

SWIFT-FEDGNN: FEDERATED GRAPH LEARNING WITH LOW COMMUNICATION AND SAMPLE COMPLEXITIES

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph neural networks (GNNs) have achieved great success in a wide variety of graph-based learning applications. To expedite training for large-scale graphs, distributed GNN training has been proposed using sampling-based mini-batch training. However, such a traditional distributed GNN training approach is not applicable to emerging GNN learning applications with geo-distributed input graphs, which require the data to be kept within the site where it is generated to protect privacy. On the other hand, federated learning (FL) has been widely used to enable privacy-preserving training under data parallelism. However, because of cross-client links in the aforementioned geo-distributed graph data, applying federated learning directly to GNNs incurs expensive cross-client neighbor sampling and communication costs due to the large graph size and the dependencies between nodes among different clients. To overcome these challenges, we propose a new mini-batch and sampling-based federated GNN algorithmic framework called Swift-FedGNN that primarily performs efficient parallel local training and periodically conducts time-consuming cross-client training. Specifically, in Swift-FedGNN, each client *primarily* trains a local GNN model using only its local graph data, and some randomly sampled clients *periodically* learn the local GNN models based on their local graph data and the dependent nodes across clients. We theoretically establish the convergence performance of Swift-FedGNN and show that it enjoys a convergence rate of $\mathcal{O}(T^{-1/2})$, matching the state-of-the-art (SOTA) rate of sampling-based GNN methods, despite operating in the challenging FL setting. Extensive experiments on real-world datasets show that Swift-FedGNN significantly outperforms the SOTA federated GNN approaches with comparable accuracy in terms of efficiency.

1 INTRODUCTION

1) Background and Motivation: Graph neural networks (GNNs) have received increasing attention in recent years and have been widely used across various applications, such as social networks (Deng et al., 2019; Qiu et al., 2018; Wang et al., 2019a), recommendation systems (Ying et al., 2018; Wang et al., 2019b;d), traffic prediction (Cui et al., 2019; Kumar et al., 2019; Li et al., 2019), drug discovery (Wang et al., 2022b; Do et al., 2019; Fout et al., 2017), and disease prediction (Ghorbani et al., 2022; Kazi et al., 2023; Li & Zhang, 2024). GNN learns the high-level graph representations by iteratively aggregating the neighboring features of each node, which is then used for downstream tasks, such as node classification (Kipf & Welling, 2017; Hamilton et al., 2017), link prediction (Yao et al., 2023b; Zhang & Chen, 2018), and graph classification (Zhang et al., 2018; Bacciu et al., 2018).

However, real-world graph datasets can be extensive in scale (*e.g.*, Microsoft Academic Graph (Wang et al., 2020) with over 100 million nodes) and generated in a geo-distributed fashion (Yao et al., 2023a). Similar to traditional datasets (*e.g.*, images), graph datasets may be collected across multiple geo-distributed sites/devices and stored locally. Collecting the entire dataset onto a single site/device not only incurs prohibitively high communication costs, but may also violate data protection regulations. Unlike the traditional datasets where data is mutually independent, the nodes in graph data are usually dependent (shown as the links between nodes in Figure 1). Such large-scale real-world graph datasets often exceed the memory and computational capabilities of a single device (*e.g.*, GPU), which leads to a *compelling need* for developing distributed GNN training with multiple devices or machines (Fey & Lenssen, 2019; Zheng et al., 2020). A common paradigm in distributed GNN training involves subgraph sampling (Zeng et al., 2020) and mini-batch training on each device (Luo et al., 2022). In

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

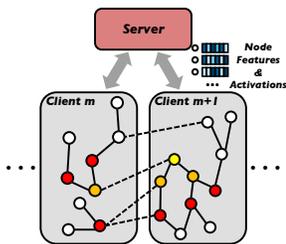


Figure 1: Federated GNN setting. Dashed lines show graph dependency cross-clients.

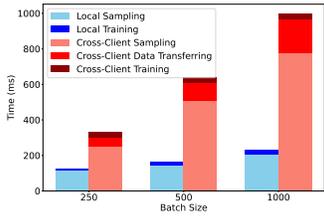


Figure 2: Per-iteration time breakdown: local vs. cross-client training.

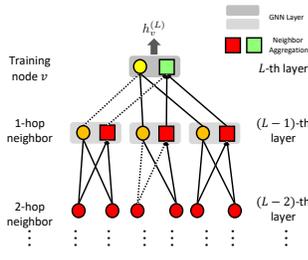


Figure 3: Federated GNN model. Dashed lines show communication between clients.

this approach, workers on each device select training nodes, perform cross-device sampling to gather multi-hop neighbor features as subgraphs, construct subgraphs as mini-batches, and then train on these mini-batches in parallel. However, due to the aforementioned data privacy constraints, this distributed GNN training paradigm is not directly applicable to geo-distributed graphs, as features *cannot* be directly shared among clients.

Meanwhile, federated learning (FL) (McMahan et al., 2017; Yang et al., 2021; Karimireddy et al., 2020), which has emerged as a promising learning paradigm, enables collaborative training of a model using geo-distributed traditional datasets under the coordination of a central server. However, applying FL to geo-distributed graph data is highly non-trivial due to the dependencies between the nodes in a graph and the fact that the neighbors of the node may be located on different clients, which we refer to as “*cross-client neighbors*” (shown as the dashed links between nodes in Figure 1). Ignoring the cross-client neighbors as in (Wang et al., 2022a; He et al., 2021b) would degrade the performance of the models and prevent them from reaching the same accuracy as the models trained on a single device/machine, which is due to the information loss of the cross-client neighbors.

2) Technical Challenges: A naive solution to federated GNN is to leverage the server as an intermediary to perform subgraph sampling and part of the training operation (*i.e.*, neighbor aggregation) to protect data privacy. For example, consider a country that has many hospitals and one medical administrative center and needs to investigate a healthcare problem in the whole country (*e.g.*, infection prediction) (Zhang et al., 2021). The residents may go to different hospitals for healthcare because of various reasons, *e.g.*, the locations of the hospitals. Their healthcare data (*e.g.*, personal information, and patient interactions) would be stored locally only at the hospitals they visit, and such data *cannot* be directly shared among different hospitals due to privacy concerns and conflicts of interest. Note that in a graph, the patients here are the nodes and the patient interactions are the edges, and the graphs located at different hospitals may have cross-hospital edges and thus have cross-hospital neighbors. In this situation, the medical administrative center can serve as the server in federated GNN because it is trusted and thus has access to the graph data located at different hospitals.

However, this method introduces significant sampling and communication overhead (as shown in Figure 1), as the server needs to communicate with all clients to perform subgraph sampling and neighbor aggregation for each client sequentially. Figure 2 illustrates the per-iteration time breakdown of local training (*i.e.*, training only on the local graph of the client) versus cross-client training (*i.e.*, training using both the client’s local graph and the cross-client neighbors) on the Amazon product co-purchasing dataset (Leskovec et al., 2007) ¹. As observed from the figure, cross-client sampling and data communication time dominate the total time for cross-client training, making it *five times* slower than local training. Therefore, it is critical to mitigate the *communication overhead* in cross-client training to enable efficient federated learning on GNNs.

While some prior works ignore the information of the cross-client neighbors (He et al., 2021a) or assume overlapping nodes between different clients (Wu et al., 2021) (may not hold for geo-distributed graphs), other prior works (Zhang et al., 2021; Du & Wu, 2022; Yao et al., 2023a) address the information loss of the cross-client neighbors by facilitating the exchange of such information between clients. However, these approaches may lead to significant sampling and communication

¹Figure 2 uses a two-layer GNN, with sampling fanout values being 15 and 10 for the two layers, and the network bandwidth being 1 Gbps. We ignore model synchronization time, as the model (of size 0.3 MB) takes less than 10 ms to synchronize.

108 overhead, which is due to cross-client sampling and cross-client neighbor information transferring.
 109 In addition, they may impose heavy memory burdens, which is attributed to the information (graph
 110 structure and node features) storage of the cross-client neighbors on the clients. (see detailed
 111 discussions in Section 2).

112 **3) Our Contributions:** The key contribution of this paper is that, by addressing the above challenges,
 113 we develop a mini-batch-based and sampling-based federated GNN framework called Swift-FedGNN.
 114 The main results and technical contributions of this paper are as follows:

- 115 • We develop a new communication- and sample-efficient mini-batch sampling-based federated GNN
 116 algorithm called Swift-FedGNN to train GNNs on geo-distributed graphs in a federated fashion.
 117 To reduce the sampling and communication overhead, the clients in Swift-FedGNN *primarily*
 118 conduct the efficient local training in parallel and some sampled clients *only occasionally and*
 119 *periodically* perform the time-consuming cross-client training. The information loss of the cross-
 120 client neighbors in the federated setting is alleviated via cross-client training. Thanks to our use of
 121 different mini-batch training nodes at each iteration, the clients do *not* need to store the cross-client
 122 neighbor information, which significantly reduces the memory overhead. To further reduce the
 123 communication cost, the cross-client neighbor information is aggregated at remote clients before
 124 communicating to the server and accumulated one more time before transferring to the training
 125 client. This special design further offers the benefit of helping preserve data privacy since the
 126 information of each node is not leaked.
- 127 • We conduct rigorous theoretical convergence performance analysis for Swift-FedGNN. It is worth
 128 noting that the convergence analysis of our Swift-FedGNN is highly non-trivial. Unlike deep
 129 neural networks (DNNs), the stochastic gradients in GNNs are *biased*, which poses significant
 130 challenges on the theoretical analysis of Swift-FedGNN’s convergence guarantees. Moreover, the
 131 structural entanglement in GNNs (*i.e.*, the interleaving of neighbor aggregations and non-linear
 132 transformations across multiple layers) further complicates the performance analysis. In stark
 133 contrast to existing works in the literature that made strong assumptions on the biases of stochastic
 134 gradients (e.g., the unbiased stochastic gradient assumption in (Chen et al., 2018) and the consistent
 135 stochastic gradient assumption in (Chen & Luss, 2018), etc.), for the *first time* in the literature,
 136 we are able to *bound* the stochastic gradient approximation errors rather than resorting to these
 137 unrealistic assumptions in practice. Such results could also be of independent theoretical interests.
- 138 • Given the *biased* stochastic gradients in GNNs that arise from the missing cross-client neighbors
 139 and the neighbor sampling process, we reveal an interesting theoretical insight that the stochastic
 140 gradient approximation errors are correlated with the structure of GNNs. More specifically, our
 141 theoretical analysis quantifies and characterizes a *positive correlation* with the number of layers in
 142 the networks. We note that this is a new finding that is unique to federated GNN training. Lastly,
 143 by putting the above insights together, we show that Swift-FedGNN achieves a convergence rate of
 144 $\mathcal{O}(T^{-1/2})$, which *matches* the state-of-the-art (SOTA) convergence rate of sampling-based GNN
 145 methods (hence low communication and sample complexities), *despite* operating in the far more
 146 challenging FL setting with much less frequent information exchanges among the clients.

147 2 RELATED WORK

148 In this section, we provide an overview on distributed GNNs and offer a comprehensive comparison
 149 with the most relevant work on federated GNNs.

150 **1) Distributed Graph Neural Networks:** Distributed GNN training framework (*e.g.*, DGL’s Dist-
 151 DGL (Wang et al., 2019c; Zheng et al., 2020), Pytorch Geometric (Fey & Lenssen, 2019), Ali-
 152 Graph (Zhao et al., 2019) and Dorylus (Thorpe et al., 2021)) have been developed to train large-scale
 153 graph datasets that exceed the storage capacities of a single device. Each worker on the device
 154 constructs mini-batches via cross-device sampling and communication, trains in parallel, and syn-
 155 chronizes the model. However, in distributed GNN training, extensive graph sampling and data
 156 communication can account for up to 80% of the total training time, substantially slowing the training
 157 process (Gandhi & Iyer, 2021). Considerable efforts have been made to optimize distributed GNN
 158 training, including employing strategic graph partitioning to minimize edge cuts between graph
 159 partitions (Zheng et al., 2020), implementing static or dynamic node feature caching (Liu et al., 2023;
 160 Zhang et al., 2023), enhancing communication strategies (Cai et al., 2021; Luo et al., 2022), and
 161 utilizing various parallel training schemes (Gandhi & Iyer, 2021; Wan et al., 2022; Du et al., 2024).
 However, the majority of these techniques are not directly applicable to geo-distributed graphs due

to privacy concerns, as they typically involve operations that require access to or the transfer of graph data across different training devices. To our knowledge, LLCG (Ramezani et al., 2022) is the only distributed GNN training framework that avoids transferring node features between workers, making it potentially applicable to geo-distributed graphs. Every worker in LLCG trains only on its local graph partitions. To address the missing information from cross-device neighbors, LLCG employs central server to periodically perform full-neighbor training with neighbor aggregation over all workers. However, this method suffers significant communication overhead on the server end, as the server needs to communicate with all workers to perform the full-neighbor training.

2) Federated Graph Neural Networks: To date, the research on federated GNNs remains in its infancy and results in this area are quite limited. In (He et al., 2021a), it is assumed that graphs are dispersed across multiple clients and the information of the cross-client neighbors is ignored, which does not align with the real-world scenarios and would degrade the performance of the trained model. In (Wu et al., 2021), it is assumed that the clients’ local graphs have overlapped nodes and the edges are distributed, which may not be true in real-world situations. The authors in (Zhang et al., 2021) mitigate the information loss of the cross-client neighbors by exchanging such information in each training round. However, this approach incurs considerable communication overhead and exposes private node information to other clients. Although the algorithm in (Yao et al., 2023a) exchanges the information of the full cross-client neighbors only once before the training to supplement the missing information from the cross-client neighbors, it is not applicable to large-scale graphs because it uses the full graph for training and incurs significant memory overhead on each client. Note that this one-time communication only works for full graph training and is not suitable for the situations in which clients sample different mini-batches of training nodes in each iteration since the cross-client neighbors of these training batches are different. Furthermore, the homomorphic encryption used in (Yao et al., 2023a) would significantly increase the communication cost, which is at least several times higher than communication without homomorphic encryption.

The most related work to ours can be found in (Du & Wu, 2022), where the authors used sparse cross-client neighbor sampling to supplement the lost information of the cross-client neighbors and reduce the communication overhead. Each client periodically samples the cross-client neighbors and exchanges the information of the sampled neighbors with other clients. In the remaining iterations, the clients reuse the most recent sampled cross-client neighbors, which requires additional cache for saving transferred graph data, thereby increasing memory overhead. However, as training progresses, the frequency of information exchange increases, leading to higher communication costs. Additionally, they relaxed the privacy constraint to allow the transfer of the graph data between clients directly. Reusing the same sampled data across multiple iterations would cause additional bias, and thus degrading the performance of the model. In contrast, our proposed Swift-FedGNN method limits cross-client training to a *subset* of sampled clients and avoids direct graph data exchange between clients by offloading certain operations to the central server. Before communication with the training clients, cross-client neighbor information is aggregated twice: first at the remote clients and then on the server—helping to preserve data privacy and significantly reduce communication costs. Since clients perform local training in the remaining iterations, cross-client neighbor information does *not* need to be stored, thereby reducing memory overhead.

3 FEDERATED GRAPH LEARNING: PRELIMINARIES

In this section, we provide the background of the mathematical formulation for training GNNs in a federated setting. For convenience, we provide a list of key notations used in this paper in Appendix A. In order for this paper to be self-contained and to facilitate easy comparisons, we provide the background for training GNNs on a single machine in Appendix B.

Consider a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, where \mathcal{V} is a set of nodes with $N = |\mathcal{V}|$ and \mathcal{E} is a set of edges. We consider a standard federated setting that has a central server and a set of \mathcal{M} clients with $M = |\mathcal{M}|$. The graph \mathcal{G} is geographically distributed over these clients, and each client m contains a subgraph represented by $\mathcal{G}^m(\mathcal{V}^m, \mathcal{E}^m)$. Note that $\bigcup_{m=1}^M \mathcal{G}^m \neq \mathcal{G}$ due to the missing cross-client edges between clients ($\bigcup_{m=1}^M \mathcal{E}^m \neq \mathcal{E}$). In addition, we assume that the nodes are *disjointly* partitioned across clients, *i.e.*, $\bigcup_{m=1}^M \mathcal{V}^m = \mathcal{V}$ and $\bigcap_{m=1}^M \mathcal{V}^m = \emptyset$. Each node $v \in \mathcal{V}^m$ has a feature vector $\mathbf{x}_v^m \in \mathbb{R}^d$, and each node $v \in \mathcal{V}_{train}^m$ corresponds to a label y_v^m , where $\mathcal{V}_{train}^m \subseteq \mathcal{V}^m$.

In federated GNN training, the clients collaboratively learn a model with distributed graph data and under the coordination of the central server. Typically, clients in FL receive the model from the server, compute local model updates iteratively, and then send the updated model to the server. The server periodically aggregates the models and then sends the aggregated model back to the clients. The goal in federated GNN training is to solve the following optimization problem:

$$\min \mathcal{L}(\boldsymbol{\theta}) := \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} F^m(\boldsymbol{\theta}) = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \frac{1}{|\mathcal{V}_B^m|} \sum_{v \in \mathcal{V}_B^m} \ell^m \left(\mathbf{h}_v^{(L),m}, y_v^m \right), \quad (1)$$

where ℓ^m is a loss function (e.g., cross-entropy loss) at client m , \mathcal{V}_B^m denotes a mini-batch of training nodes uniformly sampled from \mathcal{V}^m , and $\boldsymbol{\theta} := \{\mathbf{W}^{(l)}\}_{l=1}^L$ corresponds to all model parameters.

GNNs aim to generate representations (embeddings) for each node in the graph by combining information from its neighboring nodes. Recall that in federated GNNs, the neighbors of node v may be located on its local client $m(v)$ or on remote clients $\bar{m}(v) \in \bar{\mathcal{M}}(v)$, where $\bar{\mathcal{M}}(v)$ represents a set of the remote clients that host the neighbors of node v , and $\mathcal{M}(v) \subseteq \mathcal{M} \setminus \{m(v)\}$. As shown in Figure 3, to compute the embedding of node v at the l -th layer in a GNN with L layers, the client $m(v)$ first aggregates the neighbor information from both itself and the remote clients $\bar{m}(v) \in \bar{\mathcal{M}}(v)$, and then updates the embedding of node v , as follows:

$$\mathbf{h}_{\mathcal{N}(v)}^{(l)} = \text{AGG} \left(\left\{ \mathbf{h}_u^{(l-1),m(v)} \mid u \in \mathcal{N}^{m(v)}(v) \right\} \cup \left\{ \bigcup_{\bar{m}(v) \in \bar{\mathcal{M}}(v)} \left\{ \mathbf{h}_u^{(l-1),\bar{m}(v)} \mid u \in \mathcal{N}^{\bar{m}(v)}(v) \right\} \right\} \right),$$

$$\mathbf{h}_v^{(l),m(v)} = \sigma \left(\mathbf{W}^{(l)} \cdot \text{COMBINE}(\mathbf{h}_v^{(l-1),m(v)}, \mathbf{h}_{\mathcal{N}(v)}^{(l)}) \right), \quad (2)$$

where $\mathcal{N}^{m(v)}(v)$ is a set of the neighbors of node v located on its local client $m(v)$, $\mathcal{N}^{\bar{m}(v)}(v)$ is a set of the neighbors of node v located on remote client $\bar{m}(v)$, $\mathbf{h}_{\mathcal{N}(v)}^{(l)}$ is the aggregated embedding from node v 's neighbors, $\mathbf{h}_v^{(l),m(v)}$ is the embedding of node v located on client $m(v)$ and is initialized as $\mathbf{h}_v^{(0),m(v)} = \mathbf{x}_v^{m(v)}$, $\mathbf{W}^{(l)}$ represents the weight matrix at l -th layer, $\sigma(\cdot)$ corresponds to an activation function (e.g., ReLU), $\text{AGG}(\cdot)$ is an aggregation function (e.g., mean), and $\text{COMBINE}(\cdot)$ is a combination function (e.g., concatenation). Compared to distributed GNNs where clients can directly transfer node features, the *key difference* in federated GNNs is that clients *cannot* do so due to privacy concerns, requiring additional modifications.

4 THE Swift-FedGNN ALGORITHM

In this section, we propose a new algorithmic framework called Swift-FedGNN, designed to efficiently solve Problem (1) by reducing both sampling and communication costs in federated GNN training. The overall algorithmic framework of Swift-FedGNN is illustrated in Algorithms 1-3. Rather than each client performing cross-client training in every round, the clients in Swift-FedGNN primarily conduct the *efficient* local training in parallel, and a set of randomly selected clients periodically carry out the time-consuming cross-client training. By offloading part of the graph operation to the server and remote clients, Swift-FedGNN eliminates the need for sharing graph features among clients.

Algorithm 1 outlines the main framework of Swift-FedGNN. Specifically, it performs parallel local training across clients for every $I - 1$ iterations, followed by one iteration of cross-client training involving randomly selected clients. In the local training iterations (t), every client m updates the local GNN model only using its local graph, as presented in Algorithm 3. Client m samples a mini-batch of training nodes \mathcal{B}_v^m and a subset of L -hop neighbors for the training nodes in \mathcal{B}_v^m , denote as $\tilde{\mathcal{S}} = \{\tilde{\mathcal{S}}^{(l)}\}_{l=0}^{L-1}$, all from the local graph data. To compute the embedding of node v in the l -th GNN layer ($v \in \mathcal{B}_v^m$ if $l = L$, otherwise $v \in \tilde{\mathcal{S}}^{(l)}$), client m first conducts the neighbor aggregation for node v based on the sampled neighbors using:

$$\tilde{\mathbf{h}}_{\tilde{\mathcal{N}}(v)}^{(l)} = \text{AGG} \left(\left\{ \tilde{\mathbf{h}}_u^{(l-1),m} \mid u \in \tilde{\mathcal{N}}^m(v) \right\} \right), \quad (3)$$

where $\tilde{\mathcal{N}}^m(v)$ denotes a set of the sampled neighbors located on client m for node v , $\tilde{\mathcal{N}}^m(v) \subseteq \tilde{\mathcal{S}}^{(l-1)}$, and $\tilde{\mathcal{N}}^m(v) \subseteq \mathcal{N}^m(v)$. Then, client m updates the embedding of node v in the l -th GNN layer based on the aggregated neighbor information and the embedding of node v from the $(l-1)$ -th layer, as follows:

Algorithm 1: Swift-FedGNN Algorithm.

Input: Initial parameters θ_0 , learning rate α , and correction frequency I

for $t = 0$ **to** $T - 1$ **do**

if $t \bmod I = 0$ **then**

 Randomly sample $|\mathcal{K}|$ clients

for $m \in \mathcal{M}$ **in parallel do**

if $m \in \mathcal{K}$ **then**

 Client update with local graph data and cross-client neighbors using Algorithm 2

else

 Client update with local graph data according to Algorithm 3

else

for $m \in \mathcal{M}$ **in parallel do**

 Client update with local graph data based on Algorithm 3

Server:

 Aggregate and update global model parameter as:

$\theta_{t+1} = \theta_t - \alpha \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\theta_t^m)$

Algorithm 2: Client m in the t -th iteration: update with local graph data and cross-client neighbors.

Receive global parameter $\theta_t^m = \theta_t$

Construct a mini-batch \mathcal{B}_v^m of nodes

Server samples a subset of L -hop neighbors $\tilde{\mathcal{S}} = \{\tilde{\mathcal{S}}^{(l)}\}_{l=0}^{L-1}$ for the training nodes in \mathcal{B}_v^m

for $l = 1$ **to** L **do**

 /* Derive l -th layer embedding of node $v \in \mathcal{B}_v^m$ if $l = L$, otherwise $v \in \tilde{\mathcal{S}}^{(l)}$ */

for Remote client $\tilde{m}(v) \in \tilde{\mathcal{M}}(v)$ in parallel do

 Aggregate the neighbor embeddings using Eq. (5)

 Send the aggregated embedding $\tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l), \tilde{m}(v)}$ to the server

Server:

 Aggregate the neighbor embeddings from the remote clients using Eq. (6)

 Send the aggregated cross-client neighbor embedding $\tilde{\mathbf{r}}_{\mathcal{N}(v)}^{(l)}$ to Client $m(v)$

Client $m(v)$: Compute node embeddings using Eq. (7) and (8)

 Compute the stochastic gradient as $\nabla \tilde{F}^m(\theta_t^m)$ and send to the server

Algorithm 3: Client m in the t -th iteration: update with local graph data.

Receive global parameter $\theta_t^m = \theta_t$

Construct a mini-batch \mathcal{B}_v^m of nodes

Sample a subset of L -hop neighbors $\tilde{\mathcal{S}} = \{\tilde{\mathcal{S}}^{(l)}\}_{l=0}^{L-1}$ for the training nodes in \mathcal{B}_v^m

for $l = 1$ **to** L **do**

 /* Derive l -th layer embedding of node $v \in \mathcal{B}_v^m$ if $l = L$, otherwise $v \in \tilde{\mathcal{S}}^{(l)}$ */

 Compute node embeddings using Eq. (3) and (4)

 Compute the stochastic gradient $\nabla \tilde{F}^m(\theta_t^m)$ and send to the server

$$\tilde{\mathbf{h}}_v^{(l), m} = \sigma \left(\mathbf{W}_t^{(l), m} \cdot \text{COMBINE}(\tilde{\mathbf{h}}_v^{(l-1), m}, \tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l)}) \right). \quad (4)$$

At every I -th iteration, Swift-FedGNN allows a set of K clients, uniformly sampled from \mathcal{M} , to conduct cross-client training that trains the local GNN models using both their local graph data and the cross-client neighbors. We use \mathcal{K} to denote the set of K clients, where $\mathcal{K} \subset \mathcal{M}$. The remaining clients perform local training as shown in Algorithm 3. Algorithm 2 details the cross-client training process for client $m \in \mathcal{K}$. Rather than directly exchanging node features between clients, Swift-FedGNN partitions GNN training between the clients and the server. We offload² the aggregation of node features and intermediate activations at each GNN layer to the server and remote clients corresponding to node v , thus reducing the communication overhead and eliminating the need for graph data sharing. This procedure helps preserve data privacy because the clients are unaware of the locations of neighbor nodes, and the embeddings of these neighbor nodes are aggregated before being transmitted to the clients. Operations performed on the server and the remote clients are colored using `server` and `remote client` respectively.

Specifically, client $m \in \mathcal{K}$ samples a mini-batch of training nodes \mathcal{B}_v^m . Then, with the cooperation of the server, a subset of L -hop neighbors for the training nodes in \mathcal{B}_v^m is sampled and represented as

²Note that the operation offloading in Swift-FedGNN only supports element-wise (e.g., mean, sum, max) operations, e.g., GCN Kipf & Welling (2017) and SGCN Wu et al. (2019). To support non-element-wise operation, e.g., GAT Veličković et al. (2017), each remote client can transfer the raw graph features or activations to the server for aggregation, instead of performing the locally partial aggregation first with Eq. (5)

324 $\tilde{\mathcal{S}} = \{\tilde{\mathcal{S}}^{(l)}\}_{l=0}^{L-1}$. The nodes $v \in \mathcal{B}_v^m$ are on client m , while for $v \in \tilde{\mathcal{S}}^{(l)}$ with $l < L$, the nodes may
 325 be on clients other than m , denoting the client storing v as $m(v)$. The set $\tilde{\mathcal{M}}(v)$ represents remote
 326 clients with respect to $m(v)$, i.e., $\tilde{\mathcal{M}}(v) \subseteq \mathcal{M} \setminus \{m(v)\}$, where the sampled cross-client neighbors
 327 of the training node v are located. Each remote client $\tilde{m}(v) \in \tilde{\mathcal{M}}(v)$ may contain multiple sampled
 328 neighbors of the training node v , and the numbers of the sampled neighbors can vary across clients.
 329

330 Computing the l -th layer embedding of node v consists of four steps. Steps 1 to 3 below are used to
 331 aggregate the neighbor information of node v , and Step 4 is used to update the node v 's embedding at
 332 l -th GNN layer.

333 **Step 1)** Each remote client $\tilde{m}(v)$ aggregates its sampled neighbors of node v in *parallel*, using
 334

$$335 \tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l), \tilde{m}(v)} = \text{AGG} \left(\{ \tilde{\mathbf{h}}_u^{(l-1), \tilde{m}(v)} \mid u \in \tilde{\mathcal{N}}^{\tilde{m}(v)}(v) \} \right). \quad (5)$$

337 We send only the aggregated results from each remote client $\tilde{m}(v)$ to the server, which can help
 338 preserve data privacy and reduce communication overhead.

339 **Step 2)** Upon receiving the aggregated neighbor information from all the remote clients $\tilde{m}(v) \in$
 340 $\tilde{\mathcal{M}}(v)$, the server aggregates this information from different remote clients before sending it to client
 341 $m(v)$ as follows:
 342

$$343 \tilde{\mathbf{r}}_{\mathcal{N}(v)}^{(l)} = \text{AGG} \left(\{ \tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l), \tilde{m}(v)} \mid \tilde{m}(v) \in \tilde{\mathcal{M}}(v) \} \right). \quad (6)$$

344 This approach not only helps maintain data privacy but also reduces communication costs by mini-
 345 mizing the amount of data transmitted between clients and the server.
 346

347 **Step 3)** Neighbor information of node v for both the sampled local neighbors and the sampled
 348 cross-client neighbors is aggregated as follows:

$$349 \tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l)} = \text{AGG} \left(\{ \tilde{\mathbf{h}}_u^{(l-1), m(v)} \mid u \in \tilde{\mathcal{N}}^{m(v)}(v) \} \cup \{ \tilde{\mathbf{r}}_{\mathcal{N}(v)}^{(l)} \} \right). \quad (7)$$

351 The cross-client neighbor information used here helps mitigate the information loss and reduce the
 352 performance degradation caused by connected nodes being distributed across different clients.
 353

354 **Step 4)** The embedding of node v in the l -th GNN layer is updated using the aggregated neighbor
 355 information and the embedding of node v from the $(l-1)$ -th layer as:

$$356 \tilde{\mathbf{h}}_v^{(l), m(v)} = \sigma \left(\mathbf{W}_t^{(l), m(v)} \cdot \text{COMBINE}(\tilde{\mathbf{h}}_v^{(l-1), m(v)}, \tilde{\mathbf{h}}_{\mathcal{N}(v)}^{(l)}) \right). \quad (8)$$

358 Using the embeddings of the training nodes in the mini-batch and the model parameters, the local
 359 stochastic gradients $\nabla \tilde{F}^m(\boldsymbol{\theta}_t^m)$ are computed and then used in the update of the global model
 360 parameters shown as $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m)$, where α is the learning rate.
 361

362 5 THEORETICAL PERFORMANCE ANALYSIS

363 In this section, we establish the theoretical convergence guarantees for Swift-FedGNN using Graph
 364 Convolutional Network (GCN) (Kipf & Welling, 2017) as the GNN architecture to solve Problem
 365 (1). The analysis of GNN convergence is significantly more challenging compared to the existing
 366 literature on deep neural networks (DNNs). The key difficulties stem from the fact that, unlike in
 367 DNNs, the stochastic gradients in GNNs are inherently biased. This bias is primarily caused by the
 368 presence of cross-client neighbors and the neighbor sampling process. The errors from missing or
 369 unsampled neighbors propagate across layers, gradually getting amplified from the input layer to the
 370 output layer, complicating the overall convergence behavior.
 371

372 For a graph \mathcal{G} , the structure can be represented by its adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$, where $\mathbf{A}_{vu} = 1$
 373 if $(v, u) \in \mathcal{E}$, otherwise $\mathbf{A}_{vu} = 0$. The propagation matrix can be computed as $\mathbf{P} = \mathbf{D}^{-1/2} \hat{\mathbf{A}} \mathbf{D}^{-1/2}$,
 374 where $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$, and $\mathbf{D} \in \mathbb{R}^{N \times N}$ corresponds to the degree matrix and $\mathbf{D}_{vv} = \sum_u \hat{\mathbf{A}}_{vu}$.
 375

376 For subgraph \mathcal{G}^m located on client m , the adjacency matrix \mathbf{A}^m can be denoted as $\mathbf{A}^m = \mathbf{A}_{local}^m +$
 377 \mathbf{A}_{remote}^m , where \mathbf{A}_{local}^m corresponds to the nodes located on client m , and \mathbf{A}_{remote}^m corresponds to
 their cross-client neighbors located on the remote clients other than m . Then, the propagation matrix

can be calculated as $\mathbf{P}^m = \mathbf{D}_m^{-1/2} (\mathbf{A}^m + \mathbf{I}^m) \mathbf{D}_m^{-1/2}$, and can be represented as $\mathbf{P}^m = \mathbf{P}_{local}^m + \mathbf{P}_{remote}^m$, where $\mathbf{P}_{local}^m = \mathbf{D}_m^{-1/2} (\mathbf{A}_{local}^m + \mathbf{I}^m) \mathbf{D}_m^{-1/2}$ and $\mathbf{P}_{remote}^m = \mathbf{D}_m^{-1/2} (\mathbf{A}_{remote}^m) \mathbf{D}_m^{-1/2}$.

Given GCN as the GNN architecture, for client m training using only the local graph data, Eq. (3) and (4) are equivalent to $\widetilde{\mathbf{H}}_t^{(l),m} = \sigma \left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \mathbf{W}_t^{(l),m} \right)$. For client m training based on both the local graph data and the cross-client neighbors, Eq. (5)–(8) are equivalent to $\widetilde{\mathbf{H}}_t^{(l),m} = \sigma \left(\left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} + \widetilde{\mathbf{P}}_{remote}^{(l),m} \widetilde{\mathbf{H}}_{remote}^{(l-1),m} \right) \mathbf{W}_t^{(l),m} \right)$.

Before proceeding with the convergence analysis, we make the following standard assumptions.

Assumption 5.1. The loss function $\ell^m(\cdot, \cdot)$ is C_l -Lipschitz continuous and L_l -smooth with respect to the node embedding $\mathbf{h}^{(L)}$, i.e., $\|\ell^m(\mathbf{h}_1^{(L)}, y) - \ell^m(\mathbf{h}_2^{(L)}, y)\|_2 \leq C_l \|\mathbf{h}_1^{(L)} - \mathbf{h}_2^{(L)}\|_2$ and $\|\nabla \ell^m(\mathbf{h}_1^{(L)}, y) - \nabla \ell^m(\mathbf{h}_2^{(L)}, y)\|_2 \leq L_l \|\mathbf{h}_1^{(L)} - \mathbf{h}_2^{(L)}\|_2$.

Assumption 5.2. The activation function $\sigma(\cdot)$ is C_σ -Lipschitz continuous and L_σ -smooth, i.e., $\|\sigma(\mathbf{z}_1^{(l)}) - \sigma(\mathbf{z}_2^{(l)})\|_2 \leq C_\sigma \|\mathbf{z}_1^{(l)} - \mathbf{z}_2^{(l)}\|_2$ and $\|\nabla \sigma(\mathbf{z}_1^{(l)}) - \nabla \sigma(\mathbf{z}_2^{(l)})\|_2 \leq L_\sigma \|\mathbf{z}_1^{(l)} - \mathbf{z}_2^{(l)}\|_2$.

Assumption 5.3. For any $l \in [L]$, the norm of weight matrices, the propagation matrix, and the node feature matrix are bounded by B_W , B_P and B_X , respectively, i.e., $\|\mathbf{W}^{(l)}\|_F \leq B_W$, $\|\mathbf{P}\|_F \leq B_P$, and $\|\mathbf{X}\|_F \leq B_X$. Note that this assumption is commonly used in the analysis of GNNs, e.g., (Chen et al., 2018; Liao et al., 2020; Garg et al., 2020; Cong et al., 2021; Wan et al., 2022)

Different from DNNs with unbiased stochastic gradients, the stochastic gradients in sampling-based GNNs are *biased* due to neighbor sampling of the training nodes. This is one of the **key challenges** in the convergence performance analysis of Swift-FedGNN. Some existing works used strong assumptions to deal with these biased stochastic gradients in their analysis, e.g., the authors in (Chen et al., 2018) adopted the unbiased stochastic gradient assumption, and the authors in (Chen & Luss, 2018) used the consistent stochastic gradient assumption. However, these assumptions may not hold in reality. In this paper, without using the aforementioned strong assumptions, we are able to bound the errors between the stochastic gradients and the full gradients in the following lemma.

Lemma 5.4. *Under Assumptions 5.1–5.3, the errors between the stochastic gradients and the full gradients are bounded as follows:*

$$\left\| \nabla F_{local}^m(\boldsymbol{\theta}^m) - \nabla \widetilde{F}_{local}^m(\boldsymbol{\theta}^m) \right\|_F \leq LB_{\Delta G}^l, \quad \left\| \nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla \widetilde{F}_{full}^m(\boldsymbol{\theta}^m) \right\|_F \leq LB_{\Delta G}^f,$$

where $\nabla F_{local}^m(\boldsymbol{\theta}^m)$ and $\nabla \widetilde{F}_{local}^m(\boldsymbol{\theta}^m)$ correspond to the full and stochastic gradients computed with only the local graph data, respectively. $\nabla F_{full}^m(\boldsymbol{\theta}^m)$ and $\nabla \widetilde{F}_{full}^m(\boldsymbol{\theta}^m)$ represent the full and stochastic gradients computed with both the local graph data and the cross-client neighbors of the training nodes, respectively. $B_{\Delta G}^l$ and $B_{\Delta G}^f$ are defined in Eq. (12) and (13) in Appendix D.

Furthermore, the dependencies of the nodes located on different clients can lead to additional errors in the gradient computations when client m is updated only with its local graph data, since the cross-client neighbors are missed. This becomes another **key challenge** in the analysis of the convergence of Swift-FedGNN. We prove that such an error is upper-bounded as shown in the following lemma.

Lemma 5.5. *Under Assumptions 5.1–5.3, the error between the full gradient computed with both the local graph data and the cross-client neighbors of the training nodes (denote as $\nabla F_{full}^m(\boldsymbol{\theta}^m)$) and the full gradient computed with only the local graph data (denote as $\nabla F_{local}^m(\boldsymbol{\theta}^m)$) is upper-bounded as follows:*

$$\left\| \nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla F_{local}^m(\boldsymbol{\theta}^m) \right\|_F \leq LB_{\Delta G}^r,$$

where $B_{\Delta G}^r$ is defined in Eq. (14) in Appendix D.

We note that all the errors mentioned in Lemmas 5.4 and 5.5 are correlated with the structure of GNNs, specifically showing a positive correlation with the number of layers in the networks. This finding is unique to GNNs, where each layer involves both neighbor aggregation and non-linear transformation. As these two operations are interleaved across multiple layers, they create a structural entanglement that complicates the analysis.

Using Lemmas 5.4 and 5.5, we state the main convergence result of Swift-FedGNN solving an L -layer GNN in the following theorem:

Theorem 5.6. Under Assumptions 5.1–5.3, choose step-size $\alpha = \min\left\{\frac{\sqrt{M}}{\sqrt{T}}, \frac{1}{L_F}\right\}$, where L_F is the smoothness constant in Lemma D.2. The output of Swift-FedGNN solving an L -layer GNN satisfies:

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla \mathcal{L}(\theta_t)\|^2 \leq \frac{2}{\sqrt{MT}} (\mathcal{L}(\theta_0) - \mathcal{L}(\theta^*)) + L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2 + \frac{K}{IM} L^2 \left((B_{\Delta G}^f)^2 - (B_{\Delta G}^l + B_{\Delta G}^r)^2 \right).$$

The detailed proof of Theorem 5.6 can be found in Appendix D. We can see from Theorem 5.6 that the convergence rate of Swift-FedGNN is $\mathcal{O}(T^{-1/2})$ to a neighborhood of the exact solution, which matches the SOTA convergence rate of sampling-based GNN algorithms, e.g., (Chen et al., 2018; Cong et al., 2021; Ramezani et al., 2022; Du & Wu, 2022), even though Swift-FedGNN operates in the far more challenging federated setting.

Three important remarks on Theorem 5.6 are in order: (1) When choosing $I = 1$ and $K = M$, Swift-FedGNN performs fully cross-client training, ensuring no information loss in the graph data. In this scenario, Swift-FedGNN experiences minimal residual error. Such error is caused by sampling and is inevitable. However, Swift-FedGNN suffers from maximum sampling and communication overhead; (2) When choosing $K = 0$, Swift-FedGNN conducts fully local training, resulting in the information loss of all the cross-client neighbors. Consequently, Swift-FedGNN encounters maximum residual error. Nonetheless, the sampling and communication overhead is minimized; and (3) It can be shown that the last term of the convergence rate bound in Theorem 5.6 is negative. Hence, increasing I or decreasing K would increase the residual error due to more information loss of the cross-client neighbors. However, this would reduce the sampling and communication overhead. Thus, there is a trade-off between the information loss and the sampling and communication overhead.

6 NUMERICAL RESULTS

In this section, we conduct experiments to evaluate the performance of Swift-FedGNN. Table 1: Benchmark datasets and key parameters.

DATASET	# OF NODES	# OF EDGES
OGBN-PRODUCTS	2.4 M	61.8 M
REDDIT	0.2 M	114.6 M

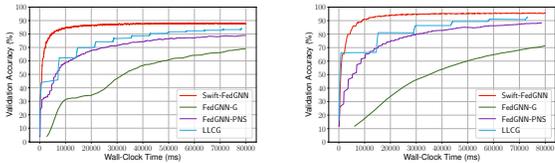
1) Experiment Settings: We train a representative GNN model, GraphSAGE (Hamilton et al., 2017), in the FL settings on two real-world node

classification datasets: 1) ogbn-products (Hu et al., 2020), which is an Amazon product co-purchasing graph derived from (Leskovec et al., 2007); and 2) Reddit (Hamilton et al., 2017), which consists of online forum posts within a month, where posts commented on by the same user are connected by an edge. Table 1 summarizes the key statistics of the datasets. Note that Ogbn-products dataset is the largest dataset one can find in the federated GNN literature, while the Reddit dataset is known for its density. These datasets were chosen for their distinct and representative characteristics, ensuring a thorough evaluation that addresses diverse scenarios in federated GNN training. In our FL simulations, we use 20 clients for the experiments with ogbn-products dataset and 10 clients for the experiments with Reddit dataset. Both graphs are partitioned with METIS partitioning (Karypis & Kumar, 1998). Due to space limitations, additional experimental details and results are provided in Appendix C.

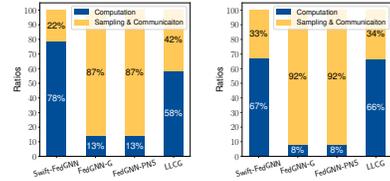
2) Baselines: Since the goal of Swift-FedGNN is to reduce the sampling and communication time, we compare Swift-FedGNN with the algorithms most closely related to Swift-FedGNN, which mitigates the information loss of cross-client neighbors through periodical (sampling-based) full-neighbor training: **1) LCG** (Ramezani et al., 2022): A distributed GNN training framework that performs local training on each client independently, with periodic full-neighbor training conducted on a central server; **2) FedGNN-PNS** (Du & Wu, 2022): A federated GNN training framework where each client periodically samples cross-client neighbors with an increasing sampling frequency. In the remaining iterations, clients reuse the most recently sampled cross-client neighbors; and **3) FedGNN-G:** Naive federated GNN training where cross-client training is performed on each client in every iteration.

3) Convergence Performance Comparisons: In Figure 4, we can see that for both the ogbn-products dataset and the Reddit dataset, Swift-FedGNN demonstrates the fastest convergence speed compared to the baseline algorithms, which verifies the effectiveness of Swift-FedGNN. In addition, the validation accuracy of Swift-FedGNN is comparable to that of FedGNN-G, which trains a GNN model on the dataset without any information loss. Specifically, when Swift-FedGNN converges, the validation accuracy is 87.73% on the ogbn-products dataset and 95.60% on the Reddit dataset. When FedGNN-G converges, the validation accuracy is 87.93% on ogbn-products dataset and 96.03% on Reddit dataset. Although LCG performs periodic cross-client training on the server, it requires

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539



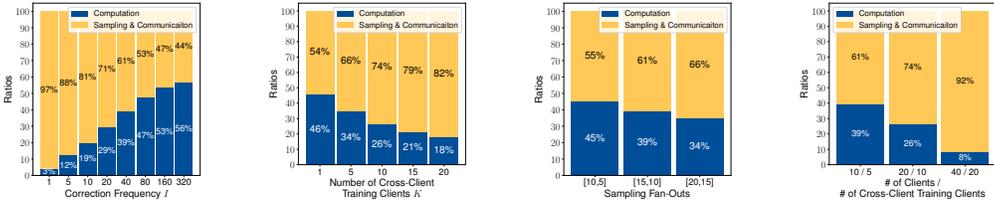
(a) ogbn-products (b) Reddit



(a) ogbn-products (b) Reddit

Figure 4: Convergence performance in terms of validation accuracy of different algorithms.

Figure 5: Comp.-(sampling & comm.) ratio of different algorithms.



(a) correction frequencies (I) (b) # of cross-client training clients (K) (c) # of sampled neighbors (d) # of clients (50% for cross-client training)

Figure 6: Comp.-(sampling & comm.) ratio of Swift-FedGNN on ogbn-products dataset.

training over the full set of neighbors of the training nodes, leading to significant sampling and communication overhead. For instance, when training the ogbn-products dataset, LLCC takes over 5000 ms to perform cross-client training on the server, whereas Swift-FedGNN completes cross-client training within 200 ms due to neighbor sampling. FedGNN-PNS employs a dynamic cross-client sampling interval throughout training, gradually reducing the interval as training progresses. Consequently, FedGNN-PNS incurs extensive sampling and communication overhead during the later stages of training, slowing down the convergence process.

4) Communication and Sample Costs Analysis: Figure 5 illustrates the comparison between the ratios of the computation time and the sampling and communication time for Swift-FedGNN and the baseline algorithms. It can be seen that Swift-FedGNN significantly reduces the computation-(sampling & communication) ratio on the ogbn-products dataset. On the Reddit dataset, Swift-FedGNN also significantly reduces this ratio compared to FedGNN-PNS and FedGNN-G. While Swift-FedGNN achieves a comparable ratio to LLCC, it converges much faster and achieves higher validation accuracy than LLCC.

5) Hyperparameter Sensitivity Analysis: We explore the impact of the important hyperparameters in Swift-FedGNN. Figure 6a shows that when the correction frequency I increases, the computation-(sampling & communication) ratio increases. Figure 6b and 6c indicate that as the number of cross-client training clients K , and the number of sampled neighbors increase, the computation-(sampling & communication) ratio decreases. Figure 6d evaluates Swift-FedGNN with different numbers of clients. In this experiment, 50% of clients periodically conduct cross-client training on both local and cross-client neighbors. We can see that as the number of clients increases, the computation-(sampling & communication) ratio decreases. These findings align with our expectations since sampling and communication overhead is significantly greater than computation overhead in GNN training.

7 CONCLUSION

In this paper, we proposed the Swift-FedGNN algorithm, which is a mini-batch-based and sampling-based federated GNN framework, for efficient federated GNN training. Swift-FedGNN reduces the cross-client neighbor sampling and communication overhead by *periodically* sampling a set of clients to conduct the local GNN training on local graph data and cross-client neighbors, which is time-consuming. The rest clients in these periodical iterations and all the clients in the remaining iterations perform efficient parallel local GNN training using only local graph data. We theoretically proved that the convergence rate of Swift-FedGNN is $\mathcal{O}(T^{-1/2})$, matching the SOTA rate of sampling-based GNN methods, even in the more challenging federated settings. We conducted extensive numerical experiments on real-world graph datasets and verified the effectiveness of Swift-FedGNN.

REFERENCES

- 540
541
542 Davide Bacciu, Federico Errica, and Alessio Micheli. Contextual graph markov model: A deep and
543 generative approach to graph processing. In *International conference on machine learning*, pp.
544 294–303. PMLR, 2018.
- 545 Zhenkun Cai, Xiao Yan, Yidi Wu, Kaihao Ma, James Cheng, and Fan Yu. DGCL: An Efficient
546 Communication Library for Distributed GNN Training. In *Proc. of EuroSys*, pp. 130–144, 2021.
- 547 Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with
548 variance reduction. In *International Conference on Machine Learning*, pp. 942–950. PMLR, 2018.
- 549 Jie Chen and Ronny Luss. Stochastic gradient descent with biased but consistent gradient estimators.
550 *arXiv preprint arXiv:1807.11880*, 2018.
- 551
552 Weilin Cong, Morteza Ramezani, and Mehrdad Mahdavi. On the importance of sampling in training
553 gcns: Tighter analysis and variance reduction. *arXiv preprint arXiv:2103.02696*, 2021.
- 554
555 Zhiyong Cui, Kristian Henrickson, Ruimin Ke, and Yinhai Wang. Traffic graph convolutional
556 recurrent neural network: A deep learning framework for network-scale traffic learning and
557 forecasting. *IEEE Transactions on Intelligent Transportation Systems*, 21(11):4883–4894, 2019.
- 558 Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Learning dynamic context graphs for predicting
559 social events. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge
560 Discovery & Data Mining*, pp. 1007–1016, 2019.
- 561 Kien Do, Truyen Tran, and Svetha Venkatesh. Graph transformation policy network for chemical
562 reaction prediction. In *Proceedings of the 25th ACM SIGKDD international conference on
563 knowledge discovery & data mining*, pp. 750–760, 2019.
- 564
565 Bingqian Du and Chuan Wu. Federated graph learning with periodic neighbour sampling. In *2022
566 IEEE/ACM 30th International Symposium on Quality of Service (IWQoS)*, pp. 1–10. IEEE, 2022.
- 567
568 Bingqian Du, Jun Liu, Ziyue Luo, Chuan Wu, Qiankun Zhang, and Hai Jin. Expediting Distributed
569 GNN Training with Feature-only Partition and Optimized Communication Planning. In *Proc. of
570 IEEE INFOCOM*, 2024.
- 571 Matthias Fey and Jan Eric Lenssen. Fast Graph Representation Learning with PyTorch Geometric. In
572 *Proc. of ICLR*, 2019.
- 573
574 Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. Protein interface prediction using graph
575 convolutional networks. *Advances in neural information processing systems*, 30, 2017.
- 576 Swapnil Gandhi and Anand Padmanabha Iyer. P3: Distributed Deep Graph Learning at Scale. In
577 *Proc. of OSDI*, pp. 551–568, 2021.
- 578
579 Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits of
580 graph neural networks. In *International Conference on Machine Learning*, pp. 3419–3430. PMLR,
581 2020.
- 582 Mahsa Ghorbani, Anees Kazi, Mahdieh Soleymani Baghshah, Hamid R Rabiee, and Nassir Navab.
583 Ra-gcn: Graph convolutional network for disease prediction problems with imbalanced data.
584 *Medical image analysis*, 75:102272, 2022.
- 585
586 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
587 *Advances in neural information processing systems*, 30, 2017.
- 588
589 Chaoyang He, Keshav Balasubramanian, Emir Ceyani, Carl Yang, Han Xie, Lichao Sun, Lifang
590 He, Liangwei Yang, Philip S Yu, Yu Rong, et al. Fedgraphnn: A federated learning system and
591 benchmark for graph neural networks. *arXiv preprint arXiv:2104.07145*, 2021a.
- 592
593 Chaoyang He, Emir Ceyani, Keshav Balasubramanian, Murali Annavaram, and Salman Avestimehr.
Spreadgnn: Serverless multi-task federated learning for graph neural networks. *arXiv preprint
arXiv:2106.02743*, 2021b.

- 594 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,
595 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in*
596 *neural information processing systems*, 33:22118–22133, 2020.
- 597
- 598 Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and
599 Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In
600 *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.
- 601 George Karypis and Vipin Kumar. A Fast and High Quality Multilevel Scheme for Partitioning
602 Irregular Graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.
- 603
- 604 Anees Kazi, Soroush Farghadani, Iman Aganj, and Nassir Navab. Ia-gcn: Interpretable attention
605 based graph convolutional network for disease prediction. In *International Workshop on Machine*
606 *Learning in Medical Imaging*, pp. 382–392. Springer, 2023.
- 607 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
608 In *International Conference on Learning Representations*, 2017.
- 609
- 610 Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in
611 temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international conference*
612 *on knowledge discovery & data mining*, pp. 1269–1278, 2019.
- 613
- 614 Jure Leskovec, Lada A Adamic, and Bernardo A Huberman. The Dynamics of Viral Marketing.
615 *ACM Transactions on the Web (TWEB)*, 1(1):5–es, 2007.
- 616 Jia Li, Zhichao Han, Hong Cheng, Jiao Su, Pengyun Wang, Jianfeng Zhang, and Lujia Pan. Predicting
617 path failure in time-evolving graphs. In *Proceedings of the 25th ACM SIGKDD international*
618 *conference on knowledge discovery & data mining*, pp. 1279–1289, 2019.
- 619
- 620 Sihui Li and Rui Zhang. A novel interactive deep cascade spectral graph convolutional network with
621 multi-relational graphs for disease prediction. *Neural Networks*, pp. 106285, 2024.
- 622 Renjie Liao, Raquel Urtasun, and Richard Zemel. A pac-bayesian approach to generalization bounds
623 for graph neural networks. *arXiv preprint arXiv:2012.07690*, 2020.
- 624
- 625 Tianfeng Liu, Yangrui Chen, Dan Li, Chuan Wu, Yibo Zhu, Jun He, Yanghua Peng, Hongzheng
626 Chen, Hongzhi Chen, and Chuanxiong Guo. BGL: GPU-Efficient GNN Training by Optimizing
627 Graph Data I/O and Preprocessing. In *Proc. of NSDI*, pp. 103–118, 2023.
- 628 Ziyue Luo, Yixin Bao, and Chuan Wu. Optimizing Task Placement and Online Scheduling for
629 Distributed GNN Training Acceleration. In *Proc. of IEEE INFOCOM*, 2022.
- 630
- 631 Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas.
632 Communication-efficient learning of deep networks from decentralized data. In *Artificial intelli-*
633 *gence and statistics*, pp. 1273–1282. PMLR, 2017.
- 634 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
635 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An Imperative Style,
636 High-Performance Deep Learning Library. In *Proc. of NeurIPS*, 2019.
- 637
- 638 Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. Deepinf: Modeling
639 influence locality in large social networks. In *Proceedings of the 24th ACM SIGKDD International*
640 *Conference on Knowledge Discovery and Data Mining (KDD’18)*, 2018.
- 641 Morteza Ramezani, Weilin Cong, Mehrdad Mahdavi, Mahmut Kandemir, and Anand Sivasubrama-
642 niam. Learn locally, correct globally: A distributed algorithm for training graph neural networks.
643 In *International Conference on Learning Representations*, 2022.
- 644
- 645 John Thorpe, Yifan Qiao, Jonathan Eyolfson, Shen Teng, Guanzhou Hu, Zhihao Jia, Jinliang Wei,
646 Keval Vora, Ravi Netravali, Miryung Kim, et al. Dorylus: Affordable, Scalable, and Accurate
647 GNN Training with Distributed CPU Servers and Serverless Threads. In *Proc. of USENIX OSDI*,
2021.

- 648 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
649 Bengio. Graph Attention Networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 650
- 651 Cheng Wan, Youjie Li, Cameron R. Wolfe, Anastasios Kyrillidis, Nam Sung Kim, and Yingyan Lin.
652 PipeGCN: Efficient full-graph training of graph convolutional networks with pipelined feature
653 communication. In *International Conference on Learning Representations*, 2022.
- 654 Binghui Wang, Ang Li, Meng Pang, Hai Li, and Yiran Chen. Graphfl: A federated learning framework
655 for semi-supervised node classification on graphs. In *2022 IEEE International Conference on Data
656 Mining (ICDM)*, pp. 498–507. IEEE, 2022a.
- 657 Hao Wang, Tong Xu, Qi Liu, Defu Lian, Enhong Chen, Dongfang Du, Han Wu, and Wen Su.
658 Mcne: An end-to-end framework for learning multiple conditional network representations of
659 social network. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge
660 discovery & data mining*, pp. 1064–1072, 2019a.
- 661 Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and
662 Zhongyuan Wang. Knowledge-aware graph neural networks with label smoothness regularization
663 for recommender systems. In *Proceedings of the 25th ACM SIGKDD international conference on
664 knowledge discovery & data mining*, pp. 968–977, 2019b.
- 665
- 666 Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia.
667 Microsoft Academic Graph: When Experts Are Not Enough. *Quantitative Science Studies*, 1(1):
668 396–413, 2020.
- 669
- 670 Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang,
671 Chao Ma, et al. Deep Graph Library: Towards Efficient and Scalable Deep Learning on Graphs. In
672 *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019c.
- 673 Xiang Wang, Xiangnan He, Yixin Cao, Meng Liu, and Tat-Seng Chua. Kgat: Knowledge graph
674 attention network for recommendation. In *Proceedings of the 25th ACM SIGKDD international
675 conference on knowledge discovery & data mining*, pp. 950–958, 2019d.
- 676 Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive
677 learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3):279–287,
678 2022b.
- 679
- 680 Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. Fedgcn: Federated graph
681 neural network for privacy-preserving recommendation. *arXiv preprint arXiv:2102.04925*, 2021.
- 682 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-
683 plifying graph convolutional networks. In *International conference on machine learning*, pp.
684 6861–6871. PMLR, 2019.
- 685
- 686 Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation
687 in non-IID federated learning. In *International Conference on Learning Representations*, 2021.
- 688 Yuhang Yao, Weizhao Jin, Srivatsan Ravi, and Carlee Joe-Wong. Fedgcn: Convergence-
689 communication tradeoffs in federated training of graph convolutional networks. *Advances in
690 Neural Information Processing Systems*, 36, 2023a.
- 691 Yuhang Yao, Mohammad Mahdi Kamani, Zhongwei Cheng, Lin Chen, Carlee Joe-Wong, and
692 Tianqiang Liu. Fedrule: Federated rule recommendation system with graph neural networks. In
693 *Proceedings of the 8th ACM/IEEE Conference on Internet of Things Design and Implementation*,
694 pp. 197–208, 2023b.
- 695
- 696 Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec.
697 Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the
698 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 974–983,
699 2018.
- 700 Binhang Yuan, Yongjun He, Jared Davis, Tianyi Zhang, Tri Dao, Beidi Chen, Percy S Liang,
701 Christopher Re, and Ce Zhang. Decentralized training of foundation models in heterogeneous
environments. volume 35, pp. 25464–25477, 2022.

702 Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graph-
703 SAINT: Graph Sampling Based Inductive Learning Method. In *Proc. of ICLR*, 2020.
704

705 Ke Zhang, Carl Yang, Xiaoxiao Li, Lichao Sun, and Siu Ming Yiu. Subgraph federated learning with
706 missing neighbor generation. *Advances in Neural Information Processing Systems*, 34:6671–6682,
707 2021.

708 Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in neural
709 information processing systems*, 31, 2018.
710

711 Muhan Zhang, Zhicheng Cui, Marion Neumann, and Yixin Chen. An end-to-end deep learning
712 architecture for graph classification. In *Proceedings of the AAAI conference on artificial intelligence*,
713 volume 32, 2018.

714 Zhe Zhang, Ziyue Luo, and Chuan Wu. Two-Level Graph Caching for Expediting Distributed GNN
715 Training. In *Proc. of INFOCOM*, pp. 1–10. IEEE, 2023.
716

717 Kun Zhao, Wencong Xiao, Baole Ai, Wenting Shen, Xiaolin Zhang, Yong Li, and Wei Lin. AliGraph:
718 An Industrial Graph Neural Network Platform. In *Proc. of SOSIP Workshop on AI Systems*, 2019.

719 Da Zheng, Chao Ma, Minjie Wang, Jinjing Zhou, Qidong Su, Xiang Song, Quan Gan, Zheng Zhang,
720 and George Karypis. DistDGL: Distributed Graph Neural Network Training for Billion-Scale
721 Graphs. In *IEEE/ACM Workshop on Irregular Applications: Architectures and Algorithms*, 2020.
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

A LIST OF NOTATIONS

$\mathcal{G}(\mathcal{V}, \mathcal{E})$	Graph
\mathcal{V}	Set of nodes
\mathcal{E}	Set of edges
$N = \mathcal{V} $	Number of nodes
\mathcal{M}	Set of clients
$M = \mathcal{M} $	Number of clients
$\mathcal{G}^m(\mathcal{V}^m, \mathcal{E}^m)$	Subgraph at client m
\mathcal{V}^m	Set of nodes at client m
\mathcal{E}^m	Set of edges at client m
$\mathbf{x}_v^m \in \mathbb{R}^d$	Feature vector of node v at client m
y_v^m	Label of node v at client m
ℓ^m	Loss function (<i>e.g.</i> , cross-entropy loss) at client m
\mathcal{V}_B^m	Mini-batch of training nodes
$\boldsymbol{\theta} = \{\mathbf{W}^{(l)}\}_{l=1}^L$	Set of trainable model parameters
$m(v)$	Local client of node v
$\bar{m}(v)$	Remote client of node v
$\bar{\mathcal{M}}(v)$	Set of the remote clients that host the neighbors of node v
$\mathcal{N}^{m(v)}(v)$	Set of the neighbors of node v located on local client $m(v)$
$\mathcal{N}^{\bar{m}(v)}(v)$	Set of the neighbors of node v located on remote client $\bar{m}(v)$
$\mathbf{h}_v^{(l), m(v)}$	Embedding of node v located on client $m(v)$
$\mathbf{h}_{\mathcal{N}(v)}^{(l)}$	Aggregated embedding from node v 's neighbors
$\mathbf{W}^{(l)}$	Weight matrix at l -th layer
$\sigma(\cdot)$	Activation function (<i>e.g.</i> , ReLU)
AGG(\cdot)	Aggregation function (<i>e.g.</i> , mean)
COMBINE(\cdot)	Combination function (<i>e.g.</i> , concatenation)
\mathcal{B}_v^m	Mini-batch of training nodes at client m
$\tilde{\mathcal{S}} = \{\tilde{\mathcal{S}}^{(l)}\}_{l=0}^{L-1}$	Subset of L -hop neighbors for the training nodes in \mathcal{B}_v^m
$\tilde{\mathcal{N}}^m(v)$	Set of the sampled neighbors located on client m for node v
\mathcal{K}	Set of sampled clients for cross-client training
$K = \mathcal{K} $	Number of sampled clients for cross-client training
$\tilde{\mathcal{M}}(v)$	Set of remote clients that host the sampled cross-client neighbors of the training node v
$\tilde{m}(v)$	Remote client with respect to $m(v)$
$\nabla F^m(\boldsymbol{\theta}_t^m)$	Stochastic gradient
α	Learning rate
$\mathbf{A} \in \mathbb{R}^{N \times N}$	Adjacency matrix of graph \mathcal{G}
\mathbf{P}	Propagation matrix

810	D	Degree matrix
811	A^m	Adjacency matrix of subgraph \mathcal{G}^m
812	A_{local}^m	Adjacency matrix corresponds to the nodes located on client m
813	A_{remote}^m	Adjacency matrix corresponds to the cross-client neighbors located on the remote clients other than m
814	D^m	Degree matrix of client m
815	P^m	Propagation matrix of client m
816	P_{local}^m	Propagation matrix corresponds to the nodes located on client m
817	P_{remote}^m	Propagation matrix corresponds to the cross-client neighbors located on the remote clients other than m

826 B SINGLE-MACHINE GRAPH NEURAL NETWORKS TRAINING

827 We consider a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, where \mathcal{V} is a set of nodes with $N = |\mathcal{V}|$ and \mathcal{E} is a set of edges. Each node $v \in \mathcal{V}$ is associated with a feature vector $x_v \in \mathbb{R}^d$, where d is the dimension of the feature vector. Each node $v \in \mathcal{V}_{train}$ has a corresponding label y_v , where $\mathcal{V}_{train} \subseteq \mathcal{V}$.

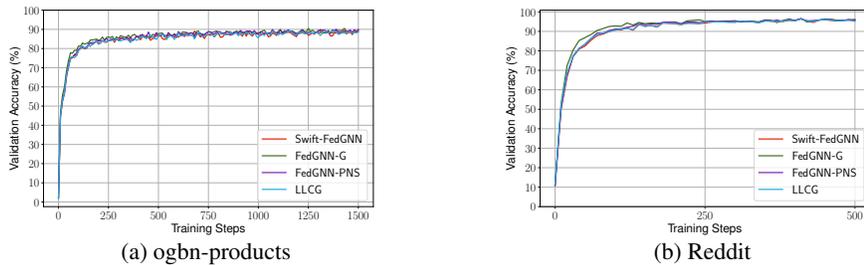
831 GNNs aim to generate representations (embeddings) for each node in the graph by combining information from its neighboring nodes. Consider a GNN that consists of L layers. The embedding of node v at l -th layer, which is represented by $h_v^{(l)}$, can be obtained through neighbor aggregation and node update, which are formulated as follows:

$$832 \quad h_{\mathcal{N}(v)}^{(l)} = \text{AGG} \left(\{h_u^{(l-1)} \mid u \in \mathcal{N}(v)\} \right), \quad h_v^{(l)} = \sigma \left(\mathbf{W}^{(l)} \cdot \text{COMBINE}(h_v^{(l-1)}, h_{\mathcal{N}(v)}^{(l)}) \right),$$

833 where $h_v^{(0)}$ is initialized as the feature vector x_v , $\mathcal{N}(v)$ denotes the set of neighbors of node v , $h_{\mathcal{N}(v)}^{(l)}$ is the aggregated embedding from node v 's neighbors aggregated neighbor embedding for node v , $\mathbf{W}^{(l)}$ represents the weight matrix at l -th layer, $\sigma(\cdot)$ corresponds to an activation function (e.g., ReLU), $\text{AGG}(\cdot)$ is an aggregation function (e.g., mean), and $\text{COMBINE}(\cdot)$ is a combination function (e.g., concatenation).

845 C ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

847 C.1 ADDITIONAL EXPERIMENTAL RESULTS



857 Figure 7: Convergence performance (validation accuracy versus training steps) of different algorithms.

860 Table 2: Communication overhead per iteration when communication occurs.

	Swift-FedGNN	LLCG	FedGNN-PNS	FedGNN-G
861 OGBN-PRODUCTS	19.5 MB	378.3 MB	78.0 MB	78.0 MB
862 REDDIT	90.4 MB	619.6 MB	180.7 MB	180.7 MB

Figure 7 shows that the numbers of iterations required by all algorithms in comparison are similar, which is due to the fact that they share the same convergence rate result. However, since Swift-FedGNN minimizes the sampling and communication overhead, it achieves the lowest wall-clock time for convergence, making it the most efficient in terms of practical implementation.

Table 2 shows the communication overhead per iteration when cross-client sampling and communication occur for different algorithms. We can see that Swift-FedGNN significantly reduces the communication overhead compared to all baselines across both datasets. Specifically, on the ogbn-products dataset, Swift-FedGNN incurs 19.5 MB of overhead per iteration, which is approximately 20 times less than LLCG and 4 times less than both FedGNN-PNS and FedGNN-G. Similarly, for the Reddit dataset, due to its dense inter-node connections and larger feature size, Swift-FedGNN’s overhead is 90.4 MB, which is still about 7 times less than LLCG and 2 times less than both FedGNN-PNS and FedGNN-G. This highlights the efficiency of Swift-FedGNN in reducing communication costs during cross-client training.

C.2 ADDITIONAL EXPERIMENTAL DETAILS

Implementation and testbed. We implement Swift-FedGNN using Python on DGL 2.0.0 (Wang et al., 2019c) and PyTorch 2.2.1 (Paszke et al., 2019) with 1241 LoC. Our implementation includes a custom GPU-based sampler built on top of DGL’s native sampler, which is designed to sequentially sample local and remote neighbors for each client at every layer. Additionally, we customized the GraphSAGE layer (Hamilton et al., 2017) to facilitate model-parallel training within Swift-FedGNN. In this setup, the server handles the sampling and aggregation of node features and intermediate activations, while the clients are responsible for executing the nonlinear computations associated with the GraphSAGE layer.

We simulate a real-world federated learning scenario using a single machine equipped with NVIDIA Tesla V100 GPUs and 64GB memory. In our setup, both the clients and the server operate on the GPU, and data communication between them is simulated using shared memory. We monitor the data transfer size between the server and clients and set a simulated cross-client network bandwidth at 1Gbps, aligning with real-world measurements reported in (Yuan et al., 2022).

GNN Model. We train a two-layer GraphSAGE model with a hidden dimension of 256. Uniform sampling is employed for neighbor sampling, with fan-outs—*i.e.*, the number of sampled neighbors—set according to the official training script provided by the DGL team. The fan-outs for both the ogbn-products dataset and the Reddit dataset are set to be [15, 10]. The training mini-batch size is set at 256. For optimization, we use the Adam optimizer with a learning rate of 0.001 and a weight decay of 5×10^{-4} .

D PROOF OF THEOREM 5.6

D.1 GRADIENT COMPUTATIONS IN Swift-FedGNN

Recall that Swift-FedGNN uses GCN (Kipf & Welling, 2017) as the architecture of GNN to prove the convergence performance. When client m performs local training that updates the local GNN model using only the local graph data, Each sampling-based GCN layer executes one feature propagation step, defined as:

$$\widetilde{\mathbf{H}}_{local}^{(l),m} = \left[\widetilde{f}^{(l),m} \left(\widetilde{\mathbf{H}}_{local}^{(l-1),m}, \mathbf{W}^{(l),m} \right) \triangleq \sigma \left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \mathbf{W}^{(l),m} \right) \right].$$

Using the chain rule, the stochastic gradient can be computed as $\nabla \widetilde{F}^m(\boldsymbol{\theta}^m) = \{ \widetilde{\mathbf{G}}_{local}^{(l),m} \}_{l=1}^L$, where

$$\begin{aligned} \widetilde{\mathbf{G}}_{local}^{(l),m} &= \left[\nabla_{\mathbf{W}} \widetilde{f}^{(l),m} \left(\widetilde{\mathbf{D}}_{local}^{(l),m}, \widetilde{\mathbf{H}}_{local}^{(l-1),m}, \mathbf{W}^{(l),m} \right) \right. \\ &\quad \left. \triangleq \left[\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \right]^{\top} \widetilde{\mathbf{D}}_{local}^{(l),m} \circ \nabla \sigma \left(\widetilde{\mathbf{Z}}_{local}^{(l),m} \right) \right], \\ \widetilde{\mathbf{D}}_{local}^{(l),m} &= \left[\nabla_{\mathbf{H}} \widetilde{f}^{(l+1),m} \left(\widetilde{\mathbf{D}}_{local}^{(l+1),m}, \widetilde{\mathbf{H}}_{local}^{(l),m}, \mathbf{W}^{(l+1),m} \right) \right] \end{aligned}$$

$$\triangleq \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top,$$

in which $\tilde{\mathbf{Z}}_{local}^{(l),m} = \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} \mathbf{W}^{(l),m}$, $\tilde{\mathbf{D}}_{local}^{(L),m} = \partial \ell^m \left(\tilde{\mathbf{H}}_{local}^{(L),m}, \mathbf{Y}_{local}^m \right) / \partial \tilde{\mathbf{H}}_{local}^{(L),m}$, and \circ represents Hadamard product.

Similarly, when client m conducts cross-client training that updates the local GNN model based on the local graph data and the cross-client neighbors, each sampling-based GNN layer can be defined as:

$$\tilde{\mathbf{H}}_{full}^{(l),m} = \left[\tilde{f}^{(l),m} \left(\tilde{\mathbf{H}}_{full}^{(l-1),m}, \mathbf{W}^{(l),m} \right) \triangleq \sigma \left(\left(\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} + \tilde{\mathbf{P}}_{remote}^{(l),m} \tilde{\mathbf{H}}_{remote}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right) \right].$$

Using the chain rule, the stochastic gradient can be calculated as $\nabla \tilde{F}^m(\boldsymbol{\theta}^m) = \{ \tilde{\mathbf{G}}_{full}^{(l),m} \}_{l=1}^L$, where

$$\begin{aligned} \tilde{\mathbf{G}}_{full}^{(l),m} &= \left[\nabla_W \tilde{f}^{(l),m} \left(\tilde{\mathbf{D}}_{full}^{(l),m}, \tilde{\mathbf{H}}_{full}^{(l-1),m}, \mathbf{W}^{(l),m} \right) \right. \\ &\quad \left. \triangleq \left[\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} + \tilde{\mathbf{P}}_{remote}^{(l),m} \tilde{\mathbf{H}}_{remote}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) \right], \\ \tilde{\mathbf{D}}_{full}^{(l),m} &= \left[\nabla_H \tilde{f}^{(l+1),m} \left(\tilde{\mathbf{D}}_{full}^{(l+1),m}, \tilde{\mathbf{H}}_{full}^{(l),m}, \mathbf{W}^{(l+1),m} \right) \right. \\ &\quad \left. \triangleq \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} + \tilde{\mathbf{P}}_{remote}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right], \end{aligned}$$

in which $\tilde{\mathbf{Z}}_{full}^{(l),m} = \left(\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} + \tilde{\mathbf{P}}_{remote}^{(l),m} \tilde{\mathbf{H}}_{remote}^{(l-1),m} \right) \mathbf{W}^{(l),m}$, and $\tilde{\mathbf{D}}_{full}^{(L),m} = \partial \ell^m \left(\tilde{\mathbf{H}}_{full}^{(L),m}, \mathbf{Y}_{full}^m \right) / \partial \tilde{\mathbf{H}}_{full}^{(L),m}$.

D.2 USEFUL PROPOSITIONS AND LEMMAS

Proposition D.1. *Under Assumption 5.3, the inequalities in Table 3 and Table 4 are hold.*

Table 3: Upper-bound for the norms of the propagation matrix and the node feature matrix.

	PROPAGATION MATRIX	NODE FEATURE MATRIX
FULL GRAPH	$\ \mathbf{P}_{full}\ _F \leq B_P$	$\ \mathbf{X}_{full}\ _F \leq B_X$
LOCAL GRAPH	$\ \mathbf{P}_{local}\ _F \leq B_P^l \leq B_P$	$\ \mathbf{X}_{local}\ _F \leq B_X^l \leq B_X$
CROSS-CLIENT NEIGHBORS	$\ \mathbf{P}_{remote}\ _F \leq B_P^r \leq B_P$	$\ \mathbf{X}_{remote}\ _F \leq B_X^r \leq B_X$

Table 4: Relationships for the norms of the propagation matrix and the node feature matrix before and after sampling.

	PROPAGATION MATRIX	NODE FEATURE MATRIX
FULL GRAPH	$\ \tilde{\mathbf{P}}_{full} - \mathbf{P}_{full}\ _F \leq B_{\Delta P}^f$	$\ \tilde{\mathbf{X}}_{full} - \mathbf{X}_{full}\ _F \leq B_{\Delta X}^f$
LOCAL GRAPH	$\ \tilde{\mathbf{P}}_{local} - \mathbf{P}_{local}\ _F \leq B_{\Delta P}^l$	$\ \tilde{\mathbf{X}}_{local} - \mathbf{X}_{local}\ _F \leq B_{\Delta X}^l$
CROSS-CLIENT NEIGHBORS	$\ \tilde{\mathbf{P}}_{remote} - \mathbf{P}_{remote}\ _F \leq B_{\Delta P}^r$	$\ \tilde{\mathbf{X}}_{remote} - \mathbf{X}_{remote}\ _F \leq B_{\Delta X}^r$

Lemma D.2. [Lemma 1 in (Cong et al., 2021)] *An L -later GCN is L_F -Lipschitz smooth, i.e., $\|\nabla \mathcal{L}(\boldsymbol{\theta}_1) - \nabla \mathcal{L}(\boldsymbol{\theta}_2)\|_F \leq L_F \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|_F$.*

Lemma D.3. *Under Assumptions 5.1–5.3, and for any $l \in [L]$, the Frobenius norm of node embedding matrices, gradient passing from the l -th layer node embeddings to the $(l-1)$ -th are bounded, i.e.,*

$$\begin{aligned} \left\| \mathbf{H}_{local}^{(l),m} \right\|_F, \left\| \tilde{\mathbf{H}}_{local}^{(l),m} \right\|_F &\leq B_H^l, & \left\| \mathbf{H}_{full}^{(l),m} \right\|_F, \left\| \tilde{\mathbf{H}}_{full}^{(l),m} \right\|_F &\leq B_H^f, \\ \left\| \mathbf{D}_{local}^{(l),m} \right\|_F, \left\| \tilde{\mathbf{D}}_{local}^{(l),m} \right\|_F &\leq B_D^l, & \left\| \mathbf{D}_{full}^{(l),m} \right\|_F, \left\| \tilde{\mathbf{D}}_{full}^{(l),m} \right\|_F &\leq B_D^f, \end{aligned}$$

where

$$B_H^l, B_H^f = \max_{1 \leq l \leq L} (C_\sigma B_P B_W)^l B_X, \quad B_D^l, B_D^f = \max_{1 \leq l \leq L} (B_P B_W C_\sigma)^{L-l} C_l.$$

972 *Proof.*

$$\begin{aligned}
973 \quad \left\| \mathbf{H}_{local}^{(l),m} \right\|_F &= \left\| \sigma \left(\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \mathbf{W}^{(l),m} \right) \right\|_F \\
974 \quad &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \leq C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \right\| \left\| \mathbf{H}_{local}^{(l-1),m} \right\|_F \\
975 \quad &\stackrel{(b)}{\leq} C_\sigma B_W B_P \left\| \mathbf{H}_{local}^{(l-1),m} \right\|_F \leq (C_\sigma B_W B_P)^l \left\| \mathbf{X}^m \right\|_F \\
976 \quad &\stackrel{(c)}{\leq} (C_\sigma B_W B_P)^l B_X \leq \max_{1 \leq l \leq L} (C_\sigma B_W B_P)^l B_X,
\end{aligned}$$

977 where (a)–(c) results from Assumptions 5.2 and 5.3.

$$\begin{aligned}
982 \quad \left\| \widetilde{\mathbf{H}}_{local}^{(l),m} \right\|_F &= \left\| \sigma \left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \mathbf{W}^{(l),m} \right) \right\|_F \\
983 \quad &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F \leq C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} \right\| \left\| \widetilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F \\
984 \quad &\stackrel{(b)}{\leq} C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F \leq (C_\sigma B_W B_P)^l \left\| \mathbf{X}^m \right\|_F \\
985 \quad &\stackrel{(c)}{\leq} (C_\sigma B_W B_P)^l B_X \leq \max_{1 \leq l \leq L} (C_\sigma B_W B_P)^l B_X,
\end{aligned}$$

986 where (a)–(c) follow from Assumptions 5.2 and 5.3.

$$\begin{aligned}
991 \quad \left\| \mathbf{H}_{full}^{(l),m} \right\|_F &= \left\| \sigma \left(\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \mathbf{W}^{(l),m} \right) \right\|_F \\
992 \quad &\stackrel{(a)}{\leq} C_\sigma B_P B_W \left\| \mathbf{H}_{full}^{(l-1),m} \right\|_F \leq (C_\sigma B_P B_W)^l \left\| \mathbf{X}^m \right\|_F \\
993 \quad &\stackrel{(b)}{\leq} (C_\sigma B_P B_W)^l B_X \leq \max_{1 \leq l \leq L} (C_\sigma B_P B_W)^l B_X,
\end{aligned}$$

994 where (a) and (b) are because of Assumptions 5.2 and 5.3.

$$\begin{aligned}
1000 \quad \left\| \widetilde{\mathbf{H}}_{full}^{(l),m} \right\|_F &= \left\| \sigma \left(\left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} + \widetilde{\mathbf{P}}_{remote}^{(l),m} \widetilde{\mathbf{H}}_{remote}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right) \right\|_F \\
1001 \quad &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} + \widetilde{\mathbf{P}}_{remote}^{(l),m} \widetilde{\mathbf{H}}_{remote}^{(l-1),m} \right\|_F = C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{full}^{(l),m} \widetilde{\mathbf{H}}_{full}^{(l-1),m} \right\|_F \\
1002 \quad &\stackrel{(b)}{\leq} C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{full}^{(l-1),m} \right\|_F \leq (C_\sigma B_W B_P)^l \left\| \mathbf{X}^m \right\|_F \\
1003 \quad &\stackrel{(c)}{\leq} (C_\sigma B_W B_P)^l B_X \leq \max_{1 \leq l \leq L} (C_\sigma B_W B_P)^l B_X,
\end{aligned}$$

1004 where (a)–(c) follow from Assumptions 5.2 and 5.3.

$$\begin{aligned}
1008 \quad \left\| \mathbf{D}_{local}^{(l),m} \right\|_F &= \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
1009 \quad &\stackrel{(a)}{\leq} B_W C_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} \right\|_F \left\| \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
1010 \quad &\stackrel{(b)}{\leq} B_P B_W C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} \right\|_F \leq (B_P B_W C_\sigma)^{L-l} \left\| \mathbf{D}_{local}^{(L),m} \right\|_F \\
1011 \quad &\stackrel{(c)}{\leq} (B_P B_W C_\sigma)^{L-l} C_l \leq \max_{1 \leq l \leq L} (B_P B_W C_\sigma)^{L-l} C_l,
\end{aligned}$$

1026 where (a)–(c) are because of Assumptions 5.1–5.3.

$$\begin{aligned}
1027 \quad & \left\| \tilde{\mathbf{D}}_{local}^{(l),m} \right\|_F = \left\| \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
1028 \quad & \\
1029 \quad & \stackrel{(a)}{\leq} B_W C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l+1),m} \right\|_F \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} \right\|_F \\
1030 \quad & \stackrel{(b)}{\leq} B_P B_W C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} \right\|_F \leq (B_P B_W C_\sigma)^{L-l} \left\| \tilde{\mathbf{D}}_{local}^{(L),m} \right\|_F \\
1031 \quad & \stackrel{(c)}{\leq} (B_P B_W C_\sigma)^{L-l} C_l \leq \max_{1 \leq l \leq L} (B_P B_W C_\sigma)^{L-l} C_l,
\end{aligned}$$

1032 where (a)–(c) follow from Assumptions 5.1–5.3.

$$\begin{aligned}
1033 \quad & \left\| \mathbf{D}_{full}^{(l),m} \right\|_F = \left\| \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
1034 \quad & \stackrel{(a)}{\leq} B_P B_W C_\sigma \left\| \mathbf{D}_{full}^{(l+1),m} \right\|_F \leq (B_P B_W C_\sigma)^{L-l} \left\| \mathbf{D}_{full}^{(L),m} \right\|_F \\
1035 \quad & \stackrel{(b)}{\leq} (B_P B_W C_\sigma)^{L-l} C_l \leq \max_{1 \leq l \leq L} (B_P B_W C_\sigma)^{L-l} C_l,
\end{aligned}$$

1036 where (a) and (b) use Assumptions 5.1–5.3.

$$\begin{aligned}
1037 \quad & \left\| \tilde{\mathbf{D}}_{full}^{(l),m} \right\|_F = \left\| \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} + \tilde{\mathbf{P}}_{remote}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
1038 \quad & = \left\| \left[\tilde{\mathbf{P}}_{full}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
1039 \quad & \stackrel{(a)}{\leq} B_P B_W C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} \right\|_F \leq (B_P B_W C_\sigma)^{L-l} \left\| \tilde{\mathbf{D}}_{full}^{(L),m} \right\|_F \\
1040 \quad & \stackrel{(b)}{\leq} (B_P B_W C_\sigma)^{L-l} C_l \leq \max_{1 \leq l \leq L} (B_P B_W C_\sigma)^{L-l} C_l,
\end{aligned}$$

1041 where (a) and (b) utilize Assumptions 5.1–5.3.

□

1042 **Lemma D.4.** Under Assumptions 5.1–5.3, and for any $l \in [L]$, the errors caused by sampling are bounded, i.e.,

$$\begin{aligned}
1043 \quad & \left\| \tilde{\mathbf{H}}_{local}^{(l),m} - \mathbf{H}_{local}^{(l),m} \right\|_F \leq B_{\Delta H}^l, & \left\| \tilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F \leq B_{\Delta H}^f, \\
1044 \quad & \left\| \tilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \leq B_{\Delta D}^l, & \left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \leq B_{\Delta D}^f,
\end{aligned}$$

1045 where

$$\begin{aligned}
1046 \quad & B_{\Delta H}^l = \max_{1 \leq l \leq L} \left((C_\sigma B_W B_H^l B_{\Delta P}^l)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^l \right), \\
1047 \quad & B_{\Delta H}^f = \max_{1 \leq l \leq L} \left((C_\sigma B_W B_H^f B_{\Delta P}^f)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^f \right), \\
1048 \quad & B_{\Delta D}^l = \max_{1 \leq l \leq L} \left((B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l)^{L-l} \right. \\
1049 \quad & \quad \left. + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^l \right),
\end{aligned}$$

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104

$$B_{\Delta D}^f = \max_{1 \leq l \leq L} \left(\left(B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \right)^{L-l} + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^f \right).$$

Proof.

$$\begin{aligned} & \left\| \widetilde{\mathbf{H}}_{local}^{(l),m} - \mathbf{H}_{local}^{(l),m} \right\|_F \\ &= \left\| \sigma \left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \mathbf{W}^{(l),m} \right) - \sigma \left(\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right\|_F \\ &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ &\leq C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F + C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ &\stackrel{(b)}{\leq} C_\sigma B_W B_H^l \left\| \widetilde{\mathbf{P}}_{local}^{(l),m} - \mathbf{P}_{local}^{(l),m} \right\|_F + C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ &\stackrel{(c)}{\leq} C_\sigma B_W B_H^l B_{\Delta P}^l + C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ &\leq (C_\sigma B_W B_H^l B_{\Delta P}^l)^l + (C_\sigma B_W B_P)^l \left\| \widetilde{\mathbf{X}}_{local}^m - \mathbf{X}_{local}^m \right\|_F \\ &\stackrel{(d)}{\leq} (C_\sigma B_W B_H^l B_{\Delta P}^l)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^l \\ &\leq \max_{1 \leq l \leq L} \left((C_\sigma B_W B_H^l B_{\Delta P}^l)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^l \right), \end{aligned} \tag{9}$$

where (a) uses Assumptions 5.2 and 5.3, (b) is because of Assumption 5.3 and Lemma D.3, and (c) and (d) follow from Proposition D.1.

1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126

$$\begin{aligned} & \left\| \widetilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F \\ &= \left\| \sigma \left(\left(\widetilde{\mathbf{P}}_{local}^{(l),m} \widetilde{\mathbf{H}}_{local}^{(l-1),m} + \widetilde{\mathbf{P}}_{remote}^{(l),m} \widetilde{\mathbf{H}}_{remote}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right) - \sigma \left(\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right\|_F \\ &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{full}^{(l),m} \widetilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\leq C_\sigma B_W \left\| \widetilde{\mathbf{P}}_{full}^{(l),m} \widetilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \widetilde{\mathbf{H}}_{full}^{(l-1),m} \right\|_F + C_\sigma B_W \left\| \mathbf{P}_{full}^{(l),m} \widetilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\stackrel{(b)}{\leq} C_\sigma B_W B_H^f \left\| \widetilde{\mathbf{P}}_{full}^{(l),m} - \mathbf{P}_{full}^{(l),m} \right\|_F + C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\stackrel{(c)}{\leq} C_\sigma B_W B_H^f B_{\Delta P}^f + C_\sigma B_W B_P \left\| \widetilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\leq (C_\sigma B_W B_H^f B_{\Delta P}^f)^l + (C_\sigma B_W B_P)^l \left\| \widetilde{\mathbf{X}}_{full}^m - \mathbf{X}_{full}^m \right\|_F \\ &\stackrel{(d)}{\leq} (C_\sigma B_W B_H^f B_{\Delta P}^f)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^f \\ &\leq \max_{1 \leq l \leq L} \left((C_\sigma B_W B_H^f B_{\Delta P}^f)^l + (C_\sigma B_W B_P)^l B_{\Delta X}^f \right), \end{aligned} \tag{10}$$

where (a) follows from Assumptions 5.2 and 5.3, (b) is due to Assumption 5.3 and Lemma D.3, and (c) and (d) are because of Proposition D.1.

1127
1128
1129
1130
1131
1132
1133

$$\begin{aligned} & \left\| \widetilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \\ &= \left\| \left[\widetilde{\mathbf{P}}_{local}^{(l+1),m} \right]^\top \widetilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\widetilde{\mathbf{Z}}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \end{aligned}$$

$$\begin{aligned}
& - \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \Big\|_F \\
& \stackrel{(a)}{\leq} B_W \left\| \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \right\|_F \\
& \leq B_W \left\| \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) \right\|_F \\
& + B_W \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) \right\|_F \\
& + B_W \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \right\|_F \\
& \stackrel{(b)}{\leq} B_W B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l+1),m} - \mathbf{P}_{local}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} - \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
& + B_W B_P B_D^l \left\| \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l+1),m} \right) - \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \right\|_F \\
& \stackrel{(c)}{\leq} B_W B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l+1),m} - \mathbf{P}_{local}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} - \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^l L_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\
& \leq B_W B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l+1),m} - \mathbf{P}_{local}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} - \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^l L_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F \\
& + B_W^2 B_P B_D^l L_\sigma \left\| \mathbf{P}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\
& \stackrel{(d)}{\leq} B_W B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l+1),m} - \mathbf{P}_{local}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} - \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^l L_\sigma B_H^l \left\| \tilde{\mathbf{P}}_{local}^{(l),m} - \mathbf{P}_{local}^{(l),m} \right\|_F + B_W^2 B_P^2 B_D^l L_\sigma \left\| \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{H}_{local}^{(l-1),m} \right\|_F \\
& \stackrel{(e)}{\leq} B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l \\
& + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l+1),m} - \mathbf{D}_{local}^{(l+1),m} \right\|_F \\
& \leq \left(B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l \right)^{L-l} \\
& + \left(B_W B_P C_\sigma \right)^{L-l} \left\| \tilde{\mathbf{D}}_{local}^{(L),m} - \mathbf{D}_{local}^{(L),m} \right\|_F \\
& \stackrel{(f)}{\leq} \left(B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l \right)^{L-l} \\
& + \left(B_W B_P C_\sigma \right)^{L-l} L_l \left\| \tilde{\mathbf{H}}_{local}^{(L),m} - \mathbf{H}_{local}^{(L),m} \right\|_F \\
& \stackrel{(g)}{\leq} \left(B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l \right)^{L-l} \\
& + \left(B_W B_P C_\sigma \right)^{L-l} L_l B_{\Delta H}^l \\
& \leq \max_{1 \leq l \leq L} \left(\left(B_W B_D^l C_\sigma B_{\Delta P}^l + B_W^2 B_P B_D^l L_\sigma B_H^l B_{\Delta P}^l + B_W^2 B_P^2 B_D^l L_\sigma B_{\Delta H}^l \right)^{L-l} \right. \\
& \left. + \left(B_W B_P C_\sigma \right)^{L-l} L_l B_{\Delta H}^l \right),
\end{aligned}$$

where (a) uses Assumption 5.3, (b) is because of Assumptions 5.2 and 5.3 and Lemma D.3, (c) follows from Assumptions 5.2 and 5.3, (d) utilizes Assumption 5.3 and Lemma D.3, (e) results from Eq. (9) and Proposition D.1, (f) is because of Assumption 5.1, and (g) is due to Eq. (9).

$$\left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F$$

$$\begin{aligned}
&= \left\| \left[\tilde{\mathbf{P}}_{local}^{(l+1),m} + \tilde{\mathbf{P}}_{remote}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right. \\
&\quad \left. - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\
&\stackrel{(a)}{\leq} B_W \left\| \left[\tilde{\mathbf{P}}_{full}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\
&\leq B_W \left\| \left[\tilde{\mathbf{P}}_{full}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \right\|_F \\
&\quad + B_W \left\| \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) \right\|_F \\
&\quad + B_W \left\| \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\
&\stackrel{(b)}{\leq} B_W B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
&\quad + B_W B_P B_D^f \left\| \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l+1),m} \right) - \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\
&\stackrel{(c)}{\leq} B_W B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
&\quad + B_W^2 B_P B_D^f L_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} \tilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{full}^{(l),m} \right\|_F \\
&\leq B_W B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
&\quad + B_W^2 B_P B_D^f L_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} \tilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \tilde{\mathbf{H}}_{full}^{(l),m} \right\|_F \\
&\quad + B_W^2 B_P B_D^f L_\sigma \left\| \mathbf{P}_{full}^{(l+1),m} \tilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{full}^{(l),m} \right\|_F \\
&\stackrel{(d)}{\leq} B_W B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
&\quad + B_W^2 B_P B_D^f L_\sigma B_H^f \left\| \tilde{\mathbf{P}}_{full}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W^2 B_P^2 B_D^f L_\sigma \left\| \tilde{\mathbf{H}}_{full}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F \\
&\stackrel{(e)}{\leq} B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \\
&\quad + B_W B_P C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
&\leq \left(B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \right)^{L-l} \\
&\quad + (B_W B_P C_\sigma)^{L-l} \left\| \tilde{\mathbf{D}}_{full}^{(L),m} - \mathbf{D}_{full}^{(L),m} \right\|_F \\
&\stackrel{(f)}{\leq} \left(B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \right)^{L-l} \\
&\quad + (B_W B_P C_\sigma)^{L-l} L_l \left\| \tilde{\mathbf{H}}_{full}^{(L),m} - \mathbf{H}_{full}^{(L),m} \right\|_F \\
&\stackrel{(g)}{\leq} \left(B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \right)^{L-l} \\
&\quad + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^f \\
&\leq \max_{1 \leq l \leq L} \left(\left(B_W B_D^f C_\sigma B_{\Delta P}^f + B_W^2 B_P B_D^f L_\sigma B_H^f B_{\Delta P}^f + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^f \right)^{L-l} \right. \\
&\quad \left. + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^f \right),
\end{aligned}$$

where (a) is because of Assumption 5.3, (b) results from Assumptions 5.2 and 5.3 and Lemma D.3, (c) uses Assumptions 5.2 and 5.3, (d) is due to Assumption 5.3 and Lemma D.3, (e) follows from Eq. (10) and Proposition D.1, (f) utilizes Assumption 5.1, and (g) is because of Eq. (10). \square

Lemma D.5. *Under Assumptions 5.1–5.3, and for any $l \in [L]$, the errors caused by the information loss of the cross-client neighbors are bounded, i.e.,*

$$\left\| \mathbf{H}_{local}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F \leq B_{\Delta H}^r, \quad \left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \leq B_{\Delta D}^r,$$

where

$$\begin{aligned} B_{\Delta H}^r &= \max_{1 \leq l \leq L} \left((C_\sigma B_W B_P)^l B_X^r + (C_\sigma B_W B_H^f B_P)^l \right), \\ B_{\Delta D}^r &= \max_{1 \leq l \leq L} \left((B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r)^{L-l} \right. \\ &\quad \left. + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^r \right). \end{aligned}$$

Proof.

$$\begin{aligned} \left\| \mathbf{H}_{local}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F &= \left\| \sigma \left(\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right) \mathbf{W}^{(l),m} - \sigma \left(\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right) \mathbf{W}^{(l),m} \right\|_F \\ &\stackrel{(a)}{\leq} C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\leq C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\quad + C_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\stackrel{(b)}{\leq} C_\sigma B_W B_P \left\| \mathbf{H}_{local}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F + C_\sigma B_W B_H^f \left\| \mathbf{P}_{local}^{(l),m} - \mathbf{P}_{full}^{(l),m} \right\|_F \\ &\leq C_\sigma B_W B_P \left\| \mathbf{H}_{local}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F + C_\sigma B_W B_H^f \left\| \mathbf{P}_{remote}^{(l),m} \right\|_F \\ &\stackrel{(c)}{\leq} C_\sigma B_W B_P \left\| \mathbf{H}_{local}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F + C_\sigma B_W B_H^f B_P \\ &\leq (C_\sigma B_W B_P)^l \left\| \mathbf{X}_{local}^m - \mathbf{X}_{full}^m \right\|_F + (C_\sigma B_W B_H^f B_P)^l \\ &\stackrel{(d)}{\leq} (C_\sigma B_W B_P)^l B_X^r + (C_\sigma B_W B_H^f B_P)^l \\ &\leq \max_{1 \leq l \leq L} \left((C_\sigma B_W B_P)^l B_X^r + (C_\sigma B_W B_H^f B_P)^l \right), \end{aligned} \tag{11}$$

where (a) uses Assumptions 5.2 and 5.3, (b) is because of Assumption 5.3 and Lemma D.3, (c) follows from Assumption 5.3, and (d) is due to Proposition D.1.

$$\begin{aligned} &\left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\ &= \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right. \\ &\quad \left. - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \left[\mathbf{W}^{(l+1),m} \right]^\top \right\|_F \\ &\stackrel{(a)}{\leq} B_W \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\ &\leq B_W \left\| \left[\mathbf{P}_{local}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \right\|_F \end{aligned}$$

$$\begin{aligned}
& + B_W \left\| \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) \right\|_F \\
& + B_W \left\| \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) - \left[\mathbf{P}_{full}^{(l+1),m} \right]^\top \mathbf{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\
& \stackrel{(b)}{\leq} B_W B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
& + B_W B_P B_D^f \left\| \nabla \sigma \left(\mathbf{Z}_{local}^{(l+1),m} \right) - \nabla \sigma \left(\mathbf{Z}_{full}^{(l+1),m} \right) \right\|_F \\
& \stackrel{(c)}{\leq} B_W B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^f L_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} \mathbf{H}_{local}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{full}^{(l),m} \right\|_F \\
& \leq B_W B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^f L_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} \mathbf{H}_{local}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{local}^{(l),m} \right\|_F \\
& + B_W^2 B_P B_D^f L_\sigma \left\| \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{local}^{(l),m} - \mathbf{P}_{full}^{(l+1),m} \mathbf{H}_{full}^{(l),m} \right\|_F \\
& \stackrel{(d)}{\leq} B_W B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W B_P C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
& + B_W^2 B_P B_D^f L_\sigma B_H^l \left\| \mathbf{P}_{local}^{(l+1),m} - \mathbf{P}_{full}^{(l+1),m} \right\|_F + B_W^2 B_P^2 B_D^f L_\sigma \left\| \mathbf{H}_{local}^{(l),m} - \mathbf{H}_{full}^{(l),m} \right\|_F \\
& \stackrel{(e)}{\leq} B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r \\
& + B_W B_P C_\sigma \left\| \mathbf{D}_{local}^{(l+1),m} - \mathbf{D}_{full}^{(l+1),m} \right\|_F \\
& \leq \left(B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r \right)^{L-l} \\
& + (B_W B_P C_\sigma)^{L-l} \left\| \mathbf{D}_{local}^{(L),m} - \mathbf{D}_{full}^{(L),m} \right\|_F \\
& \stackrel{(f)}{\leq} \left(B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r \right)^{L-l} \\
& + (B_W B_P C_\sigma)^{L-l} L_l \left\| \mathbf{H}_{local}^{(L),m} - \mathbf{H}_{full}^{(L),m} \right\|_F \\
& \stackrel{(g)}{\leq} \left(B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r \right)^{L-l} \\
& + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^r \\
& \leq \max_{1 \leq l \leq L} \left(\left(B_W B_D^l C_\sigma B_P + B_W^2 B_P^2 B_D^f L_\sigma B_H^l + B_W^2 B_P^2 B_D^f L_\sigma B_{\Delta H}^r \right)^{L-l} \right. \\
& \left. + (B_W B_P C_\sigma)^{L-l} L_l B_{\Delta H}^r \right),
\end{aligned}$$

where (a) follows from Assumption 5.3, (b) uses Assumptions 5.2 and 5.3 and Lemma D.3, (c) is because of Assumptions 5.2 and 5.3, (d) results from Assumption 5.3 and Lemma D.3, (e) is due to Assumption 5.3 and Eq. (11), (f) utilizes Assumption 5.1, and (g) uses Eq. (11). \square

D.3 ERRORS OF STOCHASTIC GRADIENTS

Lemma D.6. *Under Assumptions 5.1–5.3, the errors between the stochastic gradients and the full gradients are bounded as follows:*

$$\left\| \nabla F_{local}^m(\boldsymbol{\theta}^m) - \nabla \tilde{F}_{local}^m(\boldsymbol{\theta}^m) \right\|_F \leq L B_{\Delta G}^l, \quad \left\| \nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla \tilde{F}_{full}^m(\boldsymbol{\theta}^m) \right\|_F \leq L B_{\Delta G}^f,$$

1350 where

$$1351 B_{\Delta G}^l = \max_{1 \leq l \leq L} \left((B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_H^l B_{\Delta P}^l + B_P B_H^l C_\sigma B_{\Delta D}^l \right. \\ 1352 \left. + (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_P B_{\Delta H}^l \right), \quad (12)$$

$$1353 B_{\Delta G}^f = \max_{1 \leq l \leq L} \left((B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W) B_H^f B_{\Delta P}^f + B_P B_H^f C_\sigma B_{\Delta D}^f \right. \\ 1354 \left. + (B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W) B_P B_{\Delta H}^f \right) \quad (13)$$

1362 *Proof.*

$$1363 \left\| \tilde{\mathbf{G}}_{local}^{(l),m} - \mathbf{G}_{local}^{(l),m} \right\|_F \\ 1364 = \left\| \left[\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) - \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \right\|_F \\ 1365 \leq \left\| \left[\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) - \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) \right\|_F \\ 1366 + \left\| \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) - \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) \right\|_F \\ 1367 + \left\| \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) - \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \right\|_F \\ 1368 \stackrel{(a)}{\leq} B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F + B_P B_H^l C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \\ 1369 + B_P B_H^l B_D^l \left\| \nabla \sigma \left(\tilde{\mathbf{Z}}_{local}^{(l),m} \right) - \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \right\|_F \\ 1370 \stackrel{(b)}{\leq} B_D^l C_\sigma \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F + B_P B_H^l C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \\ 1371 + B_P B_H^l B_D^l L_\sigma B_W \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ 1372 \leq (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) \left\| \tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} \right\|_F \\ 1373 + (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) \left\| \mathbf{P}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ 1374 + B_P B_H^l C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \\ 1375 \stackrel{(c)}{\leq} (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_H^l \left\| \tilde{\mathbf{P}}_{local}^{(l),m} - \mathbf{P}_{local}^{(l),m} \right\|_F + B_P B_H^l C_\sigma \left\| \tilde{\mathbf{D}}_{local}^{(l),m} - \mathbf{D}_{local}^{(l),m} \right\|_F \\ 1376 + (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_P \left\| \tilde{\mathbf{H}}_{local}^{(l-1),m} - \mathbf{H}_{local}^{(l-1),m} \right\|_F \\ 1377 \stackrel{(d)}{\leq} (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_H^l B_{\Delta P}^l + B_P B_H^l C_\sigma B_{\Delta D}^l \\ 1378 + (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_P B_{\Delta H}^l \\ 1379 \leq \max_{1 \leq l \leq L} \left((B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_H^l B_{\Delta P}^l + B_P B_H^l C_\sigma B_{\Delta D}^l \right. \\ 1380 \left. + (B_D^l C_\sigma + B_P B_H^l B_D^l L_\sigma B_W) B_P B_{\Delta H}^l \right) := B_{\Delta G}^l, \\ 1381 \\ 1382 \\ 1383 \\ 1384 \\ 1385 \\ 1386 \\ 1387 \\ 1388 \\ 1389 \\ 1390 \\ 1391 \\ 1392 \\ 1393 \\ 1394 \\ 1395 \\ 1396 \\ 1397 \\ 1398 \\ 1399 \\ 1400 \\ 1401$$

1402 where (a) follows from Assumptions 5.2 and 5.3 and Lemma D.3, (b) is because of Assumptions 5.2
1403 and 5.3, (c) uses Assumption 5.3 and Lemma D.3, and (d) results from Lemma D.4 and Proposi-
tion D.1.

When client m performs local training with only its local data, the error between the stochastic gradient and the full-gradient can be bounded as:

$$\begin{aligned}
& \left\| \nabla F_{local}^m(\boldsymbol{\theta}^m) - \nabla \tilde{F}_{local}^m(\boldsymbol{\theta}^m) \right\|_F = \sum_{l=1}^L \left\| \mathbf{G}_{local}^{(l),m} - \tilde{\mathbf{G}}_{local}^{(l),m} \right\|_F \leq LB_{\Delta G}^l. \\
& \left\| \tilde{\mathbf{G}}_{full}^{(l),m} - \mathbf{G}_{full}^{(l),m} \right\|_F \\
& = \left\| \left[\tilde{\mathbf{P}}_{local}^{(l),m} \tilde{\mathbf{H}}_{local}^{(l-1),m} + \tilde{\mathbf{P}}_{remote}^{(l),m} \tilde{\mathbf{H}}_{remote}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) \right. \\
& \quad \left. - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\
& \leq \left\| \left[\tilde{\mathbf{P}}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) \right\|_F \\
& \quad + \left\| \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \tilde{\mathbf{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) \right\|_F \\
& \quad + \left\| \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\
& \stackrel{(a)}{\leq} B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F + B_P B_H^f C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\
& \quad + B_P B_H^f B_D^f \left\| \nabla \sigma \left(\tilde{\mathbf{Z}}_{full}^{(l),m} \right) - \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\
& \stackrel{(b)}{\leq} B_D^f C_\sigma \left\| \tilde{\mathbf{P}}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F + B_P B_H^f C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\
& \quad + B_P B_H^f B_D^f L_\sigma B_W \left\| \tilde{\mathbf{P}}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\
& \leq \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) \left\| \tilde{\mathbf{P}}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} \right\|_F \\
& \quad + \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) \left\| \mathbf{P}_{full}^{(l),m} \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\
& \quad + B_P B_H^f C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\
& \stackrel{(c)}{\leq} \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f \left\| \tilde{\mathbf{P}}_{full}^{(l),m} - \mathbf{P}_{full}^{(l),m} \right\|_F + B_P B_H^f C_\sigma \left\| \tilde{\mathbf{D}}_{full}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\
& \quad + \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P \left\| \tilde{\mathbf{H}}_{full}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F \\
& \stackrel{(d)}{\leq} \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f B_{\Delta P}^f + B_P B_H^f C_\sigma B_{\Delta D}^f \\
& \quad + \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P B_{\Delta H}^f \\
& \leq \max_{1 \leq l \leq L} \left(\left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f B_{\Delta P}^f + B_P B_H^f C_\sigma B_{\Delta D}^f \right. \\
& \quad \left. + \left(B_D^f C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P B_{\Delta H}^f \right) := B_{\Delta G}^f,
\end{aligned}$$

where (a) results from Assumptions 5.2 and 5.3 and Lemma D.3, (b) uses Assumptions 5.2 and 5.3, (c) is due to Assumption 5.3 and Lemma D.3, and (d) is because of Lemma D.4 and Proposition D.1.

When client m conducts cross-client training using its local data and the cross-client neighbors, the error between the stochastic gradient and the full-gradient can be bounded as:

$$\left\| \nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla \tilde{F}_{full}^m(\boldsymbol{\theta}^m) \right\|_F = \sum_{l=1}^L \left\| \mathbf{G}_{full}^{(l),m} - \tilde{\mathbf{G}}_{full}^{(l),m} \right\|_F \leq LB_{\Delta G}^f.$$

1458

□

1459

1460

1461

1462

Lemma D.7. *Under Assumptions 5.1–5.3, the error between the full gradient computed with both the local graph data and the cross-client neighbors and the full gradient computed with only the local graph data is upper-bounded as follows:*

1463

$$\|\nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla F_{local}^m(\boldsymbol{\theta}^m)\|_F \leq LB_{\Delta G}^r,$$

1464

1465 *where*

1466

1467

1468

1469

1470

$$\begin{aligned} B_{\Delta G}^r = \max_{1 \leq l \leq L} & \left((B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W) B_P B_{\Delta H}^r + B_P B_H^f C_\sigma B_{\Delta D}^r \right. \\ & \left. + (B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W) B_H^f B_P \right) \end{aligned} \quad (14)$$

1471

1472 *Proof.*

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

1487

1488

1489

1490

1491

1492

1493

1494

1495

1496

1497

1498

1499

1500

1501

1502

1503

1504

1505

1506

1507

1508

1509

1510

1511

$$\begin{aligned} & \left\| \mathbf{G}_{local}^{(l),m} - \mathbf{G}_{full}^{(l),m} \right\|_F \\ &= \left\| \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\ &\leq \left\| \left[\mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \right\|_F \\ &+ \left\| \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \right\|_F \\ &+ \left\| \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right]^\top \mathbf{D}_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\ &\stackrel{(a)}{\leq} B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F + B_P B_H^f C_\sigma \left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\ &+ B_P B_H^f B_D^f \left\| \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \right\|_F \\ &\stackrel{(b)}{\leq} B_D^l C_\sigma \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F + B_P B_H^f C_\sigma \left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\ &+ B_P B_H^f B_D^f L_\sigma B_W \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &\leq \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{local}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &+ \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) \left\| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &+ B_P B_H^f C_\sigma \left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\ &\stackrel{(c)}{\leq} \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P \left\| \mathbf{H}_{local}^{(l-1),m} - \mathbf{H}_{full}^{(l-1),m} \right\|_F \\ &+ \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f \left\| \mathbf{P}_{local}^{(l),m} - \mathbf{P}_{full}^{(l),m} \right\|_F + B_P B_H^f C_\sigma \left\| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \right\|_F \\ &\stackrel{(d)}{\leq} \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P B_{\Delta H}^r + B_P B_H^f C_\sigma B_{\Delta D}^r \\ &+ \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f B_P \\ &\leq \max_{1 \leq l \leq L} \left(\left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_P B_{\Delta H}^r + B_P B_H^f C_\sigma B_{\Delta D}^r \right. \\ &\left. + \left(B_D^l C_\sigma + B_P B_H^f B_D^f L_\sigma B_W \right) B_H^f B_P \right) = B_{\Delta G}^r, \end{aligned}$$

where (a) is because of Assumptions 5.2 and 5.3 and Lemma D.3, (b) uses Assumptions 5.2 and 5.3, (c) follow from Assumption 5.3 and Lemma D.3, and (d) results from Assumption 5.3 and Lemma D.5.

The error between the full gradient computed with both the local graph data and the cross-client neighbors and the full gradient computed with only the local graph data is bounded as follows:

$$\|\nabla F_{full}^m(\boldsymbol{\theta}^m) - \nabla F_{local}^m(\boldsymbol{\theta}^m)\|_F = \sum_{l=1}^L \left\| \mathbf{G}_{local}^{(l),m} - \mathbf{G}_{full}^{(l),m} \right\|_F \leq LB_{\Delta G}^r.$$

□

D.4 MAIN PROOF OF THEOREM 5.6

Theorem D.8. *Under Assumptions 5.1–5.3, choose step-size $\alpha = \min\{\sqrt{M}/\sqrt{T}, 1/L_F\}$, where L_F is the smoothness constant given in Lemma D.2. The output of Swift-FedGNN with a L -layer GNN satisfies:*

$$\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla \mathcal{L}(\boldsymbol{\theta}_t)\|^2 \leq \frac{2}{\sqrt{MT}} (\mathcal{L}(\boldsymbol{\theta}_0) - \mathcal{L}(\boldsymbol{\theta}^*)) + \left(1 - \frac{K}{IM}\right) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2 + \frac{K}{IM} L^2 (B_{\Delta G}^f)^2.$$

Proof.

$$\mathcal{L}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}(\boldsymbol{\theta}_t)$$

$$\stackrel{(a)}{\leq} \langle \nabla \mathcal{L}(\boldsymbol{\theta}_t), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t \rangle + \frac{L_F}{2} \|\boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_t\|^2$$

$$\stackrel{(b)}{=} -\alpha \left\langle \nabla \mathcal{L}(\boldsymbol{\theta}_t), \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\rangle + \frac{L_F}{2} \alpha^2 \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$\stackrel{(c)}{=} -\frac{\alpha}{2} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 - \frac{\alpha}{2} \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2 + \frac{\alpha}{2} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) - \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$+ \frac{L_F}{2} \alpha^2 \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$= -\frac{\alpha}{2} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 - \frac{\alpha}{2} \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2 + \frac{\alpha}{2} \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} (\nabla F^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m)) \right\|^2$$

$$+ \frac{L_F}{2} \alpha^2 \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$\stackrel{(d)}{\leq} -\frac{\alpha}{2} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 - \frac{\alpha}{2} \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2 + \frac{\alpha}{2} \frac{1}{M} \sum_{m \in \mathcal{M}} \left\| \nabla F^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$+ \frac{L_F}{2} \alpha^2 \left\| \frac{1}{M} \sum_{m \in \mathcal{M}} \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2$$

$$\stackrel{(e)}{\leq} -\frac{\alpha}{2} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 + \frac{\alpha}{2} \frac{1}{M} \sum_{m \in \mathcal{M}} \left\| \nabla F^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\|^2, \quad (15)$$

where (a) follows from Lemma D.2, (b) is because of the update rule in Swift-FedGNN, (c) uses $\langle \mathbf{x}, \mathbf{y} \rangle = \frac{1}{2} \|\mathbf{x}\|^2 + \frac{1}{2} \|\mathbf{y}\|^2 - \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|^2$, (d) utilizes $\|\sum_{i=1}^n \mathbf{x}_i\|^2 \leq n \sum_{i=1}^n \|\mathbf{x}_i\|^2$, and (e) is due to the choice of $\alpha \leq 1/L_F$.

When $t \in [(n_t - 1)I + 1, n_t I - 1] \cap \mathbb{Z}$, where $n_t = \{1, 2, \dots\}$, Swift-FedGNN conducts local training for all clients $m \in \mathcal{M}$. Thus,

$$\begin{aligned} \left\| \nabla F^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\| &= \left\| \nabla F_{full}^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}_{local}^m(\boldsymbol{\theta}_t^m) \right\| \\ &\leq \left\| \nabla F_{full}^m(\boldsymbol{\theta}_t^m) - \nabla F_{local}^m(\boldsymbol{\theta}_t^m) \right\| + \left\| \nabla F_{local}^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}_{local}^m(\boldsymbol{\theta}_t^m) \right\| \end{aligned}$$

$$\stackrel{(a)}{\leq} LB_{\Delta G}^r + LB_{\Delta G}^l, \quad (16)$$

where (a) follows from Lemmas D.6 and D.7.

When $t = n_t I$, where $n_t = \{1, 2, \dots\}$, Swift-FedGNN performs local training for clients $m \in \mathcal{M} \setminus \mathcal{K}$, and thus the inequality (16) holds for these clients. The randomly sampled clients $m \in \mathcal{K}$ conduct cross-client training, and thus

$$\left\| \nabla F^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}^m(\boldsymbol{\theta}_t^m) \right\| = \left\| \nabla F_{full}^m(\boldsymbol{\theta}_t^m) - \nabla \tilde{F}_{full}^m(\boldsymbol{\theta}_t^m) \right\| \stackrel{(a)}{\leq} LB_{\Delta G}^f,$$

where (a) uses Lemma D.6.

Telescoping (15) from $i = (n_t - 1)I + 1$ to $n_t I$, we have

$$\begin{aligned} & \sum_{i=(n_t-1)I+1}^{n_t I} (\mathcal{L}(\boldsymbol{\theta}_{i+1}) - \mathcal{L}(\boldsymbol{\theta}_i)) \\ & \leq -\frac{\alpha}{2} \sum_{i=(n_t-1)I+1}^{n_t I} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_i) \right\|^2 + \frac{\alpha}{2} (I-1) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2 + \frac{\alpha}{2M} K L^2 (B_{\Delta G}^f)^2 \\ & \quad + \frac{\alpha}{2M} (M-K) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2. \end{aligned}$$

Choosing $T = n_t I$ yields

$$\begin{aligned} & \sum_{t=0}^{T-1} (\mathcal{L}(\boldsymbol{\theta}_{t+1}) - \mathcal{L}(\boldsymbol{\theta}_t)) \\ & \leq -\frac{\alpha}{2} \sum_{t=0}^{T-1} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 + \frac{\alpha}{2} (T - n_t) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2 + n_t \frac{\alpha}{2M} K L^2 (B_{\Delta G}^f)^2 \\ & \quad + n_t \frac{\alpha}{2M} (M-K) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2. \end{aligned}$$

Rearranging the terms and multiplying both sides by $2/\alpha$, we get

$$\begin{aligned} & \sum_{t=0}^{T-1} \left\| \nabla \mathcal{L}(\boldsymbol{\theta}_t) \right\|^2 \\ & \leq \frac{2}{\alpha} \sum_{t=0}^{T-1} (\mathcal{L}(\boldsymbol{\theta}_t) - \mathcal{L}(\boldsymbol{\theta}_{t+1})) + (T - n_t) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2 + \frac{n_t}{M} K L^2 (B_{\Delta G}^f)^2 \\ & \quad + \frac{n_t}{M} (M-K) L^2 (B_{\Delta G}^l + B_{\Delta G}^r)^2. \end{aligned}$$

Dividing both sides by T and choosing $\alpha = \sqrt{M}/\sqrt{T}$ completes the proof of Theorem 5.6. \square

1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619