac{1}{TQ} \sum_{t=0}^{T-1} \sum_{q=0}^{Q-1} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}; \mathcal{D}) \right\|^2 \le \frac{2\Delta_{\mathcal{L}}}{\eta TQ} + \frac{28\eta M \cdot \left( C_0 + \sqrt{M+1}Q \right)}{3} \sigma,$$

where  $\Delta_{\mathcal{L}} = \mathcal{L}(\mathbf{W}^{0,0}) - \mathcal{L}^{\star}$ ,  $C_0 = G_{\ell}L_f + L_{\ell}G_f^2$ , and  $\sigma > 0$  is a function of  $\log(TQ/\delta)$ ,  $L_f$ ,  $L_g$ ,  $G_f$  and  $G_g$ . Remark 1. There are two key challenges in the analysis. (1) Owing to neighborhood sampling, the stochastic gradient is biased (i.e.,  $\mathbb{E}_{\mathcal{S}} \nabla \mathcal{L}(\mathbf{W}; \mathcal{S}) \neq \nabla \mathcal{L}(\mathbf{W}; \mathcal{D})$ ). (2) The stale updates in one communication round are correlated, as they use the same mini-batch and samples. Hence, the general unbiasedness and independence assumptions on the stochastic gradients in the analysis of SGD-type of algorithms do not apply. We borrow the technique by Ramezani et al. (2020) to bound the error of the stochastic gradient through the biasvariance decomposition and extend the analysis by Liu et al. (2022) for VFL with correlated updates to establish our proof. For details, see Appendix C.

Remark 2. To better expose the convergence rate, assuming that Q is upper bounded by  $\frac{C_0}{\sqrt{M+1}}$ , one may set  $\eta = \sqrt{\frac{3\Delta_{\mathcal{L}}}{28MC_0\sigma TQ}}$ , such that

$$\frac{1}{TQ} \sum_{t=0}^{T-1} \sum_{q=0}^{Q-1} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}; \mathcal{D}) \right\|^2 \le 8\sqrt{\frac{7\Delta_{\mathcal{L}} M C_0 \sigma}{3TQ}}.$$

Ignoring the logarithmic factor  $\log(TQ/\delta)$  in  $\sigma$ , the above bound states that the squared gradient norm decreases as  $\mathcal{O}((TQ)^{-1})$ . Note that this bound holds only when T is sufficiently large, because the choice of  $\eta$  must satisfy the condition of Theorem 1.

Remark 3. Based on the preceding remark, we see that to achieve  $\epsilon$ -stationarity, the number of model updates is  $QT = \mathcal{O}(\frac{1}{\epsilon^2})$ . That is, as long as Q obeys the upper bound, running more local updates Q reduces the amount of communications Q. To the best of our knowledge, this is the first result for VFL on graph data.

Table 1: Datasets. Each of the HeriGraph datasets (Suzhou, Venice, Amsterdam) contains three naturally formed subgraphs. For other datasets, each contains one single graph and each client holds a sampled subgraph of it.

Dataset	# Nodes	# Edges	# Feat.	# Class
Cora	2,708	10, 556	1,433	7
PubMed	19,717	88, 648	500	3
CiteSeer	3,327	9, 104	3,703	6
Suzhou	3, 137	916, 496	979	9
Venice	2, 951	534, 513	979	9
Amsterdam	3, 727	1, 271, 171	979	9
Reddit	232,965	114, 615, 892	602	41

Remark 4. While we have analyzed the impact of stale updates (Q), lazy aggregation (K) does not play a role in convergence, because it does not affect model updates. Instead, it affects model accuracy in a manner similar to how changing a neural network impacts the prediction accuracy.

Remark 5. If we consider the impact of the number of clients, the factor M in the numerator of the bound indicates a slowdown when more clients participate training. Similar results are seen in FedBCD (Liu et al., 2022), but therein one can use a large batch size S to counter the slowdown. For graphs, however, S does not appear in the bound because of the biased gradient estimation. Nevertheless, we note that unlike other federated scenarios, in VFL, M is very small because it is limited by, e.g., the feature length.

## 5 Numerical experiments

In this section, we conduct numerical experiments on a variety of datasets and demonstrate the effectiveness of GLASU in training with distributed graph data. We first compare its performance with related methods, including those tackling a different assumption on the data distribution and communication pattern. Then, we examine the communication saving owing to the use of lazy aggregation and stale updates. We further showcase the flexibility of GLASU through demonstration with different GNN backbones and varying clients. The experiments are conducted on a distributed cluster with three Tesla V100 GPUs communicated through Ethernet.

#### 5.1 Datasets

We use seven datasets (in three groups) with varying sizes and data distributions: the Planetoid collection (Yang et al., 2016), the HeriGraph collection (Bai et al., 2022), and the Reddit dataset (Hamilton et al., 2017). Each dataset in the HeriGraph collection (Suzhou, Venice, and Amsterdam) contains data readily distributed: three subgraphs and more than three feature blocks for each node. Hence, we use three clients, each of which handles one subgraph and one feature block. For the other four datasets (Cora, PubMed, and CiteSeer in the Planetoid collection; and Reddit), each contains one single graph and thus we manually construct subgraphs through randomly sampling the edges and splitting the input features into non-overlapping blocks, so that each client handles one subgraph and one feature block. The dataset statistics are summarized in Table 1 and more details are given in Appendix D.1.

#### 5.2 Accuracy

We compare GLASU with three training methods: (a) centralized training, where there is only a single client (M=1), which holds the whole dataset without any data distribution and communication; (b) standalone training (Zhou et al., 2020), where each client trains a model with its local data only and they do not communicate; (c) simulated centralized training (Ni et al., 2021), where each client possesses the full graph but only the partial features, so that it simulates centralized training through server aggregation in each GNN layer. Methods (b) and (c) are typical VFL methods; they are also special cases of our method (see Section 3.5). Except for centralized training, the number of clients M=3. The number of training rounds, T, and the learning rate  $\eta$  are optimized through grid search. See Appendix D.2 for details.

Table 2: Test accuracy (%). The compared algorithms are Centralized training (Cent.); Standalone training (StAl.); Simulated centralized training (Sim.); GLASU with no stale updates, i.e., Q = 1 (GLASU-1); and GLASU with stale updates Q = 4 (GLASU-4).

		/			
Dataset	Cent.	StAl.	Sim.	GLASU-1	GLASU-4
Cora PubMed CiteSeer	$ \begin{vmatrix} 80.9 \pm 0.6 \\ 84.9 \pm 0.6 \\ 70.2 \pm 0.8 \end{vmatrix} $	$74.6 \pm 0.5$ $77.2 \pm 0.5$ $64.4 \pm 0.5$	$80.1 \pm 1.2$ $82.7 \pm 1.2$ $70.0 \pm 1.2$	$ \begin{vmatrix} 81.0 \pm 1.3 \\ 82.3 \pm 1.6 \\ 70.0 \pm 1.7 \end{vmatrix} $	$80.3 \pm 1.2$ $83.8 \pm 1.8$ $68.8 \pm 3.3$
Suzhou Venice Amsterdam	$ \begin{vmatrix} 94.3 \pm 0.3 \\ 95.7 \pm 0.5 \\ 94.6 \pm 0.1 \end{vmatrix} $	$51.6 \pm 0.9$ $33.5 \pm 2.1$ $59.8 \pm 1.0$	$93.5 \pm 0.6$ $93.1 \pm 1.3$ $95.5 \pm 0.8$	$\begin{array}{c c} 92.7 \pm 1.4 \\ 92.2 \pm 0.6 \\ 93.1 \pm 0.8 \end{array}$	$90.4 \pm 0.8$ $91.0 \pm 1.6$ $94.9 \pm 0.4$
Reddit	$95.6 \pm 0.1$	$87.3 \pm 0.3$	$95.3 \pm 0.7$	$95.7 \pm 0.6$	$94.7 \pm 1.1$

Table 3: Test accuracy (%), runtime (seconds), and saving in runtime (%) under different numbers of lazy aggregation layers (K = 4, 2, 1). The saving is with respect to K = 4. Left: PubMed; right: Amsterdam.

# Layer	K=4	K = 2	K = 1	K=4	K = 2	K = 1
Accuracy	$82.5 \pm 1.0$	$83.8 \pm 1.8$	$82.2 \pm 0.7$	$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$94.9 \pm 0.4$	$92.0 \pm 1.7$
Runtime	$130 \pm 12$	$96.6 \pm 9.9$	$81.3 \pm 6.5$	$913 \pm 76$	$544 \pm 44$	$382 \pm 35$
				_		

We use GCNII (Chen et al., 2020a) as the backbone GNN. GCNII improves over GCN through including two skip connections, one with the current layer input and the other with the initial layer input. We set the number of layers L=4 and the mini-batch size S=16. For neighborhood sampling, the sample size is three neighbors per node on average. We set K=2; i.e., lazy aggregation is performed in the middle and the last layer.

Table 2 reports the average classification accuracy of GLASU and the compared training methods, repeated five times. As expected, standalone training produces the worst results, because each client uses only local information and misses edges and node features present in other clients. The centralized training and its simulated version lead to similar performance, also as expected, because server aggregation (or its equivalence in centralized training) on each GNN layer takes effect. Our method GLASU, which skips half of the aggregations, yields prediction accuracy rather comparable with these two methods. Using stale updates (Q=4) is generally outperformed by no stale updates (Q=1), but occasionally it is better (see PubMed and Amsterdam). The gain in using lazy aggregation and stale updates occurs in timing, as will be demonstrated next.

## 5.3 Communication saving

To further investigate how the two proposed techniques affect the model performance and save the communication, we conduct a study on (a) the lazy aggregation parameter K and (b) the stale update parameter Q.

**Lazy aggregation:** We use a 4-layer GCNII as the backbone and set K = 1, 2, 4. The aggregation layers are "uniform" across the model layers. That is, when K = 1, server aggregation is performed on the last layer; when K = 2, on the middle layer and the last layer; and when K = 4, on all layers. The test accuracy and runtime are listed in Table 3. We observe that the runtime decreases drastically when using fewer and fewer aggregation layers: from K = 4 to K = 1, the reduction is 37.5% for PubMed and 58.2% for Amsterdam. The accuracy is comparable in all cases.

**Stale updates:** We experiment with a few choices of Q: 2, 4, 8, and 16. We report the time to reach the same test accuracy threshold in Table 4. We see that stale updates help speed up training by using fewer communication rounds, corroborating Remark 3 of the theory in Section 4. This trend occurs on the Amsterdam dataset even when taking Q as large as 16. The trend is also noticeable on PubMed, but at some point (Q = 8) it is reverted, likely because it gets harder and harder to reach the accuracy threshold.

Table 4: Runtime (seconds) and the test accuracy (%) to reach the *same* accuracy threshold under different numbers of stale updates (Q = 2, 4, 6, 16). Top: PubMed (threshold: 82%); bottom: Amsterdam (threshold: 89%).

# Stale	Q=2	Q = 4	Q = 8	Q = 16
Accuracy Runtime	$\begin{vmatrix} 82.5 \pm 1.6 \\ 66.1 \pm 5.0 \end{vmatrix}$	$82.0 \pm 2.4$ $43.8 \pm 4.0$	$82.1 \pm 0.3$ $88.9 \pm 7.4$	N/A > 128
# Stale	Q=2	Q = 4	Q = 8	Q = 16

Table 5: Test accuracy (%) of Centralized training (Cent.), Standalone training (StAl.), and GLASU under different numbers of clients (M = 3, 5, 7). Top: CiteSeer; bottom: PubMed.

١.				
	# Client	M=3	M = 5	M = 7
	Cent. StAl. GLASU	$ \begin{array}{ c c } \hline 64.4 \pm 0.5 \\ 68.8 \pm 3.3 \end{array} $	$70.2 \pm 0.8$ $44.6 \pm 0.3$ $69.5 \pm 1.4$	$36.6 \pm 0.8 \\ 69.4 \pm 0.7$
	# Client	M=3	M = 5	M = 7

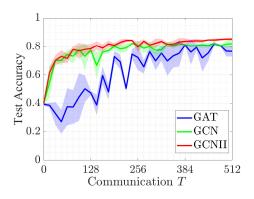


Figure 3: Test accuracy under three backbone GNNs on PubMed.

We speculate that the target 82% can never be achieved at Q = 16. This observation is consistent with Remark 2 of the theory, requiring Q to be upper bounded to claim  $\mathcal{O}((TQ)^{-1})$  convergence.

## 5.4 Flexibility

To demonstrate the flexibility of GLASU, we conduct experiments to show the performance under (a) different GNN backbones and (b) different numbers of clients, M.

Backbone model: We compare three backbones: GCN, GAT (Velickovic et al., 2018), and GCNII, which are representative GNNs. The learning rate for each backbone is tuned to its best performance. The test accuracy over training rounds is plotted in Figure 3. We see that GLASU can take different GNNs as the backbone and reach a similar prediction performance, despite that the convergence curves are not all similar. For example, the convergence histories of GCN and GCNII are quite close, whereas that of GAT experiences roughness.

Number of clients: We set M=3,5,7 and investigate the change of performance for different training methods. Hyperparameters are tuned to achieve the optimal accuracy under a fixed number of epochs. Table 5 suggests that the performance of standalone training decreases as M increases, which is expected because each client has fewer features while server aggregation is not performed. Meanwhile, the performance of GLASU is not affected and it stays comparable with that of centralized training. We note that it is unrealistic to set M arbitrarily large, because M is limited by the feature length and also in practice, it is determined by data ownership.

## 6 Conclusion

We have presented a flexible model splitting approach for VFL with vertically distributed graph data and proposed a communication-efficient algorithm, GLASU, to train the resulting GNN. Due to the graph struc-

ture, VFL on GNNs incurs heavy communication and poses an extra challenge in the convergence analysis, as the stochastic gradients are no longer unbiased. To overcome these challenges, our approach uses lazy aggregation to skip server-client communication and stale global information to update local models, leading to significant communication reduction. Our analysis makes no assumptions on unbiased gradients. We provide extensive experiments to show the flexibility of the model and the communication saving in training, without compromise on the model quality.

#### References

- Nan Bai, Pirouz Nourian, Renqian Luo, and Ana Pereira Roders. Heri-graphs: A dataset creation framework for multi-modal machine learning on graphs of heritage values and attributes with social media. *ISPRS International Journal of Geo-Information*, 11(9):469, 2022.
- Hizir Can Bayram and Islem Rekik. A federated multigraph integration approach for connectional brain template learning. In *International Workshop on Multimodal Learning for Clinical Decision Support*, pp. 36–47. Springer, 2021.
- Keith Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-preserving machine learning. In proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, pp. 1175–1191, 2017.
- Debora Caldarola, Massimiliano Mancini, Fabio Galasso, Marco Ciccone, Emanuele Rodolà, and Barbara Caputo. Cluster-driven graph federated learning over multiple domains. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2749–2758, 2021.
- Chuan Chen, Weibo Hu, Ziyue Xu, and Zibin Zheng. Fedgl: federated graph learning framework with global self-supervision. arXiv preprint arXiv:2105.03170, 2021.
- Fahao Chen, Peng Li, Toshiaki Miyazaki, and Celimuge Wu. Fedgraph: Federated graph learning with intelligent sampling. *IEEE Transactions on Parallel and Distributed Systems*, 33(8):1775–1786, 2022. doi: 10.1109/TPDS.2021.3125565.
- Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks via importance sampling. In *International Conference on Learning Representations*, 2018.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *International Conference on Machine Learning*, pp. 1725–1735. PMLR, 2020a.
- Tianyi Chen, Xiao Jin, Yuejiao Sun, and Wotao Yin. Vafl: a method of vertical asynchronous federated learning. arXiv preprint arXiv:2007.06081, 2020b.
- Bin Gu, An Xu, Zhouyuan Huo, Cheng Deng, and Heng Huang. Privacy-preserving asynchronous vertical federated learning algorithms for multiparty collaborative learning. *IEEE transactions on neural networks and learning systems*, 2021.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. Advances in neural information processing systems, 30, 2017.
- Stephen Hardy, Wilko Henecka, Hamish Ivey-Law, Richard Nock, Giorgio Patrini, Guillaume Smith, and Brian Thorne. Private federated learning on vertically partitioned data via entity resolution and additively homomorphic encryption. arXiv preprint arXiv:1711.10677, 2017.
- 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.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. arXiv preprint arXiv:1609.02907, 2016.

- Anusha Lalitha, Osman Cihan Kilinc, Tara Javidi, and Farinaz Koushanfar. Peer-to-peer federated learning on graphs. arXiv preprint arXiv:1901.11173, 2019.
- Fengjun Li, Bo Luo, and Peng Liu. Secure information aggregation for smart grids using homomorphic encryption. In 2010 first IEEE international conference on smart grid communications, pp. 327–332. IEEE, 2010.
- Yang Liu, Xinwei Zhang, Yan Kang, Liping Li, Tianjian Chen, Mingyi Hong, and Qiang Yang. Fedbcd: A communication-efficient collaborative learning framework for distributed features. *IEEE Transactions on Signal Processing*, 2022.
- Chuizheng Meng, Sirisha Rambhatla, and Yan Liu. Cross-node federated graph neural network for spatiotemporal data modeling. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery* & Data Mining, pp. 1202–1211, 2021.
- 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.
- Morteza Ramezani, Weilin Cong, Mehrdad Mahdavi, Anand Sivasubramaniam, and Mahmut Kandemir. Gcn meets gpu: Decoupling "when to sample" from "how to sample". Advances in Neural Information Processing Systems, 33:18482–18492, 2020.
- Elsa Rizk and Ali H Sayed. A graph federated architecture with privacy preserving learning. In 2021 IEEE 22nd International Workshop on Signal Processing Advances in Wireless Communications (SPAWC), pp. 131–135. IEEE, 2021.
- Daniele Romanini, Adam James Hall, Pavlos Papadopoulos, Tom Titcombe, Abbas Ismail, Tudor Cebere, Robert Sandmann, Robin Roehm, and Michael A Hoeh. Pyvertical: A vertical federated learning framework for multi-headed splitnn. arXiv preprint arXiv:2104.00489, 2021.
- Chandra Thapa, Pathum Chamikara Mahawaga Arachchige, Seyit Camtepe, and Lichao Sun. Splitfed: When federated learning meets split learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 8485–8493, 2022.
- Joel A Tropp. An introduction to matrix concentration inequalities. Foundations and Trends® in Machine Learning, 8(1-2):1–230, 2015.
- Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *stat*, 1050:4, 2018.
- 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.
- Kang Wei, Jun Li, Ming Ding, Chuan Ma, Howard H Yang, Farhad Farokhi, Shi Jin, Tony QS Quek, and H Vincent Poor. Federated learning with differential privacy: Algorithms and performance analysis. *IEEE Transactions on Information Forensics and Security*, 15:3454–3469, 2020.
- Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. Fedgnn: Federated graph neural network for privacy-preserving recommendation. arXiv preprint arXiv:2102.04925, 2021.
- Han Xie, Jing Ma, Li Xiong, and Carl Yang. Federated graph classification over non-iid graphs. *Advances in Neural Information Processing Systems*, 34:18839–18852, 2021.
- Runhua Xu, Nathalie Baracaldo, Yi Zhou, Ali Anwar, James Joshi, and Heiko Ludwig. Fedv: Privacy-preserving federated learning over vertically partitioned data. In *Proceedings of the 14th ACM Workshop on Artificial Intelligence and Security*, pp. 181–192, 2021.
- Kai Yang, Tao Fan, Tianjian Chen, Yuanming Shi, and Qiang Yang. A quasi-newton method based vertical federated learning framework for logistic regression. arXiv preprint arXiv:1912.00513, 2019a.

- Shengwen Yang, Bing Ren, Xuhui Zhou, and Liping Liu. Parallel distributed logistic regression for vertical federated learning without third-party coordinator. arXiv preprint arXiv:1911.09824, 2019b.
- Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, pp. 40–48. PMLR, 2016.
- Yuhang Yao and Carlee Joe-Wong. Fedgcn: Convergence and communication tradeoffs in federated training of graph convolutional networks. arXiv preprint arXiv:2201.12433, 2022.
- Huanding Zhang, Tao Shen, Fei Wu, Mingyang Yin, Hongxia Yang, and Chao Wu. Federated graph learning—a position paper. arXiv preprint arXiv:2105.11099, 2021a.
- 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, 34:6671–6682, 2021b.
- Longfei Zheng, Jun Zhou, Chaochao Chen, Bingzhe Wu, Li Wang, and Benyu Zhang. Asfgnn: Automated separated-federated graph neural network. *Peer-to-Peer Networking and Applications*, 14(3):1692–1704, 2021.
- Jun Zhou, Chaochao Chen, Longfei Zheng, Huiwen Wu, Jia Wu, Xiaolin Zheng, Bingzhe Wu, Ziqi Liu, and Li Wang. Vertically federated graph neural network for privacy-preserving node classification. arXiv preprint arXiv:2005.11903, 2020.

#### A Related works

Vertical federated learning is a learning paradigm where the features of the data are distributed across clients. The model is split among the clients and therefore, the key challenge is the heavy communication burden in exchanging the partial sample information.

Gu et al. (2021); Chen et al. (2020b); Xu et al. (2021); Yang et al. (2019b) consider a linear combination of the features (e.g., SVM, linear regression, logistic regression) and develop asynchronous communication-efficient protocols to aggregate the partial features. These methods can hardly be generalized to complex models in a highly nonlinear nature (e.g., multi-layer perceptrons and convolutional neural networks), because they require multiple rounds of communications. Yang et al. (2019a) use the second-order Taylor expansion to linearize the training objective for efficient communication. All these methods assume that one round of communication is sufficient for exchanging partial information, too simplistic and ineffective for graph neural networks.

**Federated learning with graphs** is a broad subject due to the flexibility of graphs. There are generally four scenarios. Some are related to graph-structured data while some are to graph-structured clients. See Figure 4 for an illustration.

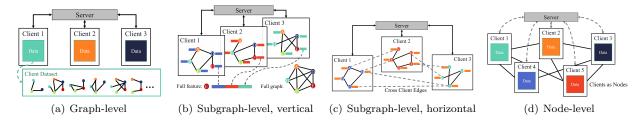


Figure 4: Four scenarios of federated learning with graphs.

(1) Graph-level (horizontal). In this scenario, each client possesses a number of graphs and all clients collaborate to train a unified model. Xie et al. (2021) apply federated learning to molecular graph property prediction with decentralized data. He et al. (2021) design a coding platform for training and evaluating graph-level federated learning algorithms. Bayram & Rekik (2021) study the scenario that multiple hospitals collaborate to train a GNN model for connectional brain templates. Wang et al. (2020) use a model-agnostic meta-learning (MAML) approach to deal with heterogeneous graphs and train a common model among clients. Zheng et al. (2021) study the setting where each client holds a heterogeneous graph dataset, a heterogeneous GNN feature extractor, and a common classifier. In this case, the clients train the local model with local datasets and upload the parameters of the classifier to the server. The server performs model averaging and optimizes the hyperparameters for local training.

Different from the scenario addressed by our paper, each client's data are independent; hence, each client can separately perform local training without gradient sharing.

(2) Subgraph-level, vertical. In this scenario, each client possesses a fraction of the graph information, particularly, partial edges and partial node features. Each client holds a partial model, which is responsible for the features and the induced subgraph it possesses. We focus on this scenario in the paper.

One way for the clients to collaboratively train the model is that they each uses a separate local GNN to extract node representations (Zhou et al., 2020). Then, the server aggregates these representations to make the prediction. In this approach, the local node representations are obtained by using a local graph, failing to capture the information from neighboring nodes in other clients' graphs. Another way of collaboration is to simulate centralized training (Ni et al., 2021). In this approach, the clients communicate the node representations after each layer computation, causing a large communication overhead.

(3) Subgraph-level, horizontal. In this scenario, each client holds a partition of the graph and all features of the nodes in the partition. Because the features are not split across clients, this scenario is not vertical; if we consider each node to be a data sample, it fits the usual horizontal definition of federated learning. A

critical challenge to address is how to economically communicate information of neighboring nodes crossing partitions, hop by hop.

In Chen et al. (2021), the clients send their local model, the predicted labels, and the node representations to the server for global semi-supervised learning that assists local model training. In Chen et al. (2022), an active node sampling method is proposed along with reinforcement learning, to improve the performance of local stochastic training. In Yao & Joe-Wong (2022), the clients first communicate the raw aggregated node features necessary for local training and then perform local training with federated model averaging. In Zhang et al. (2021b), each client holds an ego-graph with missing neighbors and trains a neural network to predict the information of these missing nodes, by federated learning. Then, each client utilizes the generated pseudo neighbor information for GNN training. In Wu et al. (2021), each client holds one user-item ego-graph, part of the entire recommendation system graph. The clients jointly train a unified model with privacy-preserving methods, using differentially private SGD updates and encrypted node sampling.

(4) Node-level. In this scenario, the clients are connected by a (communication) graph. The setting is akin to decentralized learning. The data at hand may be graph-structured or not.

In Lalitha et al. (2019); Meng et al. (2021), the clients run federated learning algorithms locally and communicate models/gradients through the graph. In the latter work, the clients possess interconnected time-series data. Each of them processes the local data and the server utilizes the underlying graph structure to run a GNN to improve global prediction. In Caldarola et al. (2021); Rizk & Sayed (2021), the underlying graph structure is kept by the server, which uses a GNN to perform aggregation on the model or data collected by the clients.

# B Modified subroutines when only one client holds the labels

Without loss of generality, we assume that the labels are held by client m = 1. Algorithms 2–4 are modified to Algorithms 5–7, respectively.

For sampling (Algorithm 5), client 1, rather than the server, generates the sampled node indices for training.

For JointInference (Algorithm 6), client 1 computes the mini-batch loss and the partial gradient with respect to the node representation of the last layer,  $\mathbf{H}_1[L]$ , and broadcasts it to all other clients through the server. We require that server aggregation is performed on the last layer; that is, for other nodes,  $\mathbf{H}_m[L] = \mathrm{Agg}(\mathbf{H}_1^+[L-1], \ldots, \mathbf{H}_M^+[L-1])$ . Then, the partial gradient with respect to  $\mathbf{H}_1[L]$  (that is,  $\nabla_{\mathbf{H}_1[L]}\mathcal{L}_1$ ) is sufficient for computing the gradients of the model parameters on all these clients. To see this, note that the clients m > 1 do not have the classifier  $\mathbf{W}_m[L]$ . Hence, we have

$$\nabla_{\mathbf{W}_{m}} \mathcal{L}_{m} = \frac{\partial}{\partial \mathbf{W}_{m}} \ell\left(\mathbf{y}[\mathcal{S}_{1}[L]], f_{1}(\mathbf{H}_{m}[L], \mathbf{W}_{1}[L])\right)$$

$$= \left\langle \frac{\partial}{\partial \mathbf{H}_{m}[L]} \ell\left(\mathbf{y}[\mathcal{S}_{1}[L]], f_{1}(\mathbf{H}_{m}[L], \mathbf{W}_{1}[L])\right), \frac{\partial}{\partial \mathbf{W}_{m}} \mathbf{H}_{m}[L] \right\rangle$$

$$= \left\langle \frac{\partial \mathcal{L}_{1}}{\partial \mathbf{H}_{1}[L]}, \frac{\partial \mathbf{H}_{1}[L]}{\partial \mathbf{W}_{m}[L]} \right\rangle. \tag{3}$$

Therefore, client 1 broadcasts  $\nabla_{\mathbf{H}_1[L]} \mathcal{L}_1$  to all other clients through the server.

For LocalUpdate (Algorithm 7), clients  $m \geq 2$  follow the chain rule to compute the second term inside the inner product of (3) and update the local model parameters.

## Algorithm 5 Sampling Procedure (Modified)

```
1: Client:
                                                             1: Server:
2: if m = 1 then
                                                            2: Recieve sample indices S[L] from client 1.
      Uniformly and independently sample
                                                            3: Broadcast(\mathcal{S}[L]).
        indices S[L] from training set.
                                                            4: for k = K - 1, ..., 2 do
                                                                  Aggregate(S_m[l_k]).
4: end if
                                                                  Compute S[l_k] = \bigcup_{m=1}^{M} S_m[l_k].
5: for k = K, ..., 2 do
                                                            6:
      Receive(S[l_k]).
6:
                                                                  Broadcast(S[l_k]).
                                                             7:
      Set S_m[l_k] = S[l_k].
                                                            8: end for
      for l = l_k - 1, \dots, l_{k-1} do
         Uniformly randomly sample indices S_m[l]
        from neighbors of S_m[l+1].
      end for
10:
      Send(S_m[l_{k-1}]) if k > 2.
11:
12: end for
13: Output: \{S_m[l]\}_{l=0}^L
```

## Algorithm 6 JointInference (Modified)

```
1: Client:
                                                                                                                         1: Server:
  2: Input: \mathbf{W}_m, \mathcal{D}_m, \{\mathcal{S}_m[l]\}_{l=0}^L
                                                                                                                         2: for l \in \mathcal{I} do
                                                                                                                                   \mathbf{H}[l+1] = \text{Agg}(\mathbf{H}_{1}^{+}[l], \dots, \mathbf{H}_{M}^{+}[l]).
  3: Set \mathbf{H}_{m}[0] = \mathbf{X}_{m}[S_{m}[0]]
  4: for l = 0, ..., L - 1 do
                                                                                                                                  Broadcast \mathbf{H}[l+1].
  5:
           \mathbf{H}_{m}^{+}[l] = \sigma(\mathbf{A}(\mathcal{E}(\mathcal{S}_{m}[l+1], \mathcal{S}_{m}[l]))\mathbf{H}_{m}[l]\mathbf{W}_{m}[l])
                                                                                                                         5: end for
                                                                                                                         6: Receive \nabla_{\mathbf{H}_{1}^{t,0}[L]} \mathcal{L}_{1}^{t,0} from client 1
           if l \in \mathcal{I} then
  6:
                Send \mathbf{H}_{m}^{+}[l] to server
  7:
                                                                                                                         7: Broadcast \nabla_{\mathbf{H}_{1}^{t,0}[L]} \mathcal{L}_{1}^{t,0}.
                Receive \mathbf{H}_m[l+1]
  8:
                \mathbf{H}_{-m}[l+1] = \text{Extract}(\mathbf{H}_m[l+1], \mathbf{H}_m^+[l])
 9:
10:
                Set \mathbf{H}_m[l+1] = \mathbf{H}_m^+[l]
11:
12:
            end if
13: end for
           Compute loss \mathcal{L}_1^{t,0} = \ell\left(\mathbf{y}[\mathcal{S}_1^t[L]], f_1(\mathbf{H}_1^{t,0}[L], \mathbf{W}_1[L])\right)
15:
           Send \nabla_{\mathbf{H}_{1}^{t,0}[L]} \mathcal{L}_{1}^{t,0} to server
16:
17: else
           Receive \nabla_{\mathbf{H}_{1}^{t,0}[L]} \mathcal{L}_{1}^{t,0} from server
18:
19: end if
20: Output: \{\mathbf{H}_{-m}[l+1]\}_{l \in \mathcal{I}}, \nabla_{\mathbf{H}_{1}^{t,0}[L]} \mathcal{L}_{1}^{t,0}
```

## Algorithm 7 LocalUpdate (Modified)

```
1: Input: \mathbf{W}_m^{t,q}, \mathcal{D}_m, \{\mathcal{S}_m^t[l]\}_{l=0}^L, \{\mathbf{H}_{-m}^t[l+1]\}_{l\in\mathcal{I}}, \nabla_{\mathbf{H}_1^{t,0}[L]}\mathcal{L}_1^{t,0}
  2: Set \mathbf{H}_{m}^{t,q}[0] = \mathbf{X}_{m}[\mathcal{S}_{m}^{t}[0]]
  3: for l = 0, \dots, L - 1 do
             \mathbf{H}_{m}^{t,q,+}[l] = \sigma(\mathbf{A}(\mathcal{E}(\mathcal{S}_{m}^{t}[l+1], \mathcal{S}_{m}^{t}[l]))\mathbf{H}_{m}^{t,q}[l]\mathbf{W}_{m}^{t,q}[l])
  5:
                  Set \mathbf{H}_m^{t,q}[l+1] = \mathrm{Agg}(\mathbf{H}_{-m}^t[l+1], \mathbf{H}_m^{t,q,+}[l])
  6:
  7:
                  Set \mathbf{H}_{m}^{t,q}[l+1] = \mathbf{H}_{m}^{t,q,+}[l]
  8:
 9:
10: end for
11: if m = 1 then
             Compute loss \mathcal{L}_{1}^{t,q} = \ell\left(\mathbf{y}[\mathcal{S}_{m}^{t}[L]], f_{m}(\mathbf{H}_{m}^{t,q}[L], \mathbf{W}_{m}[L])\right)
13:
             Compute partial loss based on \nabla_{\mathbf{H}_{m}^{t,0}[L]}\mathcal{L}_{1}^{t,0} by using the chain rule: \mathcal{L}_{m}^{t,q} = \nabla_{\mathbf{H}_{m}^{t,q}[L]}\mathcal{L}_{1}^{t,0}
14:
16: Output: \mathbf{W}_m^{t,q+1} = \mathbf{W}_m^{t,q} - \eta^{t,q} \nabla_{\mathbf{W}_m^{t,q}} \mathcal{L}_m^{t,q}
```

# C Proofs for section 4

## C.1 Assumptions

**Assumption 1** (Smooth function and Lipschitz gradient). The loss function  $\ell$  is  $G_{\ell}$ -smooth with  $L_{\ell}$ -Lipschitz gradient, i.e.,

$$\|\ell(\mathbf{y}, \mathcal{S}, \mathbf{W}) - \ell(\mathbf{y}, \mathcal{S}, \mathbf{W}')\| \le G_{\ell} \|\mathbf{W} - \mathbf{W}'\|$$
$$\|\nabla_{\mathbf{W}}\ell(\mathbf{y}, \mathcal{S}, \mathbf{W}) - \nabla_{\mathbf{W}'}\ell(\mathbf{y}, \mathcal{S}, \mathbf{W}')\| \le L_{\ell} \|\mathbf{W} - \mathbf{W}'\|, \quad \forall \mathbf{W}, \mathbf{W}'$$

and each client's prediction function  $f_m$  is  $G_f$ -smooth with  $L_f$ -Lipschitz gradient, i.e.,

$$\|f_{m}(\mathcal{S}, \mathbf{W}_{m}) - f_{m}(\mathcal{S}, \mathbf{W}'_{m})\| \leq G_{f} \|\mathbf{W}_{m} - \mathbf{W}'_{m}\|$$
$$\|\nabla_{\mathbf{W}_{m}} f_{m}(\mathcal{S}, \mathbf{W}_{m}) - \nabla_{\mathbf{W}'_{m}} f_{m}(\mathcal{S}, \mathbf{W}_{m})\| \leq L_{f} \|\mathbf{W}_{m} - \mathbf{W}'_{m}\|, \quad \forall \ \mathbf{W}_{m}, \mathbf{W}'_{m}, \forall \ m.$$

**Assumption 2** (Lower-bounded objective). The training objective is bounded below; that is, there exists a constant  $\mathcal{L}^* > -\infty$  such that for all  $\{\mathbf{W}_m\}$ , it satisfies that

$$\mathcal{L}(\{\mathbf{W}_m\}) \geq \mathcal{L}^*$$
.

**Assumption 3** (Uniform sampling). At each iteration t, the server and the clients uniformly sample nodes  $\{S_m[l]\}_{l=0}^L$ , with |S[L]| = S, according to Algorithm 2.

#### C.2 Proof of Theorem 1

We first note the following useful relation:

$$\|a+b\|^2 = \|a-c+c-b\|^2 \le (1+\alpha) \|a-c\|^2 + (1+\frac{1}{\alpha}) \|c-b\|^2, \quad \forall \alpha > 0.$$
 (4)

For notation simplicity, let us denote the expectation conditioned on all the information before iteration t as

$$\mathbb{E}^t[~\cdot~] = \mathbb{E}_{\mathcal{S}^t}[~\cdot~|\mathbf{W}^{t-1,Q},\ldots,\mathbf{W}^{0,0},\mathcal{S}^{t-1},\ldots,\mathcal{S}^0];$$

denote the "all-but-m" vector as  $(\cdot)_{-m}$ , (e.g., the collection of all client parameters except for client m is  $\mathbf{W}_{-m} = \{\mathbf{W}_{m'}\}_{m'\neq m}$ ); denote the client model updated with data  $\mathcal{S}$  as  $\mathbf{W}_m(\mathcal{S})$ ; denote the gradient

evaluated with data S on parameter  $\mathbf{W}_m$  as  $\nabla \mathcal{L}(\mathbf{W}_m(S), S)$ ; and denote the stacked gradient of all clients as  $\mathbf{G} = [\nabla \mathcal{L}(\mathbf{W}_1(S), S), \dots, \nabla \mathcal{L}(\mathbf{W}_M(S), S)]$ . Then, the update rule can be rewritten as:

$$\mathbf{W}^{t,q+1}(\mathcal{S}^t) = \mathbf{W}^{t,q}(\mathcal{S}^t) - \eta \mathbf{G}^{t,q}. \tag{5}$$

In addition, let us define a virtual model sequence updated with full data as  $\mathbf{W}(\mathcal{D})$ , i.e.,

$$\mathbf{W}^{t,q+1}(\mathcal{D}) = \mathbf{W}^{t,q}(\mathcal{D}) - \eta \nabla \mathcal{L}(\mathbf{W}^{t,q}(\mathcal{D}), \mathcal{D}). \tag{6}$$

We can bound the variance of the stochastic gradient at any round t and iteration q = 0 with the following lemma:

**Lemma 1** (Bounded variance). Under Assumptions 1–3, with probability at least  $p = 1 - \delta$ , the variance of the stochastic gradient is bounded by:

$$\mathbb{E}^{t} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}; \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}; \mathcal{D}) \right\|^{2} \right] \leq \sigma, \quad \forall \mathbf{W} \text{ independent of } \mathcal{S}^{t}, \tag{7}$$

where

$$\sigma = 64G_{\ell}^2 L_f^2 \log \left(\frac{2d}{\delta}\right) + 128L_{\ell}^2 \left(G_f^4 + \frac{1}{S}\right) \left(\log \left(\frac{2d}{\delta}\right) + \frac{1}{4}\right). \tag{8}$$

The main technique for proving this lemma is to use the matrix Bernstein inequality (Tropp, 2015) to bound the variance of the stochastic gradients and the variance of the expectation for each client. The proof steps of Lemma 1 follows the same steps in the proofs for Lemmas 5 and 6 of Ramezani et al. (2020), so we omit them here.

Further, we bound the Lipschitz constant of the total loss function in the following lemma:

**Lemma 2** (Lipschitz gradient). Under Assumptions 1–3, the full gradient and each partial gradient of the objective  $\mathcal{L}(\mathbf{W}, \mathcal{S})$  are Lipschitz continuous with uniform constant  $C_0 = G_\ell L_f + G_f^2 L_\ell$ :

$$\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}, \mathcal{S}) - \nabla_{\mathbf{W}'} \mathcal{L}(\mathbf{W}', \mathcal{S})\| \le C_0 \|\mathbf{W} - \mathbf{W}'\|, \quad \forall \mathbf{W}, \mathbf{W}'$$
$$\|\nabla_{\mathbf{W}_m} \mathcal{L}(\mathbf{W}, \mathcal{S}) - \nabla_{\mathbf{W}_m'} \mathcal{L}(\mathbf{W}', \mathcal{S})\| \le C_0 \|\mathbf{W} - \mathbf{W}'\|, \quad \forall \mathbf{W}, \mathbf{W}', \forall m.$$

The proof of Lemma 2 is given below in Section C.3.

With the above results, we begin our proof for Theorem 1. First, applying Lemma 2, we have:

$$\mathcal{L}(\mathbf{W}^{t,q+1}, \mathcal{D}) - \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) \leq \left\langle \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}), \mathbf{W}^{t,q+1} - \mathbf{W}^{t,q} \right\rangle + \frac{C_0}{2} \left\| \mathbf{W}^{t,q+1} - \mathbf{W}^{t,q} \right\|^{2}$$

$$\stackrel{(a)}{=} -\eta \left\langle \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}), \mathbf{G}^{t,q} \right\rangle + \frac{C_0 \eta^{2}}{2} \left\| \mathbf{G}^{t,q} \right\|^{2}$$

$$\stackrel{(b)}{=} -\frac{\eta}{2} \left( \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) \right\|^{2} + \left\| \mathbf{G}^{t,q} \right\|^{2} - \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbf{G}^{t,q} \right\|^{2} \right) + \frac{C_0 \eta^{2}}{2} \left\| \mathbf{G}^{t,q} \right\|^{2}$$

$$= -\frac{\eta}{2} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) \right\|^{2} - \frac{\eta}{2} (1 - \eta C_0) \left\| \mathbf{G}^{t,q} \right\|^{2} + \frac{\eta}{2} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbf{G}^{t,q} \right\|^{2}, \tag{9}$$

where step (a) applies the update rule of Algorithm 4 and step (b) uses the fact that  $\langle a,b\rangle=\frac{1}{2}\left(\|a\|^2+\|b\|^2-\|a-b\|^2\right)$ . Taking expectation, we have:

$$\mathbb{E}^{t}[\mathcal{L}(\mathbf{W}^{t,q+1},\mathcal{D}) - \mathcal{L}(\mathbf{W}^{t,q},\mathcal{D})] \leq -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q},\mathcal{D})\|^{2}$$

$$-\frac{\eta}{2}(1 - \eta C_{0}) \mathbb{E}^{t} \|\mathbf{G}^{t,q}\|^{2} + \frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q},\mathcal{D}) - \mathbf{G}^{t,q}\|^{2}$$

$$\stackrel{(a)}{=} -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q},\mathcal{D})\|^{2} + \frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q},\mathcal{D}) - \mathbf{G}^{t,q}\|^{2}$$

$$-\frac{\eta}{2}(1 - \eta C_{0})(\|\mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} + \mathbb{E}^{t} \|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2})$$

$$\stackrel{(b)}{\leq} -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D})\|^{2} - \frac{\eta}{2} (1 - \eta C_{0}) (\|\mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} + \mathbb{E}^{t} \|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2}) 
+ \frac{\eta}{2} \left( (1 + \frac{1}{\eta C_{0}}) \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} + (1 + \eta C_{0}) \mathbb{E}^{t} \|\mathbb{E}^{t} \mathbf{G}^{t,q} - \mathbf{G}^{t,q}\|^{2} \right) 
= -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D})\|^{2} - \frac{\eta}{2} (1 - \eta C_{0}) \|\mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} + \eta^{2} C_{0} \underbrace{\mathbb{E}^{t} \|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2}}_{\text{Term 1}} 
+ \frac{1 + \eta C_{0}}{2C_{0}} \underbrace{\mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2}}_{\text{Term 2}}, \tag{10}$$

where step (a) uses the fact that  $\mathbb{E}(X)^2 = \mathbb{E}(X^2) + \mathbb{E}(X - \mathbb{E}(X))^2$  and step (b) uses (4) with  $\alpha = \eta C_0$ . Next, we bound Term 1 and Term 2 in the above inequality separately.

#### C.2.1 Bound of term 1

First,

we can rewrite  $\mathbb{E}^{t}[\|\mathbf{G}^{t,q} - \mathbb{E}^{t}\mathbf{G}^{t,q}\|^{2}]$  as:

$$\mathbb{E}^{t}[\|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2}] = \sum_{m=1}^{M} \mathbb{E}^{t} \left[ \|\nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \mathbb{E}_{\mathcal{S}} \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}), \mathcal{S}) \|^{2} \right]$$

$$\stackrel{(a)}{\leq} \sum_{m=1}^{M} \mathbb{E}^{t} \left[ \|\nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{D}) \|^{2} \right],$$

$$\stackrel{\triangleq A_{m,q}^{t,q}}{=} (11)$$

where step (a) uses the fact that  $\mathbb{E}(X - \mathbb{E}(X))^2 \leq \mathbb{E}(X - Y)^2$  for all constant Y. Then, we can bound  $A_m^{t,q}$  as follows. When q = 0, by Lemma 1, we obtain that  $A_m^{t,0} \leq \sigma$  holds with probability  $1 - \delta$ . In general, when  $q \geq 1$ , we have:

$$A_{m}^{t,q} \overset{(4)}{\leq} 2 \mathbb{E}^{t} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{S}^{t}) \right\|^{2} \right]$$

$$+ 2 \mathbb{E}^{t} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{D}) \right\|^{2} \right]$$

$$\overset{(a)}{\leq} 2C_{0}^{2} \mathbb{E}^{t} \left[ \left\| \mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}) - \mathbf{W}_{m}^{t,q}(\mathcal{D}) \right\|^{2} \right] + 2 \mathbb{E}^{t} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{D}), \mathcal{D}) \right\|^{2} \right]$$

$$\overset{(b)}{\leq} 2C_{0}^{2} \mathbb{E}^{t} \left[ \left\| \mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}) - \mathbf{W}_{m}^{t,q}(\mathcal{D}) \right\|^{2} \right] + 2\sigma,$$

$$(12)$$

which holds with probability  $1 - \delta$ . Here, step (a) applies Lemma 2 to the first term and step (b) applies Lemma 1 to the second term. Then, we bound  $\mathbb{E}^t \left[ \|\mathbf{W}_m^{t,q}(\mathcal{S}^t) - \mathbf{W}_m^{t,q}(\mathcal{D})\|^2 \right]$  in the above equation as:

$$\mathbb{E}^{t} \left[ \left\| \mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}) - \mathbf{W}_{m}^{t,q}(\mathcal{D}) \right\|^{2} \right]$$

$$\stackrel{(a)}{=} \mathbb{E}^{t} \left[ \left\| \mathbf{W}_{m}^{t,0} - \eta \sum_{q'=0}^{q-1} \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \left( \mathbf{W}_{m}^{t,0} - \eta \sum_{q'=0}^{q-1} \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{D}), \mathcal{D}) \right) \right\|^{2} \right]$$

$$\stackrel{(b)}{=} \eta^{2} \mathbb{E}^{t} \left[ \left\| \sum_{q'=0}^{q-1} \left( \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{D}), \mathcal{D}) \right) \right\|^{2} \right]$$

$$\stackrel{(c)}{\leq} \eta^{2} Q \sum_{q'=0}^{q-1} \mathbb{E}^{t} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q'}(\mathcal{D}), \mathcal{D}) \right\|^{2} \right]$$

$$= \eta^2 q \sum_{q'=0}^{q-1} A_m^{t,q'}, \tag{13}$$

where in step (a) we expand the updates to  $\mathbf{W}_m^{t,0}$  with (5) and (6); step (b) cancels  $\mathbf{W}_m^{t,0}$  and rearrange the terms; and step (c) applies the Cauchy–Schwarz inequality. At this point, we have the following relations:

$$\mathbb{E}^{t}[\left\|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\right\|^{2}] \leq \sum_{m=1}^{M} A_{m}^{t,q}, \quad A_{0}^{t,0} \leq \sigma, \quad A_{m}^{t,q} \leq 2C_{0}^{2}\eta^{2}q \sum_{q'=0}^{q-1} A_{m}^{t,q'} + 2\sigma, \forall \ q \geq 1.$$

Note that  $q \leq Q$ . By choosing  $2\eta^2 C_0^2 Q^2 \leq 1$ , which implies that  $\eta \leq \frac{1}{\sqrt{2}QC_0}$ , and by recursively substituting the terms, we have the following bounds:

$$A_{m}^{t,q} \leq \left[ 2 + 4q^{2}\eta^{2}C_{0}^{2} + \frac{8}{3}q^{3}\eta^{4}C_{0}^{4} \right] \cdot \sigma \leq \frac{14}{3}\sigma,$$

$$\mathbb{E}^{t} [\|\mathbf{G}^{t,q} - \mathbb{E}^{t}\mathbf{G}^{t,q}\|^{2} \leq M \cdot \left[ 2 + 4q^{2}\eta^{2}C_{0}^{2} + \frac{8}{3}q^{3}\eta^{4}C_{0}^{4} \right] \cdot \sigma \leq \frac{14M\sigma}{3}.$$
(14)

This completes bounding the term  $\mathbb{E}[\|\mathbf{G}^{t,q} - \mathbb{E}^t \mathbf{G}^{t,q}\|^2]$ .

#### C.2.2 Bound of term 2

We have the following series of relations:

$$\begin{split} & \mathbb{E}^{t} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbb{E}^{t} \mathbf{G}^{t,q} \right\|^{2} = \sum_{m=1}^{M} \mathbb{E}^{t} \left\| \nabla_{\mathbf{W}_{m}} \mathcal{L}(\mathbf{W}^{t,q}(\mathcal{S}^{t}), \mathcal{D}) - \mathbb{E}_{\mathcal{S}} \nabla \mathcal{L}(\mathbf{W}^{t,q}_{m}(\mathcal{S}), \mathcal{S}) \right\|^{2} \\ & \stackrel{(a)}{\leq} \sum_{m=1}^{M} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left\| \nabla_{\mathbf{W}_{m}} \mathcal{L}(\mathbf{W}^{t,q}(\mathcal{S}^{t}), \mathcal{S}) - \mathbb{E}_{\mathcal{S}} \nabla \mathcal{L}(\mathbf{W}^{t,q}_{m}(\mathcal{S}), \mathcal{S}) \right\|^{2} \\ & \stackrel{(b)}{\leq} \sum_{m=1}^{M} C_{0}^{2} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left\| \mathbf{W}^{t,q}(\mathcal{S}^{t}) - [\mathbf{W}^{t,q}_{m}(\mathcal{S}), \mathbf{W}^{t,0}_{-m}] \right\|^{2} \\ & = \sum_{m=1}^{M} C_{0}^{2} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left[ \left\| \mathbf{W}^{t,q}_{m}(\mathcal{S}^{t}) - \mathbf{W}^{t,q}_{m}(\mathcal{S}) \right\|^{2} + \sum_{m' \neq m} \left\| \mathbf{W}^{t,q}_{m'}(\mathcal{S}^{t}) - \mathbf{W}^{t,0}_{m'} \right\|^{2} \right] \\ & \stackrel{(c)}{=} \eta^{2} \sum_{m=1}^{M} C_{0}^{2} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}), \mathcal{S}) \right\|^{2} \right] \\ & + \sum_{m' \neq m} \left\| \sum_{q'=0}^{q-1} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left[ \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}^{t}), \mathcal{S}^{t}) - \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}), \mathcal{S}) \right\|^{2} \right. \\ & + \sum_{m' \neq m} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}); \mathcal{S}) \right\|^{2} \right] \\ & \stackrel{(e)}{=} \eta^{2} (M+1) C_{0}^{2} q \sum_{m=1}^{M} \sum_{q'=0}^{q-1} \mathbb{E}^{t} \mathbb{E}_{\mathcal{S}} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q'}_{m}(\mathcal{S}^{t}), \mathcal{S}^{t}) \right\|^{2} \\ & \stackrel{(f)}{=} \eta^{2} (M+1) C_{0}^{2} q \sum_{q'=0}^{M} \mathbb{E}^{t} \left\| \mathbb{E}^{t} \left\| \mathbb{G}^{t,q'} \right\|^{2} \right. \end{aligned}$$

$$\stackrel{(g)}{=} \eta^2 (M+1) C_0^2 q \sum_{q'=0}^{q-1} \mathbb{E}^t \left[ \left\| \mathbf{G}^{t,q'} - \mathbb{E}^t \mathbf{G}^{t,q'} \right\|^2 + \left\| \mathbb{E}^t \mathbf{G}^{t,q'} \right\|^2 \right], \tag{15}$$

where step (a) uses Assumption 3, which states that S is uniformly sampled from D, and applies Jensen's inequality, that is

$$\begin{aligned} & \left\| \mathbb{E}_{\mathcal{S}} \, \nabla_{\mathbf{W}_{m}} \mathcal{L}(\mathbf{W}^{t,q}(\mathcal{S}^{t}); \mathcal{S}) - \mathbb{E}_{\mathcal{S}} \, \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}); \mathcal{S}) \right\|^{2} \\ & \leq \mathbb{E}_{\mathcal{S}} \, \left\| \nabla_{\mathbf{W}_{m}} \mathcal{L}(\mathbf{W}^{t,q}(\mathcal{S}^{t}); \mathcal{S}) - \nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}); \mathcal{S}) \right\|^{2}; \end{aligned}$$

step (b) applies Lemma 2 and uses the fact that  $\nabla \mathcal{L}(\mathbf{W}_{m}^{t,q}(\mathcal{S}^{t}), \mathcal{S}^{t})$  is evaluated on  $\mathbf{W}_{m}^{t,q}(\mathcal{S}^{t})$  and  $\mathbf{W}_{-m}^{t,0}$ ; in step (c) we expand the update steps until t,0 with (5); step (d) applies Cauchy-Schwarz inequality; in step (e) we reorder the sum and apply the i.i.d. Assumption 3 to  $\mathcal{S}, \mathcal{S}^{t}$ ; and in step (g) we plug in the definition of  $\mathbf{G}$ . This completes bounding the term  $\mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D}) - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2}$ .

#### C.2.3 Proof of the main result

Substituting the last term in (10) with (15), we obtain that the following holds with probability  $(1 - \delta)^Q$ :

$$\mathbb{E}^{t}[\mathcal{L}(\mathbf{W}^{t,q+1}, \mathcal{D}) - \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D})] \leq -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q}, \mathcal{D})\|^{2} - \frac{\eta}{2} (1 - \eta C_{0}) \|\mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} 
+ \eta^{2} C_{0} \mathbb{E}^{t} \|\mathbf{G}^{t,q} - \mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} 
+ \frac{1 + \eta C_{0}}{2C_{0}} \eta^{2} (M + 1) C_{0}^{2} q \sum_{q'=0}^{q-1} \mathbb{E}^{t} \left[ \|\mathbf{G}^{t,q'} - \mathbb{E}^{t} \mathbf{G}^{t,q'}\|^{2} + \|\mathbb{E}^{t} \mathbf{G}^{t,q'}\|^{2} \right] 
\leq -\frac{\eta}{2} \mathbb{E}^{t} \|\nabla \mathcal{L}(\mathbf{W}^{t,q})\|^{2} - \frac{\eta}{2} (1 - \eta L) \|\mathbb{E}^{t} \mathbf{G}^{t,q}\|^{2} 
+ \frac{1 + \eta C_{0}}{2C_{0}} \eta^{2} (M + 1) C_{0}^{2} q \sum_{q'=0}^{q-1} \mathbb{E}^{t} \|\mathbb{E}^{t} \mathbf{G}^{t,q'}\|^{2} 
+ \eta^{2} C_{0} \cdot \left( 1 + \frac{(1 + \eta C_{0}) \cdot (M + 1) \cdot \eta Q^{2}}{2} \right) \cdot \frac{14M\sigma}{3},$$

where in the second inequality, we set  $\eta \leq \frac{1}{\sqrt{2}QC_0}$ , plug in (14), and use the fact that  $q \leq Q$ . Averaging over  $t = 0, \ldots, T-1$  and  $q = 0, \ldots, Q-1$  and reorganizing the terms, we obtain:

$$\frac{1}{TQ} \sum_{t=0}^{T-1} \sum_{q=0}^{Q-1} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}; \mathcal{D}) \right\|^{2} \leq \frac{2}{\eta TQ} \mathbb{E} \left[ \mathcal{L}(\mathbf{W}^{0}) - \mathcal{L}(\mathbf{W}^{T,Q}) \right] \\
- \frac{1 - \eta C_{0} \left( 1 + (1 + \eta C_{0}) \cdot (M+1) \cdot Q^{2} \right)}{TQ} \sum_{t=0}^{T-1} \sum_{q=0}^{Q-1} \mathbb{E} \left\| \mathbb{E}^{t} \mathbf{G}^{t,q} \right\|^{2} \\
+ 2\eta C_{0} \cdot \left( 1 + \frac{(1 + \eta C_{0}) \cdot (M+1) \cdot \eta Q^{2}}{2} \right) \cdot \frac{14M\sigma}{3},$$

which holds with probability at least  $(1 - \delta)^{TQ}$ . Let  $\delta = \delta'/TQ \in (0, 1)$ ; then, the above equation holds with probability at least

$$(1-\delta'/TQ)^{TQ} \ge 1-\delta'/TQ \times TQ = 1-\delta'.$$

Let

$$1 - \eta C_0 \left( 1 + (1 + \eta C_0) \cdot (M+1) \cdot Q^2 \right) \ge 0,$$

 $(\eta \leq \frac{1}{C_0 \cdot (1+2MQ^2)})$  and apply Assumption 2. Then, we have

$$\frac{1}{TQ} \sum_{t=0}^{T-1} \sum_{q=0}^{Q-1} \mathbb{E} \left\| \nabla \mathcal{L}(\mathbf{W}^{t,q}; \mathcal{D}) \right\|^2 \le \frac{2(\mathcal{L}(\mathbf{W}^0) - \mathcal{L}^*)}{\eta TQ} + \frac{28\eta M \cdot \left( C_0 + \sqrt{M+1}Q \right)}{3} \sigma, \tag{16}$$

which holds with probability at least  $1 - \delta$ , where

$$\sigma = 64G_{\ell}^2 L_f^2 \log \left( \frac{2dTQ}{\delta} \right) + 128L_{\ell}^2 \left( G_f^4 + \frac{1}{S} \right) \left( \log \left( \frac{2dTQ}{\delta} \right) + \frac{1}{4} \right).$$

This completes the proof of Theorem 1.

# C.3 Proof for Lemma 2

In this subsection, we prove

$$\|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) - \nabla_{\mathbf{W}'} \mathcal{L}(\mathbf{W}')\| \le C_0 \|\mathbf{W} - \mathbf{W}'\|$$

and

$$\|\nabla_{\mathbf{W}_m} \mathcal{L}(\mathbf{W}) - \nabla_{\mathbf{W}_m'} \mathcal{L}(\mathbf{W}')\| \le C_0 \|\mathbf{W} - \mathbf{W}'\|.$$

Note that  $\nabla_{\mathbf{W}_m} \mathcal{L}(\mathbf{W})$  is a sub-vector of  $\nabla \mathcal{L}(\mathbf{W}')$ , so  $\|\nabla_{\mathbf{W}_m} \mathcal{L}(\mathbf{W}) - \nabla_{\mathbf{W}_m'} \mathcal{L}(\mathbf{W}')\| \le \|\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) - \nabla_{\mathbf{W}'} \mathcal{L}(\mathbf{W}')\|$ . Therefore, we only need to prove the first inequality.

The gradient  $\nabla \mathcal{L}(\mathbf{W})$  can be expanded as

$$\nabla \mathcal{L}(\mathbf{W}) = \nabla \ell(\mathbf{y}, f_m(\mathbf{S}, \mathbf{W}))$$

$$= \nabla \ell(f_m(\mathbf{S}, \mathbf{W})) \cdot \nabla_{\mathbf{W}} f_m(\mathbf{S}, \mathbf{W}) = \nabla \ell(f_m) \cdot \nabla f_m(\mathcal{W}),$$
(17)

where in the last equation we omit the irrelevant variables. Then, we have

$$\|\nabla \mathcal{L}(\mathbf{W}) - \nabla \mathcal{L}(\mathbf{W}')\| = \|\nabla \ell(f_m) \cdot \nabla f_m(\mathbf{W}) - \nabla \ell(f'_m) \cdot \nabla f_m(\mathbf{W}')\|$$

$$= \|\nabla \ell(f_m) \cdot (\nabla f_m(\mathbf{W}) - \nabla f_m(\mathbf{W}')) + (\nabla \ell(f_m) - \nabla \ell(f'_m)) \cdot \nabla f_m(\mathbf{W}')\|$$

$$\stackrel{(a)}{\leq} \|\nabla \ell(f_m) \cdot (\nabla f_m(\mathbf{W}) - \nabla f_m(\mathbf{W}'))\| + \|(\nabla \ell(f_m) - \nabla \ell(f'_m)) \cdot \nabla f_m(\mathbf{W}')\|$$

$$\stackrel{(b)}{\leq} \|\nabla \ell(f_m)\| \|\nabla f_m(\mathbf{W}) - \nabla f_m(\mathbf{W}')\| + \|\nabla \ell(f_m) - \nabla \ell(f'_m)\| \|\nabla f_m(\mathbf{W}')\|$$

$$\stackrel{(c)}{\leq} G_{\ell} L_f \|\mathbf{W} - \mathbf{W}'\| + L_{\ell} \|f_m(\mathbf{W}) - f_m(\mathbf{W}')\| \cdot G_f$$

$$\stackrel{A1}{\leq} (G_{\ell} L_f + L_{\ell} G_f^2) \cdot \|\mathbf{W} - \mathbf{W}'\|,$$

$$(18)$$

where step (a) uses the fact that  $||a+b|| \le ||a|| + ||b||$ , step (b) uses the fact that  $||ab|| \le ||a|| ||b||$ ; and in step (c) we apply Lemma 2 (that is, for any G-smooth function g, its gradient is bounded as  $||\nabla g|| \le G$ ) to the first and the fourth terms and Lipschitz gradient to the second and the third terms. This completes the proof of Lemma 2.

#### D Experiment details

#### D.1 Details of the datasets

Planetoid (Yang et al., 2016): This collection contains three citation datasets: Cora, PubMed, and CiteSeer. Each dataset contains one citation graph, where the nodes represent papers and edges represent citations. The node features are a bag of words and the classification target is the paper category. In the experiment, each client holds a non-overlapping block of node features and a subgraph that results from uniformly sampling 80% of the edges.

HeriGraph (Bai et al., 2022): This collection contains three multi-modal graph datasets, each of which is constructed from heritage data posted on social media for a particular city (Suzhou, Amsterdam, and Venice). Each post contains user information, timestamp, geolocation, image, and text annotation. The posts are connected to form three subgraphs: a social subgraph, a spatial subgraph, and a temporal subgraph. The social subgraph is formed based on friendship and common-interest relations of the users. The spatial

subgraph is formed based on the spatial proximity of the geolocations. The temporal subgraph is formed based on the temporal proximity of the posts. Each post has three blocks of image features and possibly text features; for classification, it belongs to one of nine heritage attributes. In the experiment, each client holds one of the three subgraphs and one of the three image feature blocks.

**Reddit** (Hamilton et al., 2017): Reddit is a large online community where users post and comment on different topics. Each node represents a post and the features are the text of the post. Two posts are connected if the same user comments on both. The classification target is the community (subreddit) that a post belongs to. Similar to Planetoid, in the experiment, each client holds a non-overlapping block of node features and a subgraph that results from uniformly sampling 80% of the edges.

#### D.2 Details of the experiments

Here we provide a list of hyperparameters for grid search in Table 6. The optimal set of hyperparameters for each setting is tuned according to the range listed in the table. A summary of the details of the datasets is given in Table 1.

Table 6: Hyperparameter grid search range for the numerical experiments.

Hyperparameter	Grid search range
Hidden dimension of $\mathbf{H}[l]$ Batch size $S$ Neighborhood sample size Training rounds $T$	{ 128, 192, 256, 384} { 16, 32} { 2, 3, 4, 6, 8} { 512, 640, 1024, 1152, 3200}
Learning rate $\eta$	$\left  \{1, 2, 3.5, 5, 7, 8\} \times \{10^{-1}, 10^{-2}, 10^{-3}\} \right $