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

Anonymous authors

000

001

002003004

010

012

013

014

015

016

017

018

019

021

023

024

025

026

028

029

031

033 034

036

037

038

040 041

042

043

044

046

047

051

052

Paper under double-blind review

ABSTRACT

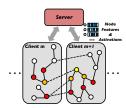
Graph neural networks (GNNs) have achieved great success in a wide variety of graph-based learning applications. While distributed GNN training with samplingbased mini-batches expedites learning on large graphs, it is not applicable to geo-distributed data that must remain on-site to preserve privacy. On the other hand, federated learning (FL) has been widely used to enable privacy-preserving training under data parallelism. However, applying FL directly to GNNs either results in cross-client neighbor information loss or 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 federated graph learning (FGL) algorithmic framework called Swift-FedGNN that primarily performs efficient parallel local training and periodically conducts 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 FGL approaches in terms of efficiency, while achieving comparable accuracy.

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), recommendation systems Ying et al. (2018); Wang et al. (2019a), and drug discovery Wang et al. (2022b); Do et al. (2019). GNN learns high-level graph representations by iteratively aggregating 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).

Real-world graph datasets can be extensive in scale (*e.g.*, Microsoft Academic Graph Wang et al. (2020) with over 100 million nodes) and often reside across geo-distributed sites where data protection laws prohibit direct data sharing Yao et al. (2023a). Single devices (*e.g.*, GPUs) often lack the capacity for training such large-scale datasets, which leads to *a compelling need* for distributed graph learning (DGL) Fey & Lenssen (2019); Zheng et al. (2020). However, the common DGL paradigm, consisting of subgraph sampling Zeng et al. (2020) and mini-batch training Luo et al. (2022), requires direct data sharing among workers, which conflicts with privacy regulations.

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



056

060 061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

078

079

081

082

083

084

085

087

880 089

091

092

094

096

098

099 100

101

102

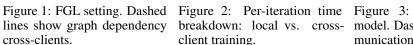
103

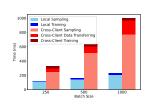
104

105

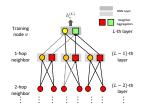
106

107





client training.



Federated GNN breakdown: local vs. cross- model. Dashed lines show communication between clients.

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: Despite the appeal of leveraging a trusted server for federated graph learning (FGL) Zhang et al. (2021), there remain several non-trivial challenges that hinder efficient and effective cross-client training. Specifically, we highlight the following major technical challenges:

Large Overhead from a Naive Design. A straightforward approach uses a trusted server to gather graph data and perform subgraph sampling and neighbor aggregation for each client Zhang et al. (2021). For instance, in a healthcare setting with multiple hospitals and one central authority, patient data stay locally due to privacy regulation. The central authority acting as the trusted server must coordinate all subgraph sampling and part of the training operation (i.e., neighbor aggregation) (see Figure 1). As shown in Figure 2, training a two-layer GNN (with sampling fanout values 15 and 10 for the two layers) on the Amazon product co-purchasing dataset Leskovec et al. (2007) under 1 Gbps network bandwidth, this approach leads to significant communication overhead: the server exchanges large amounts of node and edge information with each hospital sequentially, causing cross-client sampling and communication time to dominate the total training time, making it five times slower than purely local training.

Communication and Memory Overheads from Cross-Client Neighbors. While some methods ignore cross-client neighbors He et al. (2021a) or assume overlapping nodes Wu et al. (2021), these assumptions often fail in geo-distributed graphs (e.g., patients visiting multiple hospitals). Alternatives that preserve cross-client neighbor information Zhang et al. (2021); Du & Wu (2022); Yao et al. (2023a) require significant data transfers among clients—leading to high communication costs—and compel each client to store additional graph structure and features for these neighbors. This not only creates memory-intensive requirements but could also potentially violate data privacy constraints. Hence, mitigating these cross-client overheads (both communication and storage) is crucial to achieve efficient, privacy-preserving FGL (see detailed discussions in Section 2).

- 3) Our Contributions: The key contribution of this paper is that, by addressing the above challenges, we develop a mini-batch-based and sampling-based FGL framework called Swift-FedGNN. The main results and technical contributions of this paper are as follows:
- We develop Swift-FedGNN, a communication- and sample-efficient mini-batch FGL algorithm for geo-distributed graphs. In Swift-FedGNN, clients *primarily* conduct local training in parallel, performing cross-client training only occasionally among sampled clients, thereby reducing sampling and communication overhead while preserving minimal information loss. The cross-client neighbor information is aggregated at remote clients before communicating to the server and accumulated one more time before transferring to the training client, further minimizing data transfer cost and enhancing privacy by ensuring only aggregated neighborhood features - never raw node features - are exchanged.
- We conduct rigorous theoretical convergence analysis for Swift-FedGNN, which is highly nontrivial due to biased stochastic gradients and structural entanglement (neighbor aggregation intertwined with non-linear transformations across multiple layers) in GNNs. In stark contrast to existing works in the literature that made strong assumptions on the biases of stochastic gradients (e.g., unbiased Chen et al. (2018) or consistent Chen & Luss (2018) gradient), for the first time in the literature, we are able to bound stochastic gradient approximation errors rather than resorting to these unrealistic assumptions in practice, offering insights of independent theoretical interest.
- We show that the biased stochastic gradients in GNNs—arising from missing cross-client neighbors and neighbor sampling—are *positively correlated* with the network depth, which is *unique* to FGL.

By putting the above insights together, we show that Swift-FedGNN achieves a convergence rate of $\mathcal{O}\left(T^{-1/2}\right)$, which *matches* the state-of-the-art (SOTA) convergence rate of sampling-based GNN methods (hence low communication and sample complexities), *despite* operating in the far more challenging FL setting with much less frequent information exchanges among clients.

• We conduct extensive experiments on real-world graph datasets to evaluate the performance of Swift-FedGNN. The results show that Swift-FedGNN outperforms the SOTA FGL algorithms in terms of *efficiency*, achieving ×4 speed-up and competitive accuracy.

2 RELATED WORK

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

- 1) Distributed Graph Learning: Distributed graph learning framework (*e.g.*, DistDGL Wang et al. (2019b); Zheng et al. (2020), Pytorch Geometric Fey & Lenssen (2019), AliGraph Zhao et al. (2019) and Dorylus Thorpe et al. (2021)) have been developed to train large-scale graph datasets via cross-device sampling and direct worker-to-worker communication, and often spend up to 80% of the total training time on data communication Gandhi & Iyer (2021). Although various optimizations (graph partitioning Zheng et al. (2020), caching Liu et al. (2023); Zhang et al. (2023), communication strategies Cai et al. (2021); Luo et al. (2022), parallel training Gandhi & Iyer (2021); Wan et al. (2022); Du et al. (2024)) have been proposed to expedite DGL, they commonly require direct data sharing between workers, violating data privacy constraints in geo-distributed settings. To our knowledge, LLCG Ramezani et al. (2022) is the only DGL framework that avoids transferring node features between workers, making it potentially applicable to geo-distributed graphs. In LLCG, each worker trains only on its local graph partition. To address missing cross-device neighbor information, LLCG employs a central server to periodically perform full-neighbor training with neighbor aggregation across all workers. However, this approach imposes significant communication overhead on the server, which needs to communicate with every worker to perform the full-neighbor training.
- 2) Federated Graph Learning: To date, the research on federated graph learning 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. Zhang et al. (2021) mitigates 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. While Yao et al. (2023a) employs a one-time exchange of full cross-client neighbor information prior to training, this design relies on full-graph training and causes significant per-client memory overhead, making it impractical for large-scale graphs. Adapting it to sampling-based FGL would require per-iteration cross-client exchanges (since each mini-batch has a different training node set and sampled neighbors), further exacerbating communication overhead.

Du & Wu (2022) uses sparse cross-client neighbor sampling to supplement the lost information of the cross-client neighbors and reduce the communication overhead, which is most related to ours. Each client periodically samples and exchanges these neighbors with other clients, reusing the most recent sampled neighbors in between exchanges. However, as training progresses, the frequency of information exchange increases, leading to higher communication costs. Furthermore, privacy constraints are relaxed by allowing direct client-to-client data transfers and caching, and repeatedly reusing the same neighbor data introduces bias that degrades performance. In contrast, our 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 reduce communication costs.

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 B. 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 C.

Consider a graph $\mathcal{G}\left(\mathcal{V},\mathcal{E}\right)$, 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 $\mathcal{G}^m\left(\mathcal{V}^m,\mathcal{E}^m\right)$. 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 $\boldsymbol{x}_v^m \in \mathbb{R}^d$, 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 FGL, the clients collaboratively learn a model with distributed graph data and under the coordination of the central server. Typically, the clients 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 FGL is to solve the following optimization problem:

$$\min \mathcal{L}(\boldsymbol{\theta}) := \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} F^{m}(\boldsymbol{\theta}) = \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \frac{1}{|\mathcal{V}_{B}^{m}|} \sum_{v \in \mathcal{V}_{B}^{m}} \ell^{m} \left(\boldsymbol{h}_{v}^{(L), m}, y_{v}^{m}\right), \tag{1}$$

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 $\boldsymbol{\theta} := \left\{ \boldsymbol{W}^{(l)} \right\}_{l=1}^L$ 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 FGL, 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 $\bar{\mathcal{M}}(v) \subseteq \mathcal{M} \setminus \{m(v)\}$. 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 $\bar{m}(v)$, and then updates the embedding of node v, as follows:

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

$$\boldsymbol{h}_{v}^{(l),m(v)} = \sigma\left(\boldsymbol{W}^{(l)} \cdot \operatorname{COMB}\left(\boldsymbol{h}_{v}^{(l-1),m(v)}, \boldsymbol{h}_{\mathcal{N}(v)}^{(l)}\right)\right),$$
(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 $\bar{m}(v)$, $h_{\mathcal{N}(v)}^{(l)}$ is the aggregated embedding from node v's neighbors, $h_v^{(l),m(v)}$ is the embedding of node v located on client m(v) and is initialized as $h_v^{(0),m(v)}=x_v^{m(v)}$, $W^{(l)}$ represents the weight matrix at l-th layer, σ (·) corresponds to an activation function (e.g., ReLU), AGG (·) is an aggregation function (e.g., mean), and COMB (·) is a combination function (e.g., concatenation). Compared to DGL where clients can directly transfer node features, the key difference in FGL 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 FGL. 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

216 **Algorithm 1:** Swift-FedGNN Algorithm. 217 **Input:** Initial parameters θ_0 , learning rate α , 218 and correction frequency I219 for t = 0 to T - 1 do 220 if $t \bmod I = 0$ then 221 Randomly sample $|\mathcal{K}|$ clients 222 for $m \in \mathcal{M}$ in parallel do if $m \in \mathcal{K}$ then 223 Client update with local 224 graph and cross-client 225 neighbors using Algo-226 rithm 2 227 else 228 Client update with local graph using Algorithm 3 229 230 else 231 for $m \in \mathcal{M}$ in parallel do 232 Client update with local graph 233 using Algorithm 3 Server: 235 Aggregate and update global model pa-236 rameter as: 237 $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \alpha \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \nabla \widetilde{F}^m \left(\boldsymbol{\theta}_t^m \right)$

238239

240

241

242

243

244

245

246

247

248249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267 268

269

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

Receive global parameter $\boldsymbol{\theta}_t^m = \boldsymbol{\theta}_t$ Construct a mini-batch \mathcal{B}_v^m of nodes

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

for l=1 to L do /* Derive l-th layer embedding of node $v\in\mathcal{B}^m_v$ if l=L, otherwise $v\in\mathcal{S}^{(l)}$

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

Aggregate the neighbor embeddings using Eq. (5)

Send the aggregated embedding $h_{\mathcal{N}(v)}^{(l),\bar{m}(v)}$ to server

Server:

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

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

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

Compute stochastic gradient $\nabla \widetilde{F}^m \left(\boldsymbol{\theta}_t^m \right)$ and send to server

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

Compute stochastic gradient $\nabla \widetilde{F}^m (\boldsymbol{\theta}_t^m)$ and send to server

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 $\mathcal{S} = \left\{\mathcal{S}^{(l)}\right\}_{l=0}^{L-1}$, 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 $v \in \mathcal{S}^{(l)}$), client m first conducts the neighbor aggregation for node v based on the sampled neighbors using:

$$\boldsymbol{h}_{\mathcal{N}(v)}^{(l)} = AGG\left(\left\{\boldsymbol{h}_{u}^{(l-1),m} \mid u \in \widetilde{\mathcal{N}}^{m}\left(v\right)\right\}\right),\tag{3}$$

where $\widetilde{\mathcal{N}}^m(v)$ denotes a set of the sampled neighbors located on client m for node v, $\widetilde{\mathcal{N}}^m(v) \subseteq \mathcal{S}^{(l-1)}$, and $\widetilde{\mathcal{N}}^m(v) \subseteq \mathcal{N}^m(v)$. 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:

 $\boldsymbol{h}_{v}^{(l),m} = \sigma\left(\boldsymbol{W}_{t}^{(l),m} \cdot \text{COMB}(\boldsymbol{h}_{v}^{(l-1),m}, \boldsymbol{h}_{\mathcal{N}(v)}^{(l)})\right). \tag{4}$

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 K to denote the set of K clients, where $K \subset \mathcal{M}$. The remaining clients perform local training as shown in Algorithm 3. Algorithm 2 details the cross-client training process for client $m \in K$. Rather than directly exchanging node features between clients, Swift-FedGNN partitions GNN training between the clients and the server. We offload the aggregation of node features and intermediate activations at each GNN layer to the server and remote

¹Note that the operation offloading in Swift-FedGNN only supports element-wise (*e.g.*, mean, sum, max) operations, *e.g.*, GraphSAGE Hamilton et al. (2017), GIN Xu et al. (2019), 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

271

272

273

274

275

276

277

278

279

281

282

283

284

285 286

287

288 289

290 291

292

293

295

296

297

298

299

300

301 302

303

304

305

306

307

308

309

310

311 312 313

314 315

316

317

318

319

320

321

322

323

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

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

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

$$\boldsymbol{h}_{\mathcal{N}(u)}^{(l),\bar{m}(v)} = \operatorname{AGG}\left(\left\{\boldsymbol{h}_{u}^{(l-1),\bar{m}(v)} \mid u \in \widetilde{\mathcal{N}}^{\bar{m}(v)}(v)\right\}\right). \tag{5}$$

 $\boldsymbol{h}_{\mathcal{N}(v)}^{(l),\bar{m}(v)} = \mathrm{AGG}\left(\left\{\boldsymbol{h}_{u}^{(l-1),\bar{m}(v)} \mid u \in \widetilde{\mathcal{N}}^{\bar{m}(v)}\left(v\right)\right\}\right). \tag{5}$ We send only the aggregated results from each remote client $\bar{m}(v)$ to the server, which can help preserve data privacy and reduce communication overhead.

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

$$\boldsymbol{r}_{\mathcal{N}(v)}^{(l)} = \operatorname{AGG}\left(\left\{\boldsymbol{h}_{\mathcal{N}(v)}^{(l),\bar{m}(v)} \mid \bar{m}(v) \in \bar{\mathcal{M}}(v)\right\}\right). \tag{6}$$

This approach not only helps maintain data privacy but also reduces communication costs by minimizing the amount of data transmitted between clients and the server.

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

$$\boldsymbol{h}_{\mathcal{N}(v)}^{(l)} = \operatorname{AGG}\left(\underbrace{\left\{\boldsymbol{h}_{u}^{(l-1),m(v)} \mid u \in \widetilde{\mathcal{N}}^{m(v)}(v)\right\}}_{\text{local}} \cup \underbrace{\left\{\boldsymbol{r}_{\mathcal{N}(v)}^{(l)}\right\}}_{\text{remote}}\right). \tag{7}$$

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

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

$$\boldsymbol{h}_{v}^{(l),m(v)} = \sigma \left(\boldsymbol{W}_{t}^{(l),m(v)} \cdot \text{COMB} \left(\boldsymbol{h}_{v}^{(l-1),m(v)}, \boldsymbol{h}_{\mathcal{N}(v)}^{(l)} \right) \right). \tag{8}$$

Using the embeddings of the training nodes in the mini-batch and the model parameters, the local stochastic gradients $\nabla F^m\left(m{ heta}_t^m\right)$ are computed and used in the update of the global model parameters shown as $\theta_{t+1} = \theta_t - \alpha \frac{1}{|\mathcal{M}|} \sum_{m \in \mathcal{M}} \nabla \widetilde{F}^m(\theta_t^m)$, where α is the learning rate.

5 THEORETICAL PERFORMANCE ANALYSIS

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

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)

325

326

327

328

330 331 332

333

334 335

336 337 338

339

340

341

342

343 344 345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360 361 362

363

364

366

367

368

369 370

371

372 373 374

375

376

377

For a graph \mathcal{G} , the structure can be represented by its adjacency matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$, where $\mathbf{A}_{vu} = 1$ if $(v, u) \in \mathcal{E}$, otherwise $A_{vu} = 0$. The propagation matrix can be computed as $P = D^{-1/2} \hat{A} D^{-1/2}$, where $\hat{A} = A + I$, and $D \in \mathbb{R}^{N \times N}$ corresponds to the degree matrix and $D_{vv} = \sum_{u} \hat{A}_{vu}$. For subgraph \mathcal{G}^m located on client m, the adjacency matrix A^m can be denoted as $A^m = A_{local}^m + A_{local}^m$ subgraph ${\bf y}^m$ located on client m, the adjacency matrix ${\bf A}^m$ can be defined as ${\bf A}^m - {\bf A}_{local}^m$. A can be defined as ${\bf A}^m - {\bf A}_{local}^m$ corresponds to their cross-client neighbors located on the remote clients other than m. Then, the propagation matrix can be calculated as ${\bf P}^m = {\bf D}_m^{-1/2} \left({\bf A}^m + {\bf I}^m \right) {\bf D}_m^{-1/2}$, and can be represented as ${\bf P}^m = {\bf P}_{local}^m + {\bf P}_{local}^m$, where ${\bf P}_{local}^m = {\bf D}_m^{-1/2} \left({\bf A}_{local}^m + {\bf I}^m \right) {\bf D}_m^{-1/2}$ and ${\bf P}_{remote}^m = {\bf D}_m^{-1/2} \left({\bf A}_{remote}^m \right) {\bf D}_m^{-1/2}$. Given GCN as the GNN architecture, for client m training using only the local graph data, Eq. (3) and (4) are equivalent to $\widetilde{\boldsymbol{H}}_{t}^{(l),m} = \sigma\left(\widetilde{\boldsymbol{P}}_{local}^{(l),m}\widetilde{\boldsymbol{H}}_{local}^{(l-1),m}\boldsymbol{W}_{t}^{(l),m}\right)$. For client m training based on both the local graph data and the cross-client neighbors, Eq. (5)–(8) are equivalent to $\widetilde{\boldsymbol{H}}_t^{(l),m} = \sigma\left(\left(\widetilde{\boldsymbol{P}}_{local}^{(l),m}\widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m}\widetilde{\boldsymbol{H}}_{remote}^{(l-1),m}\right)\boldsymbol{W}_t^{(l),m}\right)$. 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 $\boldsymbol{h}^{(L)}$, i.e., $\|\ell^m(\boldsymbol{h}_1^{(L)},y) - \ell^m(\boldsymbol{h}_2^{(L)},y)\|_2 \leq C_l \|\boldsymbol{h}_1^{(L)} - \boldsymbol{h}_2^{(L)}\|_2$ and $\|\nabla \ell^m(\boldsymbol{h}_1^{(L)}, y) - \nabla \ell^m(\boldsymbol{h}_2^{(L)}, y)\|_2 \le L_l \|\boldsymbol{h}_1^{(L)} - \boldsymbol{h}_2^{(L)}\|_2.$

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

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 $||X||_F \leq B_X$. 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 analysis of Swift-FedGNN. Some existing works used strong assumptions to deal with these biased stochastic gradients in their analysis, e.g., Chen et al. (2018) adopts the unbiased stochastic gradient assumption, and Chen & Luss (2018) uses 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 $\left\| \nabla F_{local}^{m}\left(oldsymbol{ heta}^{m}\right) - \nabla \tilde{F}_{local}^{m}\left(oldsymbol{ heta}^{m}\right) \right\|_{F} \leq LB_{\Delta G}^{l}$ and $\left\|\nabla F_{full}^m\left(\boldsymbol{\theta}^m\right) - \nabla \tilde{F}_{full}^m\left(\boldsymbol{\theta}^m\right)\right\|_F \leq LB_{\Delta G}^f, \text{ where } \nabla F_{local}^m\left(\boldsymbol{\theta}^m\right) \text{ and } \nabla \widetilde{F}_{local}^m\left(\boldsymbol{\theta}^m\right) \text{ correspond} \text{ to the full and stochastic gradients computed with only local graph data, respectively. } \nabla F_{full}^m\left(\boldsymbol{\theta}^m\right)$ and $\nabla \widetilde{F}_{full}^m(\pmb{\theta}^m)$ include both local graph data and cross-client neighbors of the training nodes. $B_{\Delta G}^{l}$ and $B_{\Delta G}^{f}$ are defined in Eq. (12) and (13) in Appendix E.

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 crossclient 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 $(\nabla F_{full}^m(\boldsymbol{\theta}^m))$ and the full gradient computed with only the local graph data $(\nabla F_{local}^{m}(\boldsymbol{\theta}^{m}))$ is upper-bounded as $\left\| \nabla F_{full}^m \left(\boldsymbol{\theta}^m \right) - \nabla F_{local}^m \left(\boldsymbol{\theta}^m \right) \right\|_F \le L B_{\Delta G}^r$, where $B_{\Delta G}^r$ is defined in Eq. (14) in Appendix E.

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.

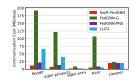


Figure 4: Convergence performance in terms of validation accuracy of different algorithms.

Figure 5: Average communication cost per step.

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 E.2. The output of Swift-FedGNN solving an L-layer GNN satisfies:

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

The detailed proof of Theorem 5.6 can be found in Appendix E. We can see from Theorem 5.6 that the convergence rate of Swift-FedGNN is $\mathcal{O}\left(T^{-1/2}\right)$ 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. Due to space limitations, additional experimental details and results are provided in Appendix D.

1) Experiment Settings: We train a representative GNN model, Graph-SAGE Hamilton et al. (2017), in the FL settings on five real-world node classification datasets: 1) ogbn-products Hu et al. (2020); 2) Reddit Hamilton

1) **Experiment Settings:** Table 1: Total communication cost (GB) when achieving a target We train a representative validation accuracy for each dataset.

	OGBN- PRODUCTS	REDDIT	OGBN- ARXIV	FLICKR	CITESEER
SWIFT-FEDGNN	0.66	5.89	0.95	1.07	35.32
FEDGNN-G	8.40	70.72	7.91	17.22	43.43
LLCG	4.47	23.89	1.46	1.30	25.80
FEDGNN-PNS	0.97	8.96	1.41	1.27	36.11

et al. (2017); 3) ogbn-arxiv Hu et al. (2020); 4) flickr Zeng et al. (2020); and 5) citeseer Giles et al. (1998). The key statistics of the datasets are summarized in Table 8 in Appendix D. Note that ogbn-products dataset is the *largest* dataset one can find in the FGL literature, while the Reddit dataset is known for its *density*. In our FL simulations, we use 20 clients for the experiments with ogbn-products dataset and 10 clients for the experiments with the other datasets. All graphs are partitioned with METIS algorithm Karypis & Kumar (1998). In addition, we evaluate Swift-FedGNN on randomly partitioned graph data and on another widely used GNN model, GIN Xu et al. (2019). The corresponding results are provided in Appendix D.

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) LLCG Ramezani et al. (2022): A DGL 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 FGL 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: A naive FGL algorithm where cross-client training is performed on each client in every iteration.

- 3) Convergence Performance Comparisons: In Figures 4a and 4b, we can see that for both the ogbnproducts dataset and the Reddit dataset, Swift-FedGNN achieves substantially faster convergence than all baseline algorithms, which verifies the effectiveness of Swift-FedGNN in handling large or dense graphs. In addition, despite less frequent cross-client training, 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. Although LLCG performs periodic cross-client training on the server, it requires 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, LLCG takes over 5000 ms to perform cross-client training on the server, whereas Swift-FedGNN completes crossclient training within 200 ms due to neighbor sampling. FedGNN-PNS employs a dynamic crossclient 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. As shown in Figure 4c, on the smaller ogbn-arxiv dataset the benefit of Swift-FedGNN is less pronounced. The dataset's limited size and sparsity reduce both neighbor sampling and communication overhead for all methods, narrowing the performance gap. This is also reflected in the following communication cost analysis. Nevertheless, Swift-FedGNN still delivers the best overall performance.
- 4) Communication Cost Analysis: Figure 5 shows the average communication cost per step for Swift-FedGNN and the baseline algorithms across five datasets, demonstrating that Swift-FedGNN consistently incurs the lowest communication cost on all of them. Specifically, our algorithm Swift-FedGNN incurs a communication cost that is $7\times$ to $21\times$ lower than that of FedGNN-G on four out of the five datasets (Reddit, ogbn-products, ogbn-arxiv, and flickr). On the smallest graph, citeseer, the gap narrows because the size of the cross-client neighbor information becomes negligible compared with the model size, yet Swift-FedGNN still maintains the lowest communication cost. For the largest dataset ogbn-products and the most dense dataset Reddit, Swift-FedGNN achieves communication costs that are $2\times$ and $5\times$ lower compared to FedGNN-PNS and LLCG, respectively. On the small datasets, ogbn-arxiv and flickr, the communication cost advantage of Swift-FedGNN remains evident, though closer to approximately $1\times$ lower than FedGNN-PNS and LLCG. These findings validate the superior communication efficiency of our proposed Swift-FedGNN algorithm.

Table 1 reports the total communication cost required to reach the same target validation accuracy on each dataset. The results demonstrate that our proposed Swift-FedGNN algorithm consistently incurs the lowest communication cost across all datasets except Citeseer. For example, to reach a target accuracy of 87% on the ogbn-products dataset, Swift-FedGNN achieves at least a 31.9% reduction in total communication cost compared to all baselines. Similarly, to reach a target accuracy of 55% on the smaller and sparser ogbn-arxiv dataset, Swift-FedGNN still delivers at least 32.2% communication savings, highlighting its robustness and efficiency across diverse graph structures.

7 CONCLUSION

In this paper, we proposed the Swift-FedGNN algorithm, which is a mini-batch-based and sampling-based federated graph learning 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 $\mathcal{O}\left(T^{-1/2}\right)$, matching the SOTA rate of sampling-based GNN methods, even in more challenging federated settings. We conducted extensive numerical experiments on real-world graph datasets and verified the effectiveness of Swift-FedGNN.

ETHICS STATEMENT

We confirm that the ICLR Code of Ethics has been thoroughly reviewed and that this work fully adheres to it. The study does not involve human subjects, sensitive data, or any foreseeable risks, and it raises no ethical, legal, or conflict-of-interest concerns.

REPRODUCIBILITY STATEMENT

We confirm the reproducibility of this work. Specifically, for the theoretical results, we state the assumptions in Section 5 and provide detailed proofs in Appendix E. For the experimental results, we submit the source code as a supplementary material and describe the implementation details in Appendix D.

REFERENCES

- Davide Bacciu, Federico Errica, and Alessio Micheli. Contextual graph markov model: A deep and generative approach to graph processing. In *International conference on machine learning*, pp. 294–303. PMLR, 2018.
- Zhenkun Cai, Xiao Yan, Yidi Wu, Kaihao Ma, James Cheng, and Fan Yu. DGCL: An Efficient Communication Library for Distributed GNN Training. In *Proc. of EuroSys*, pp. 130–144, 2021.
- Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with variance reduction. In *International Conference on Machine Learning*, pp. 942–950. PMLR, 2018.
- Jie Chen and Ronny Luss. Stochastic gradient descent with biased but consistent gradient estimators. arXiv preprint arXiv:1807.11880, 2018.
- Weilin Cong, Morteza Ramezani, and Mehrdad Mahdavi. On the importance of sampling in training gcns: Tighter analysis and variance reduction. *arXiv preprint arXiv:2103.02696*, 2021.
- Songgaojun Deng, Huzefa Rangwala, and Yue Ning. Learning dynamic context graphs for predicting social events. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1007–1016, 2019.
- Kien Do, Truyen Tran, and Svetha Venkatesh. Graph transformation policy network for chemical reaction prediction. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 750–760, 2019.
- Bingqian Du and Chuan Wu. Federated graph learning with periodic neighbour sampling. In 2022 IEEE/ACM 30th International Symposium on Quality of Service (IWQoS), pp. 1–10. IEEE, 2022.
- Bingqian Du, Jun Liu, Ziyue Luo, Chuan Wu, Qiankun Zhang, and Hai Jin. Expediting Distributed GNN Training with Feature-only Partition and Optimized Communication Planning. In *Proc. of IEEE INFOCOM*, 2024.
- Matthias Fey and Jan Eric Lenssen. Fast Graph Representation Learning with PyTorch Geometric. In *Proc. of ICLR*, 2019.
- Swapnil Gandhi and Anand Padmanabha Iyer. P3: Distributed Deep Graph Learning at Scale. In *Proc. of OSDI*, pp. 551–568, 2021.
- Vikas Garg, Stefanie Jegelka, and Tommi Jaakkola. Generalization and representational limits of graph neural networks. In *International Conference on Machine Learning*, pp. 3419–3430. PMLR, 2020.
- C Lee Giles, Kurt D Bollacker, and Steve Lawrence. Citeseer: An automatic citation indexing system. In *Proceedings of the third ACM conference on Digital libraries*, pp. 89–98, 1998.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.

- Chaoyang He, Keshav Balasubramanian, Emir Ceyani, Carl Yang, Han Xie, Lichao Sun, Lifang He, Liangwei Yang, Philip S Yu, Yu Rong, et al. Fedgraphnn: A federated learning system and benchmark for graph neural networks. *arXiv preprint arXiv:2104.07145*, 2021a.
 - 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.
 - Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems*, 33:22118–22133, 2020.
 - Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*, pp. 5132–5143. PMLR, 2020.
 - George Karypis and Vipin Kumar. A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. *SIAM Journal on Scientific Computing*, 20(1):359–392, 1998.
 - Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
 - Jure Leskovec, Lada A Adamic, and Bernardo A Huberman. The Dynamics of Viral Marketing. *ACM Transactions on the Web (TWEB)*, 1(1):5–es, 2007.
 - Renjie Liao, Raquel Urtasun, and Richard Zemel. A pac-bayesian approach to generalization bounds for graph neural networks. *arXiv preprint arXiv:2012.07690*, 2020.
 - Tianfeng Liu, Yangrui Chen, Dan Li, Chuan Wu, Yibo Zhu, Jun He, Yanghua Peng, Hongzheng Chen, Hongzhi Chen, and Chuanxiong Guo. BGL: GPU-Efficient GNN Training by Optimizing Graph Data I/O and Preprocessing. In *Proc. of NSDI*, pp. 103–118, 2023.
 - Ziyue Luo, Yixin Bao, and Chuan Wu. Optimizing Task Placement and Online Scheduling for Distributed GNN Training Acceleration. In *Proc. of IEEE INFOCOM*, 2022.
 - Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pp. 1273–1282. PMLR, 2017.
 - Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Proc. of NeurIPS*, 2019.
 - Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. Deepinf: Modeling influence locality in large social networks. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'18)*, 2018.
 - Morteza Ramezani, Weilin Cong, Mehrdad Mahdavi, Mahmut Kandemir, and Anand Sivasubramaniam. Learn locally, correct globally: A distributed algorithm for training graph neural networks. In *International Conference on Learning Representations*, 2022.
 - John Thorpe, Yifan Qiao, Jonathan Eyolfson, Shen Teng, Guanzhou Hu, Zhihao Jia, Jinliang Wei, Keval Vora, Ravi Netravali, Miryung Kim, et al. Dorylus: Affordable, Scalable, and Accurate GNN Training with Distributed CPU Servers and Serverless Threads. In *Proc. of USENIX OSDI*, 2021.
 - Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph Attention Networks. *arXiv preprint arXiv:1710.10903*, 2017.
 - Cheng Wan, Youjie Li, Cameron R. Wolfe, Anastasios Kyrillidis, Nam Sung Kim, and Yingyan Lin. PipeGCN: Efficient full-graph training of graph convolutional networks with pipelined feature communication. In *International Conference on Learning Representations*, 2022.

- Binghui Wang, Ang Li, Meng Pang, Hai Li, and Yiran Chen. Graphfl: A federated learning framework for semi-supervised node classification on graphs. In 2022 IEEE International Conference on Data Mining (ICDM), pp. 498–507. IEEE, 2022a.
 - Hongwei Wang, Fuzheng Zhang, Mengdi Zhang, Jure Leskovec, Miao Zhao, Wenjie Li, and Zhongyuan Wang. Knowledge-aware graph neural networks with label smoothness regularization for recommender systems. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 968–977, 2019a.
 - Kuansan Wang, Zhihong Shen, Chiyuan Huang, Chieh-Han Wu, Yuxiao Dong, and Anshul Kanakia. Microsoft Academic Graph: When Experts Are Not Enough. *Quantitative Science Studies*, 1(1): 396–413, 2020.
 - Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang, Chao Ma, et al. Deep Graph Library: Towards Efficient and Scalable Deep Learning on Graphs. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019b.
 - Yuyang Wang, Jianren Wang, Zhonglin Cao, and Amir Barati Farimani. Molecular contrastive learning of representations via graph neural networks. *Nature Machine Intelligence*, 4(3):279–287, 2022b.
 - Chuhan Wu, Fangzhao Wu, Yang Cao, Yongfeng Huang, and Xing Xie. Fedgnn: Federated graph neural network for privacy-preserving recommendation. *arXiv preprint arXiv:2102.04925*, 2021.
 - Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In *International conference on machine learning*, pp. 6861–6871. PMLR, 2019.
 - Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
 - Haibo Yang, Minghong Fang, and Jia Liu. Achieving linear speedup with partial worker participation in non-IID federated learning. In *International Conference on Learning Representations*, 2021.
 - Yuhang Yao, Weizhao Jin, Srivatsan Ravi, and Carlee Joe-Wong. Fedgcn: Convergence-communication tradeoffs in federated training of graph convolutional networks. *Advances in Neural Information Processing Systems*, 36, 2023a.
 - Yuhang Yao, Mohammad Mahdi Kamani, Zhongwei Cheng, Lin Chen, Carlee Joe-Wong, and Tianqiang Liu. Fedrule: Federated rule recommendation system with graph neural networks. In *Proceedings of the 8th ACM/IEEE Conference on Internet of Things Design and Implementation*, pp. 197–208, 2023b.
 - Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 974–983, 2018.
 - Binhang Yuan, Yongjun He, Jared Davis, Tianyi Zhang, Tri Dao, Beidi Chen, Percy S Liang, Christopher Re, and Ce Zhang. Decentralized training of foundation models in heterogeneous environments. volume 35, pp. 25464–25477, 2022.
 - Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graph-SAINT: Graph Sampling Based Inductive Learning Method. In *Proc. of ICLR*, 2020.
 - Ke Zhang, Carl Yang, Xiaoxiao Li, Lichao Sun, and Siu Ming Yiu. Subgraph federated learning with missing neighbor generation. *Advances in Neural Information Processing Systems*, 34:6671–6682, 2021.
 - Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in neural information processing systems*, 31, 2018.

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

- Zhe Zhang, Ziyue Luo, and Chuan Wu. Two-Level Graph Caching for Expediting Distributed GNN Training. In *Proc. of INFOCOM*, pp. 1–10. IEEE, 2023.
- Kun Zhao, Wencong Xiao, Baole Ai, Wenting Shen, Xiaolin Zhang, Yong Li, and Wei Lin. AliGraph: An Industrial Graph Neural Network Platform. In *Proc. of SOSP Workshop on AI Systems*, 2019.
- Da Zheng, Chao Ma, Minjie Wang, Jinjing Zhou, Qidong Su, Xiang Song, Quan Gan, Zheng Zhang, and George Karypis. DistDGL: Distributed Graph Neural Network Training for Billion-Scale Graphs. In *IEEE/ACM Workshop on Irregular Applications: Architectures and Algorithms*, 2020.

A THE USE OF LARGE LANGUAGE MODELS (LLMS)

LLMs were used exclusively for grammar correction and language polishing during the writing process. They did not contribute to research ideation or any substantive aspects of the work.

B LIST OF NOTATIONS

702

703 704 705

706

```
710
711
712
             \mathcal{G}(\mathcal{V},\mathcal{E})
                                          Graph
713
             \mathcal{V}
                                          Set of nodes
714
             \mathcal{E}
                                          Set of edges
715
             N = |\mathcal{V}|
                                          Number of nodes
716
717
             \mathcal{M}
                                          Set of clients
718
             M = |\mathcal{M}|
                                          Number of clients
719
             \mathcal{G}^m(\mathcal{V}^m, \mathcal{E}^m)
                                          Subgraph at client m
720
             \mathcal{V}^m
721
                                          Set of nodes at client m
             \mathcal{E}^m
722
                                          Set of edges at client m
723
             \boldsymbol{x}_v^m \in \mathbb{R}^d
                                          Feature vector of node v at client m
724
             y_v^m
                                          Label of node v at client m
725
726
             \ell^m
                                          Loss function (e.g., cross-entropy loss) at client m
727
             \mathcal{V}_B^m
                                          Mini-batch of training nodes
728
             \boldsymbol{\theta} = \left\{ \boldsymbol{W}^{(l)} \right\}_{l=1}^{L}
                                          Set of trainable model parameters
729
730
             m(v)
                                          Local client of node v
731
             \bar{m}(v)
                                          Remote client of node v
732
             \bar{\mathcal{M}}(v)
                                          Set of the remote clients that host the neighbors of node v
733
             \mathcal{N}^{m(v)}(v)
                                          Set of the neighbors of node v located on local client m(v)
734
             \mathcal{N}^{\bar{m}(v)}\left(v\right)
735
                                          Set of the neighbors of node v located on remote client \bar{m}(v)
736
             m{h}_v^{(l),m(v)}
                                          Embedding of node v located on client m(v)
737
             oldsymbol{h}_{\mathcal{N}(v)}^{(l)}
                                          Aggregated embedding from node v's neighbors
738
             W^{(l)}
739
                                          Weight matrix at l-th layer
740
             \sigma\left(\cdot\right)
                                          Activation function (e.g., ReLU)
741
             AGG(\cdot)
                                          Aggregation function (e.g., mean)
742
             COMB(\cdot)
                                          Combination function (e.g., concatenation)
743
744
                                          Mini-batch of training nodes at client m
             \mathcal{S} = \left\{ \mathcal{S}^{(l)} \right\}_{l=0}^{L-1}
745
                                          Subset of L-hop neighbors for the training nodes in \mathcal{B}_v^m
746
             \widetilde{\mathcal{N}}^m(v)
                                          Set of the sampled neighbors located on client m for node v
747
748
             \mathcal{K}
                                          Set of sampled clients for cross-client training
749
             K = |\mathcal{K}|
                                          Number of sampled clients for cross-client training
750
             \nabla \widetilde{F}^m \left( \boldsymbol{\theta}_t^m \right)
                                          Stochastic gradient
751
                                          Learning rate
752
753
              \boldsymbol{A} \in \mathbb{R}^{N \times N}
                                          Adjacency matrix of graph \mathcal{G}
754
             \boldsymbol{P}
                                          Propagation matrix
755
             D
                                          Degree matrix
```

756	$oldsymbol{A}^m$	Adjacency matrix of subgraph \mathcal{G}^m
757 758	$m{A}_{local}^m$	Adjacency matrix corresponds to the nodes located on client m
759	$m{A}_{remote}^m$	Adjacency matrix corresponds to the cross-client neighbors located on the
760		remote clients other than m
761	\boldsymbol{D}^m	Degree matrix of client m
762	$oldsymbol{P}^m$	Propagation matrix of client m
763	\mathbf{D}^m	
764	$m{P}_{local}^{m}$	Propagation matrix corresponds to the nodes located on client m
765	$m{P}_{remote}^m$	Propagation matrix corresponds to the cross-client neighbors located on the
766		remote clients other than m

C SINGLE-MACHINE GRAPH NEURAL NETWORKS TRAINING

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 $\mathbf{x}_v \in \mathbb{R}^d$, where d is the dimension of the feature vector. Each node $v \in \mathcal{V}_{train}$ has a corresponding label y_v , where $\mathcal{V}_{train} \subseteq \mathcal{V}$.

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:

$$\boldsymbol{h}_{\mathcal{N}\left(v\right)}^{\left(l\right)} = \mathrm{AGG}\left(\left\{\boldsymbol{h}_{u}^{\left(l-1\right)} \mid u \in \mathcal{N}\left(v\right)\right\}\right), \quad \boldsymbol{h}_{v}^{\left(l\right)} = \sigma\left(\boldsymbol{W}^{\left(l\right)} \cdot \mathrm{COMB}\left(\boldsymbol{h}_{v}^{\left(l-1\right)}, \boldsymbol{h}_{\mathcal{N}\left(v\right)}^{\left(l\right)}\right)\right),$$

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

D ADDITIONAL EXPERIMENTAL DETAILS AND RESULTS

D.1 ADDITIONAL EXPERIMENTAL RESULTS

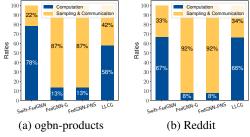
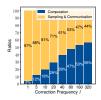


Figure 6: Ratio of computation time to sampling and communication time for different algorithms.

Communication and sample costs analysis: Figure 6 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 LLCG, it converges much faster and achieves higher validation accuracy than LLCG.

Hyperparameter sensitivity analysis: We explore the impact of the important hyperparameters in Swift-FedGNN. Figure 7a shows that when the correction frequency *I* increases, the computation-(sampling & communication) ratio increases. Figure 7b and 7c indicate that as the number of









(a) correction frequencies

(b) # of cross-client training clients (K)

(c) # of sampled neighbors

(d) # of clients (50% for cross-client training)

Figure 7: Ratio of computation time to sampling and communication time for Swift-FedGNN on ogbn-products dataset.

cross-client training clients K, and the number of sampled neighbors increase, the computation-(sampling & communication) ratio decreases. Figure 7d 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.

Table 2: Communication overhead per iteration when communication occurs.

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

Communication overhead when communication occurs: 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.

Table 3: Validation accuracy (%) of different algorithms using the GraphSAGE model.

	OGBN-PRODUCTS	REDDIT	OGBN-ARXIV	FLICKR	CITESEER
SWIFT-FEDGNN	88.88	95.47	57.17	50.19	66.00
LLCG	87.66	95.27	56.78	50.12	68.40
FEDGNN-PNS	87.89	95.46	55.86	51.47	66.27
FEDSAGE	88.15	95.30	56.55	49.75	64.39
FEDGNN-G	88.71	95.96	56.78	51.57	66.08

Validation accuracy comparisons: Table 3 shows the validation accuracy of different algorithms. To assess the impact of cross-client neighbors, we include an additional baseline FedSage Zhang et al. (2021), an FGL algorithm that entirely ignores cross-client neighbors and performs purely local training in all iterations. The results demonstrate that despite incurring lower sampling and communication overhead, our Swift-FedGNN achieves validation accuracy comparable to that of the baseline algorithms. Moreover, compared to FedSage, which completely ignores cross-client neighbors, Swift-FedGNN achieves a higher validation accuracy, highlighting the importance of incorporating cross-client neighbor information. By minimizing sampling and communication overhead, Swift-FedGNN offers the highest efficiency in practical implementation.

It is worth noting that ogbn-arxiv, flickr, and citeseer are small datasets (Table 8), where graph partitioning leads to greater information loss. As a result, baselines that frequently exchange graph data can achieve slightly higher accuracy. However, these small datasets do not require federated graph learning in practice. Federated graph learning is primarily motivated by large-scale datasets like ogbn-products, where our method achieves the best performance.

Evaluations using the GIN model: To assess the adaptivity of Swift-FedGNN to different GNN models, we conduct experiments using the GIN Xu et al. (2019) model across multiple datasets.

SWIFT-FEDGNN

LLCG

FEDGNN-PNS

FEDGNN-G

864 865

Table 4: Validation accuracy (%) of different algorithms using the GIN model.

OGBN-ARXIV

56.69

57.32

56.54

57.01

CITESEER

47.34

46.60

47.99

50.76

OGBN-PRODUCTS

81.93

80.72

78.70

83.76

866

867 868

871 872

873 874 875

876 877

878 879

882 883 884

887 888

885

889 890 891

892 893 894

895

897

903

909

910

911

912

913

914

915

916

917

870

Table 5: Performance comparison using the GIN model when achieving a target validation accuracy for each dataset. OCDN DDODUCTS OCDN-ADVIV

	OGBN-PRODUCTS		OGBN-	ARXIV
	TOTAL COMM.	WALL-CLOCK	TOTAL COMM.	Wall-Clock
	Cost (GB)	TIME (S)	Cost (GB)	TIME (S)
SWIFT-FEDGNN	0.74	65.18	2.29	75.80
LLCG	3.74	223.77	3.81	103.98
FEDGNN-PNS	5.75	113.46	3.62	131.92
FEDGNN-G	38.12	767.53	19.36	575.26

Table 4 shows that, similar to the results with the GraphSAGE model, Swift-FedGNN achieves comparable validation accuracy to the baseline algorithms while significantly reducing sampling and communication overhead.

Table 5 reports the total communication cost and wall-clock time on the ogbn-products and ogbn-arxiv datasets when reaching a target validation accuracy of 80% and 56%, respectively. In both cases, Swift-FedGNN achieves the lowest wall-clock time to reach the target accuracy. Moreover, it reduces the total communication cost by at least 80% on ogbn-products and 37% on ogbn-arxiv compared to all baselines. These results demonstrate the effectiveness of Swift-FedGNN in significantly reducing communication overhead when using the GIN model.

Table 6: Total communication cost using randomly partitioned ogbn-products dataset when achieving a target validation accuracy of 89.5%.

	SWIFT-FEDGNN	FedGNN-G	LLCG	FEDGNN-PNS
TOTAL COMMUNICATION COST (GB)	1.44	15.03	6.26	2.60

Evaluations using randomly partitioned ogbn-products dataset: To evaluate the robustness of Swift-FedGNN under less structured scenarios, we conduct additional experiments using random partitioning instead of METIS on the ogbn-products dataset. Table 6 shows the total communication cost when achieving a target validation accuracy of 89.5%. Swift-FedGNN reduces total communication cost by at least 45% compared to all baselines. Table 7 reports the validation accuracy, demonstrating that Swift-FedGNN achieves the highest accuracy among methods that do not rely on full graph training. These results confirm that Swift-FedGNN maintains both communication efficiency and competitive performance even when the data is randomly partitioned, validating its applicability beyond well-partitioned settings.

D.2 ADDITIONAL EXPERIMENTAL DETAILS

Dataset. Table 8 summarizes the key statistics of the datasets used in our experiments, including: 1) ogbn-products Hu et al. (2020), which is an Amazon product co-purchasing graph derived from Leskovec et al. (2007); 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; 3) ogbn-arxiv Hu et al. (2020), which is a citation network between arXiv papers in the field of computer science, where nodes represent papers and directed edges indicate citation links; 4) flickr Zeng et al. (2020), which is an image network where each node represents an image and edges connect images that share common properties such as tags or visual similarity; and 5) citeseer Giles et al. (1998), which is a citation graph of research papers, where each node denotes a document and edges represent citation relationships between them.

Implementation and testbed. We implement Swift-FedGNN using Python on DGL 2.0.0 Wang et al. (2019b) and PyTorch 2.2.1 Paszke et al. (2019). Our implementation includes a custom GPU-based sampler built on top of DGL's native sampler, which is designed to sequentially sample local and

Table 7: Comparison of validation accuracy (%) using randomly partitioned ogbn-products dataset.

	SWIFT-FEDGNN	FEDGNN-G	LLCG	FEDGNN-PNS
VALIDATION ACCURACY (%)	89.94	91.23	89.92	89.91

Table 8: Benchmark datasets and key parameters.

DATASET	# of Nodes	# of Edges	Edges per Node
OGBN-PRODUCTS	2.4 M	61.9 M	25.8
REDDIT	233 K	114.6 M	491.4
OGBN-ARXIV	169 K	1.2 M	7.1
FLICKR	89 K	900 K	10.1
CITESEER	3327	9228	2.8

remote neighbors for each client at every layer. Additionally, we customized the GraphSAGE layer and GIN layer 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 and a two-layer GIN 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-out values are set to [20, 15] for the ogbn-products dataset, and [15, 10] for all other datasets.

Hyperparameters. The training mini-batch size is set at 256. For optimization, we use the Adam optimizer with a weight decay of 5×10^{-4} . We use a learning rate of 0.01 for the ogbn-products dataset, 0.0001 for the flickr dataset, 0.00001 for the citeseer dataset, and 0.001 for both the ogbn-arxiv and Reddit datasets. In Swift-FedGNN, we set K=10 for the ogbn-products dataset and K=5 for all other datasets. We choose I=5 for the citeseer dataset and I=10 for the remaining datasets.

E Proof of Theorem 5.6

E.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:

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

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

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

in which $\widetilde{\boldsymbol{Z}}_{local}^{(l),m} = \widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \boldsymbol{W}^{(l),m}$, $\widetilde{\boldsymbol{D}}_{local}^{(L),m} = \partial \ell^m \left(\widetilde{\boldsymbol{H}}_{local}^{(L),m}, \boldsymbol{Y}_{local}^m \right) / \partial \widetilde{\boldsymbol{H}}_{local}^{(L),m}$, 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:

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

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

$$\begin{split} \widetilde{\boldsymbol{G}}_{full}^{(l),m} &= \left[\nabla_{W} \widetilde{\boldsymbol{f}}^{(l),m} \left(\widetilde{\boldsymbol{D}}_{full}^{(l),m}, \widetilde{\boldsymbol{H}}_{full}^{(l-1),m}, \boldsymbol{W}^{(l),m} \right) \right. \\ &\triangleq \left[\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m} \widetilde{\boldsymbol{H}}_{remote}^{(l-1),m} \right]^{\top} \widetilde{\boldsymbol{D}}_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{\boldsymbol{Z}}_{full}^{(l),m} \right) \right], \\ \widetilde{\boldsymbol{D}}_{full}^{(l),m} &= \left[\nabla_{H} \widetilde{\boldsymbol{f}}^{(l+1),m} \left(\widetilde{\boldsymbol{D}}_{full}^{(l+1),m}, \widetilde{\boldsymbol{H}}_{full}^{(l),m}, \boldsymbol{W}^{(l+1),m} \right) \right. \\ &\triangleq \left[\widetilde{\boldsymbol{P}}_{local}^{(l+1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l+1),m} \right]^{\top} \widetilde{\boldsymbol{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\widetilde{\boldsymbol{Z}}_{full}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right], \end{split}$$
which $\widetilde{\boldsymbol{Z}}^{(l),m} = \left(\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m} \widetilde{\boldsymbol{H}}^{(l-1),m} \right) \boldsymbol{W}^{(l),m} \quad \text{and} \quad \widetilde{\boldsymbol{D}}^{(L),m} \end{split}$

$$\begin{array}{lll} \text{in which} & \widetilde{\boldsymbol{Z}}_{full}^{(l),m} & = & \left(\widetilde{\boldsymbol{P}}_{local}^{(l),m}\widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m}\widetilde{\boldsymbol{H}}_{remote}^{(l-1),m}\right)\boldsymbol{W}^{(l),m}, & \text{and} & \widetilde{\boldsymbol{D}}_{full}^{(L),m} & = \\ \partial \ell^m \left(\widetilde{\boldsymbol{H}}_{full}^{(L),m}, \boldsymbol{Y}_{full}^m\right) / \partial \widetilde{\boldsymbol{H}}_{full}^{(L),m}. & \end{array}$$

E.2 USEFUL PROPOSITIONS AND LEMMAS

Proposition E.1. *Under Assumption 5.3, the inequalities in Table 9 and Table 10 are hold.*

Table 9: Upper-bound for the norms of the propagation matrix and the node feature matrix.

	PROPAGATION MATRIX	Node Feature Matrix
FULL GRAPH	$\ \boldsymbol{P}_{full}\ _F \leq B_P$	$\ \boldsymbol{X}_{full}\ _F \leq B_X$
LOCAL GRAPH	$\ \boldsymbol{P}_{local}\ _F \leq B_P^l \leq B_P$	$\ \boldsymbol{X}_{local}\ _F \leq B_X^l \leq B_X$
CROSS-CLIENT NEIGHBORS	$\ \boldsymbol{P}_{remote}\ _F \leq B_P^r \leq B_P$	$\ \boldsymbol{X}_{remote}\ _F \leq B_X^r \leq B_X$

Table 10: Relationships for the norms of the propagation matrix and the node feature matrix before and after sampling.

	PROPAGATION MATRIX	Node Feature Matrix
FULL GRAPH	$\left\ \widetilde{m{P}}_{full} - m{P}_{full} ight\ _F \leq B_{\Delta P}^f$	$\left\ \widetilde{m{X}}_{full}-m{X}_{full} ight\ _F \leq B_{\Delta X}^f$
LOCAL GRAPH	$\left\ \widetilde{P}_{local} - P_{local} \right\ _F \le B_{\Delta P}^l$	$\left\ \widetilde{\boldsymbol{X}}_{local} - \boldsymbol{X}_{local}\right\ _F \leq B_{\Delta X}^l$
CROSS-CLIENT NEIGHBORS	$\left\ \widetilde{\widetilde{P}}_{remote} - P_{remote} \right\ _F \le B_{\Delta P}^r$	$\left\ \widetilde{\widetilde{\boldsymbol{X}}}_{remote} - \boldsymbol{X}_{remote}\right\ _F \leq B_{\Delta X}^r$

Lemma E.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 E.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.,

$$\begin{split} & \left\| \boldsymbol{H}_{local}^{(l),m} \right\|_{F}, \left\| \widetilde{\boldsymbol{H}}_{local}^{(l),m} \right\|_{F} \leq B_{H}^{l}, & \left\| \boldsymbol{H}_{full}^{(l),m} \right\|_{F}, \left\| \widetilde{\boldsymbol{H}}_{full}^{(l),m} \right\|_{F} \leq B_{H}^{f}, \\ & \left\| \boldsymbol{D}_{local}^{(l),m} \right\|_{F}, \left\| \widetilde{\boldsymbol{D}}_{local}^{(l),m} \right\|_{F} \leq B_{D}^{l}, & \left\| \boldsymbol{D}_{full}^{(l),m} \right\|_{F}, \left\| \widetilde{\boldsymbol{D}}_{full}^{(l),m} \right\|_{F} \leq B_{D}^{f}, \end{split}$$

where

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

Proof.

$$\left\| \boldsymbol{H}_{local}^{(l),m} \right\|_{F} = \left\| \sigma \left(\boldsymbol{P}_{local}^{(l),m} \boldsymbol{H}_{local}^{(l-1),m} \boldsymbol{W}^{(l),m} \right) \right\|_{F}$$

$$\stackrel{(a)}{\leq} C_{\sigma} B_{W} \left\| \boldsymbol{P}_{local}^{(l),m} \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F} \leq C_{\sigma} B_{W} \left\| \boldsymbol{P}_{local}^{(l),m} \right\| \left\| \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F}$$

$$\stackrel{(b)}{\leq} C_{\sigma} B_{W} B_{P} \left\| \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F} \leq \left(C_{\sigma} B_{W} B_{P} \right)^{l} \left\| \boldsymbol{X}^{m} \right\|_{F}$$

$$\stackrel{(c)}{\leq} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X} \leq \max_{1 \leq l \leq L} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X},$$

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

$$\begin{aligned} \left\| \widetilde{\boldsymbol{H}}_{local}^{(l),m} \right\|_{F} &= \left\| \sigma \left(\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \boldsymbol{W}^{(l),m} \right) \right\|_{F} \\ &\stackrel{(a)}{\leq} C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F} \leq C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} \right\| \left\| \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F} \\ &\stackrel{(b)}{\leq} C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F} \leq \left(C_{\sigma} B_{W} B_{P} \right)^{l} \left\| \boldsymbol{X}^{m} \right\|_{F} \\ &\stackrel{(c)}{\leq} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X} \leq \max_{1 \leq l \leq L} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X}, \end{aligned}$$

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

$$\begin{aligned} \left\| \boldsymbol{H}_{full}^{(l),m} \right\|_{F} &= \left\| \sigma \left(\boldsymbol{P}_{full}^{(l),m} \boldsymbol{H}_{full}^{(l-1),m} \boldsymbol{W}^{(l),m} \right) \right\|_{F} \\ &\stackrel{(a)}{\leq} C_{\sigma} B_{P} B_{W} \left\| \boldsymbol{H}_{full}^{(l-1),m} \right\|_{F} \leq \left(C_{\sigma} B_{P} B_{W} \right)^{l} \left\| \boldsymbol{X}^{m} \right\|_{F} \\ &\stackrel{(b)}{\leq} \left(C_{\sigma} B_{P} B_{W} \right)^{l} B_{X} \leq \max_{1 \leq l \leq L} \left(C_{\sigma} B_{P} B_{W} \right)^{l} B_{X}, \end{aligned}$$

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

$$\begin{split} \left\| \widetilde{\boldsymbol{H}}_{full}^{(l),m} \right\|_{F} &= \left\| \sigma \left(\left(\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m} \widetilde{\boldsymbol{H}}_{remote}^{(l-1),m} \right) \boldsymbol{W}^{(l),m} \right) \right\|_{F} \\ &\stackrel{(a)}{\leq} C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m} \widetilde{\boldsymbol{H}}_{remote}^{(l-1),m} \right\|_{F} = C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} \right\|_{F} \\ &\stackrel{(b)}{\leq} C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} \right\|_{F} \leq \left(C_{\sigma} B_{W} B_{P} \right)^{l} \left\| \boldsymbol{X}^{m} \right\|_{F} \\ &\stackrel{(c)}{\leq} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X} \leq \max_{1 \leq l \leq l} \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{X}, \end{split}$$

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

$$\begin{aligned} \left\| \boldsymbol{D}_{local}^{(l),m} \right\|_{F} &= \left\| \left[\boldsymbol{P}_{local}^{(l+1),m} \right]^{\top} \boldsymbol{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\boldsymbol{Z}_{local}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right\|_{F} \\ &\stackrel{(a)}{\leq} B_{W} C_{\sigma} \left\| \boldsymbol{P}_{local}^{(l+1),m} \right\|_{F} \left\| \boldsymbol{D}_{local}^{(l+1),m} \right\|_{F} \\ &\stackrel{(b)}{\leq} B_{P} B_{W} C_{\sigma} \left\| \boldsymbol{D}_{local}^{(l+1),m} \right\|_{F} \leq \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} \left\| \boldsymbol{D}_{local}^{(L),m} \right\|_{F} \\ &\stackrel{(c)}{\leq} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l} \leq \max_{1 \leq l \leq L} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l}, \end{aligned}$$

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

$$\left\| \widetilde{\boldsymbol{D}}_{local}^{(l),m} \right\|_{F} = \left\| \left[\widetilde{\boldsymbol{P}}_{local}^{(l+1),m} \right]^{\top} \widetilde{\boldsymbol{D}}_{local}^{(l+1),m} \circ \nabla \sigma \left(\widetilde{\boldsymbol{Z}}_{local}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right\|_{F}$$

$$\stackrel{(a)}{\leq} B_{W} C_{\sigma} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l+1),m} \right\|_{F} \left\| \widetilde{\boldsymbol{D}}_{local}^{(l+1),m} \right\|_{F}$$

$$\stackrel{(b)}{\leq} B_{P} B_{W} C_{\sigma} \left\| \widetilde{\boldsymbol{D}}_{local}^{(l+1),m} \right\|_{F} \leq \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} \left\| \widetilde{\boldsymbol{D}}_{local}^{(L),m} \right\|_{F}$$

$$\stackrel{(c)}{\leq} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l} \leq \max_{1 \leq l \leq L} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l},$$

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

$$\begin{aligned} \left\| \boldsymbol{D}_{full}^{(l),m} \right\|_{F} &= \left\| \left[\boldsymbol{P}_{full}^{(l+1),m} \right]^{\top} \boldsymbol{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\boldsymbol{Z}_{full}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right\|_{F} \\ &\stackrel{(a)}{\leq} B_{P} B_{W} C_{\sigma} \left\| \boldsymbol{D}_{full}^{(l+1),m} \right\|_{F} \leq \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} \left\| \boldsymbol{D}_{full}^{(L),m} \right\|_{F} \\ &\stackrel{(b)}{\leq} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l} \leq \max_{1 \leq l \leq L} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l}, \end{aligned}$$

where (a) and (b) use Assumptions 5.1-5.3.

$$\begin{split} \left\| \widetilde{\boldsymbol{D}}_{full}^{(l),m} \right\|_{F} &= \left\| \left[\widetilde{\boldsymbol{P}}_{local}^{(l+1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l+1),m} \right]^{\top} \widetilde{\boldsymbol{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\widetilde{\boldsymbol{Z}}_{full}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right\|_{F} \\ &= \left\| \left[\widetilde{\boldsymbol{P}}_{full}^{(l+1),m} \right]^{\top} \widetilde{\boldsymbol{D}}_{full}^{(l+1),m} \circ \nabla \sigma \left(\widetilde{\boldsymbol{Z}}_{full}^{(l+1),m} \right) \left[\boldsymbol{W}^{(l+1),m} \right]^{\top} \right\|_{F} \\ &\stackrel{(a)}{\leq} B_{P} B_{W} C_{\sigma} \left\| \widetilde{\boldsymbol{D}}_{full}^{(l+1),m} \right\|_{F} \leq \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} \left\| \widetilde{\boldsymbol{D}}_{full}^{(L),m} \right\|_{F} \\ &\stackrel{(b)}{\leq} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l} \leq \max_{1 \leq l \leq L} \left(B_{P} B_{W} C_{\sigma} \right)^{L-l} C_{l}, \end{split}$$

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

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

$$\begin{split} & \left\| \widetilde{\boldsymbol{H}}_{local}^{(l),m} - \boldsymbol{H}_{local}^{(l),m} \right\|_{F} \leq B_{\Delta H}^{l}, & \left\| \widetilde{\boldsymbol{H}}_{full}^{(l),m} - \boldsymbol{H}_{full}^{(l),m} \right\|_{F} \leq B_{\Delta H}^{f}, \\ & \left\| \widetilde{\boldsymbol{D}}_{local}^{(l),m} - \boldsymbol{D}_{local}^{(l),m} \right\|_{F} \leq B_{\Delta D}^{l}, & \left\| \widetilde{\boldsymbol{D}}_{full}^{(l),m} - \boldsymbol{D}_{full}^{(l),m} \right\|_{F} \leq B_{\Delta D}^{f}, \end{split}$$

where

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

1134
1135
$$B_{\Delta D}^{f} = \max_{1 \leq l \leq L} \left(\left(B_{W} B_{D}^{f} C_{\sigma} B_{\Delta P}^{f} + B_{W}^{2} B_{P} B_{D}^{f} L_{\sigma} B_{H}^{f} B_{\Delta P}^{f} + B_{W}^{2} B_{P}^{2} B_{D}^{f} L_{\sigma} B_{\Delta H}^{f} \right)^{L-l} + \left(B_{W} B_{P} C_{\sigma} \right)^{L-l} L_{l} B_{\Delta H}^{f} \right).$$
1137
$$+ \left(B_{W} B_{P} C_{\sigma} \right)^{L-l} L_{l} B_{\Delta H}^{f} \right).$$
1138
1139
1140
Proof.

1141
$$\left\| \widetilde{\boldsymbol{H}}_{local}^{(l),m} - \boldsymbol{H}_{local}^{(l),m} \right\|_{F}$$
1142
$$= \left\| \sigma \left(\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \boldsymbol{W}^{(l),m} \right) - \sigma \left(\boldsymbol{P}_{local}^{(l),m} \boldsymbol{H}_{local}^{(l-1),m} \right) \boldsymbol{W}^{(l),m} \right\|_{F}$$
1144
$$\stackrel{(a)}{\leq} C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{P}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F} + C_{\sigma} B_{W} \left\| \boldsymbol{P}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{P}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F}$$
1147
$$\leq C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{P}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} \right\|_{F} + C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F}$$
1150
$$\stackrel{(b)}{\leq} C_{\sigma} B_{W} B_{H}^{l} \left\| \widetilde{\boldsymbol{P}}_{local}^{(l),m} - \boldsymbol{P}_{local}^{(l),m} \right\|_{F} + C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F}$$
1151
$$\stackrel{(c)}{\leq} C_{\sigma} B_{W} B_{H}^{l} B_{\Delta P}^{l} + C_{\sigma} B_{W} B_{P} \right) \left\| \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} - \boldsymbol{H}_{local}^{(l-1),m} \right\|_{F}$$
1152
$$\stackrel{(c)}{\leq} (C_{\sigma} B_{W} B_{H}^{l} B_{\Delta P}^{l})^{l} + (C_{\sigma} B_{W} B_{P})^{l} B_{\Delta X}^{l}$$
1155
$$\stackrel{(d)}{\leq} (C_{\sigma} B_{W} B_{H}^{l} B_{\Delta P}^{l})^{l} + (C_{\sigma} B_{W} B_{P})^{l} B_{\Delta X}^{l}$$
1156
$$\stackrel{(d)}{\leq} (C_{\sigma} B_{W} B_{H}^{l} B_{\Delta P}^{l})^{l} + (C_{\sigma} B_{W} B_{P})^{l} B_{\Delta X}^{l}$$
1157
$$\stackrel{(d)}{\leq} \max_{l} \left((C_{\sigma} B_{W} B_{H}^{l} B_{\Delta P}^{l})^{l} + (C_{\sigma} B_{W} B_{P})^{l} B_{\Delta X}^{l} \right), \qquad (9)$$

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

$$\begin{aligned} &\|\widetilde{\boldsymbol{H}}_{full}^{(l),m} - \boldsymbol{H}_{full}^{(l),m}\|_{F} \\ &= \left\| \sigma\left(\left(\widetilde{\boldsymbol{P}}_{local}^{(l),m} \widetilde{\boldsymbol{H}}_{local}^{(l-1),m} + \widetilde{\boldsymbol{P}}_{remote}^{(l),m} \widetilde{\boldsymbol{H}}_{remote}^{(l-1),m}\right) \boldsymbol{W}^{(l),m}\right) - \sigma\left(\boldsymbol{P}_{full}^{(l),m} \boldsymbol{H}_{full}^{(l-1),m}\right) \boldsymbol{W}^{(l),m} \right\|_{F} \\ &\leq C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m} \boldsymbol{H}_{full}^{(l-1),m} \right\|_{F} \\ &\leq C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} \right\|_{F} + C_{\sigma} B_{W} \left\| \boldsymbol{P}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m} \boldsymbol{H}_{full}^{(l-1),m} \right\|_{F} \\ &\leq C_{\sigma} B_{W} \left\| \widetilde{\boldsymbol{P}}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m} \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} \right\|_{F} + C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{H}_{full}^{(l-1),m} \right\|_{F} \\ &\leq C_{\sigma} B_{W} B_{H}^{f} B_{\Delta P}^{f} + C_{\sigma} B_{W} B_{P} \left\| \widetilde{\boldsymbol{H}}_{full}^{(l-1),m} - \boldsymbol{H}_{full}^{(l-1),m} \right\|_{F} \\ &\leq \left(C_{\sigma} B_{W} B_{H}^{f} B_{\Delta P}^{f} \right)^{l} + \left(C_{\sigma} B_{W} B_{P} \right)^{l} \left\| \widetilde{\boldsymbol{X}}_{full}^{m} - \boldsymbol{X}_{full}^{m} \right\|_{F} \\ &\leq \left(C_{\sigma} B_{W} B_{H}^{f} B_{\Delta P}^{f} \right)^{l} + \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{\Delta X}^{f} \\ &\leq \sum_{1 \leq l \leq L} \left(\left(C_{\sigma} B_{W} B_{H}^{f} B_{\Delta P}^{f} \right)^{l} + \left(C_{\sigma} B_{W} B_{P} \right)^{l} B_{\Delta X}^{f} \right), \end{aligned}$$

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

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

$$- \left[P_{local}^{(l+1),m} \right]^{\top} D_{local}^{(l+1),m} \circ \nabla \sigma \left(Z_{local}^{(l+1),m} \right) \left[W_{l}^{(l+1),m} \right]^{\top} \right]_{F}^{(l+1),m}$$

$$\leq B_{W} \left\| \left[\tilde{P}_{local}^{(l+1),m} \right]^{\top} \tilde{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) - \left[P_{local}^{(l+1),m} \right]^{\top} D_{local}^{(l+1),m} \circ \nabla \sigma \left(Z_{local}^{(l+1),m} \right) \right\|_{F}^{(l+1),m}$$

$$\leq B_{W} \left\| \left[\tilde{P}_{local}^{(l+1),m} \right]^{\top} \tilde{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) - \left[P_{local}^{(l+1),m} \right]^{\top} \tilde{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) \right\|_{F}^{(l+1),m}$$

$$+ B_{W} \left\| \left[P_{local}^{(l+1),m} \right]^{\top} \tilde{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) - \left[P_{local}^{(l+1),m} \right]^{\top} D_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) \right\|_{F}^{(l+1),m}$$

$$+ B_{W} \left\| \left[P_{local}^{(l+1),m} \right]^{\top} \tilde{D}_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) - \left[P_{local}^{(l+1),m} \right]^{\top} D_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) \right\|_{F}^{(l+1),m}$$

$$+ B_{W} B_{W} \left\| P_{local}^{(l+1),m} \right\|_{T}^{\top} D_{local}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} \right) \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} D_{C} \left\| \tilde{P}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right\|_{F}^{-1} + B_{W} B_{P} C_{\sigma} \left\| \tilde{D}_{local}^{(l+1),m} - \nabla \sigma \left(Z_{local}^{(l+1),m} \right) \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right) \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \nabla \sigma \left(\tilde{Z}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right) \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \tilde{D}_{\sigma} \right\|_{F}^{-1} \tilde{D}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \tilde{D}_{\sigma} \right\|_{F}^{-1} \tilde{D}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \tilde{D}_{\sigma} \right\|_{F}^{-1} \tilde{D}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \tilde{D}_{\sigma} \right\|_{F}^{-1} \tilde{D}_{local}^{(l+1),m} - P_{local}^{(l+1),m} \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} B_{D} \left\| \tilde{D}_{\sigma} \right\|_{F}^{-1} \tilde{D}_{local}^{-1} - P_{local}^{(l+1),m} \right\|_{F}^{-1}$$

$$+ B_{W} B_{P} C_{\sigma} \left\| \tilde{D}_{local}^{-1} - P_{local}^{(l+1),m} \right\|_{F}^{-1} +$$

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

$$\left\|\widetilde{m{D}}_{full}^{(l),m} - m{D}_{full}^{(l),m}
ight\|_F$$

$$\begin{aligned} & = \left\| \left[\tilde{P}_{local}^{(l+1),m} + \tilde{P}_{remote}^{(l+1),m} \right] \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) \left[W^{(l+1),m} \right]^{\top} \\ & - \left[P_{full}^{(l+1),m} \right]^{\top} D_{full}^{(l+1),m} \circ \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \left[W^{(l+1),m} \right]^{\top} \right]_{F} \\ & \leq B_{W} \left\| \left[\tilde{P}_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^{\top} D_{full}^{(l+1),m} \circ \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \right\|_{F} \\ & \leq B_{W} \left\| \left[\tilde{P}_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) \right\|_{F} \\ & \leq B_{W} \left\| \left[\tilde{P}_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) \right\|_{F} \\ & \leq B_{W} \left\| \left[P_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^{\top} \tilde{D}_{full}^{(l+1),m} \circ \nabla \sigma \left(\tilde{Z}_{full}^{(l+1),m} \right) \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} + B_{W} B_{P} C_{\sigma} \left\| \tilde{D}_{full}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} + B_{W} B_{P} C_{\sigma} \left\| \tilde{D}_{full}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} + B_{W} B_{P} C_{\sigma} \left\| \tilde{D}_{full}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} \tilde{P}_{full}^{(l),m} - P_{full}^{(l+1),m} \tilde{P}_{full}^{(l),m} - P_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} \tilde{P}_{full}^{(l),m} - P_{full}^{(l+1),m} \tilde{P}_{full}^{(l),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} B_{D}^{L} \left\| \tilde{P}_{full}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_{F} \\ & \leq B_{W} B_{D}^{L} C_{\sigma} B_$$

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

Lemma E.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.,

$$\left\| \boldsymbol{H}_{local}^{(l),m} - \boldsymbol{H}_{full}^{(l),m} \right\|_{F} \le B_{\Delta H}^{r}, \qquad \left\| \boldsymbol{D}_{local}^{(l),m} - \boldsymbol{D}_{full}^{(l),m} \right\|_{F} \le B_{\Delta D}^{r},$$

where

$$B_{\Delta H}^{r} = \max_{1 \leq l \leq L} \left((C_{\sigma} B_{W} B_{P})^{l} B_{X}^{r} + \left(C_{\sigma} B_{W} B_{H}^{f} B_{P} \right)^{l} \right),$$

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

Proof.

$$\|\boldsymbol{H}_{local}^{(l),m} - \boldsymbol{H}_{full}^{(l),m}\|_{F} = \|\sigma\left(\boldsymbol{P}_{local}^{(l),m}\boldsymbol{H}_{local}^{(l-1),m}\right)\boldsymbol{W}^{(l),m} - \sigma\left(\boldsymbol{P}_{full}^{(l),m}\boldsymbol{H}_{full}^{(l-1),m}\right)\boldsymbol{W}^{(l),m}\|_{F}$$

$$\leq C_{\sigma}B_{W} \|\boldsymbol{P}_{local}^{(l),m}\boldsymbol{H}_{local}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m}\boldsymbol{H}_{full}^{(l-1),m}\|_{F}$$

$$\leq C_{\sigma}B_{W} \|\boldsymbol{P}_{local}^{(l),m}\boldsymbol{H}_{local}^{(l-1),m} - \boldsymbol{P}_{local}^{(l),m}\boldsymbol{H}_{full}^{(l-1),m}\|_{F}$$

$$+ C_{\sigma}B_{W} \|\boldsymbol{P}_{local}^{(l),m}\boldsymbol{H}_{full}^{(l-1),m} - \boldsymbol{P}_{full}^{(l),m}\boldsymbol{H}_{full}^{(l-1),m}\|_{F}$$

$$\leq C_{\sigma}B_{W}B_{P} \|\boldsymbol{H}_{local}^{(l-1),m} - \boldsymbol{H}_{full}^{(l-1),m}\|_{F} + C_{\sigma}B_{W}B_{H}^{f} \|\boldsymbol{P}_{local}^{(l),m} - \boldsymbol{P}_{full}^{(l),m}\|_{F}$$

$$\leq C_{\sigma}B_{W}B_{P} \|\boldsymbol{H}_{local}^{(l-1),m} - \boldsymbol{H}_{full}^{(l-1),m}\|_{F} + C_{\sigma}B_{W}B_{H}^{f} \|\boldsymbol{P}_{remote}^{(l),m}\|_{F}$$

$$\leq C_{\sigma}B_{W}B_{P} \|\boldsymbol{H}_{local}^{(l-1),m} - \boldsymbol{H}_{full}^{(l-1),m}\|_{F} + C_{\sigma}B_{W}B_{H}^{f}B_{P}$$

$$\leq (C_{\sigma}B_{W}B_{P})^{l} \|\boldsymbol{X}_{local}^{m} - \boldsymbol{X}_{full}^{m}\|_{F} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l}$$

$$\leq (C_{\sigma}B_{W}B_{P})^{l} B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l}$$

$$\leq \max_{1 \leq l \leq l} \left((C_{\sigma}B_{W}B_{P})^{l}B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l} \right),$$

$$\leq \max_{1 \leq l \leq l} \left((C_{\sigma}B_{W}B_{P})^{l}B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l} \right),$$

$$\leq \max_{1 \leq l \leq l} \left((C_{\sigma}B_{W}B_{P})^{l}B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l} \right),$$

$$\leq \max_{1 \leq l \leq l} \left((C_{\sigma}B_{W}B_{P})^{l}B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l} \right),$$

$$\leq \max_{1 \leq l \leq l} \left((C_{\sigma}B_{W}B_{P})^{l}B_{X}^{r} + \left(C_{\sigma}B_{W}B_{H}^{f}B_{P}\right)^{l} \right),$$

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

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

$$\begin{aligned} & + B_W \left\| \left[P_{full}^{(l+1),m} \right]^\top D_{local}^{(l+1),m} \circ \nabla \sigma \left(Z_{local}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^\top D_{full}^{(l+1),m} \circ \nabla \sigma \left(Z_{local}^{(l+1),m} \right) \right\|_F \\ & + B_W \left\| \left[P_{full}^{(l+1),m} \right]^\top D_{full}^{(l+1),m} \circ \nabla \sigma \left(Z_{local}^{(l+1),m} \right) - \left[P_{full}^{(l+1),m} \right]^\top D_{full}^{(l+1),m} \circ \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \right\|_F \\ & + B_W B_D C_\sigma \left\| P_{local}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_F + B_W B_D C_\sigma \left\| D_{local}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_F \\ & + B_W B_D B_D^I \left\| \nabla \sigma \left(Z_{local}^{(l+1),m} \right) - \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \right\|_F \\ & + B_W B_D B_D^I \left\| \nabla \sigma \left(Z_{local}^{(l+1),m} \right) - \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \right\|_F \\ & + B_W B_D B_D^I \left\| \nabla \sigma \left(Z_{local}^{(l+1),m} \right) - \nabla \sigma \left(Z_{full}^{(l+1),m} \right) \right\|_F \\ & + B_W B_D B_D^I \left\| P_{local}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_F + B_W B_D C_\sigma \left\| D_{local}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_F \\ & + B_W^2 B_D B_D^I L_\sigma \left\| P_{local}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_F + B_W B_D C_\sigma \left\| D_{local}^{(l+1),m} - D_{full}^{(l+1),m} \right\|_F \\ & + B_W^2 B_D B_D^I L_\sigma \left\| P_{local}^{(l+1),m} + H_{local}^{(l,m)} - P_{full}^{(l+1),m} + H_{local}^{(l,m)} \right\|_F \\ & + B_W^2 B_D B_D^I L_\sigma \left\| P_{full}^{(l+1),m} + H_{local}^{(l,m)} - P_{full}^{(l+1),m} + H_{local}^{(l,m)} \right\|_F \\ & + B_W^2 B_D B_D^I L_\sigma \left\| P_{full}^{(l+1),m} + D_{full}^{(l+1),m} + B_W^2 B_D^2 B_D^I L_\sigma \left\| H_{local}^{(l+1),m} - H_{full}^{(l,m)} \right\|_F \\ & + B_W^2 B_D B_D^I L_\sigma B_H^i \left\| P_{local}^{(l+1),m} - P_{full}^{(l+1),m} \right\|_F + B_W^2 B_D^2 B_D^I L_\sigma B_\Delta^r \right\|_F \\ & + B_W^2 B_D^2 D_\sigma B_D^2 + B_W^2 B_D^2 B_D^I L_\sigma B_H^i + B_W^2 B_D^2 B_D^I L_\sigma B_\Delta^r \right)^{L-1} \\ & + B_W^2 B_D^2 C_\sigma B_D + B_W^2 B_D^2 B_D^I L_\sigma B_H^i + B_W^2 B_D^2 B_D^I L_\sigma B_\Delta^r \right)^{L-1} \\ & + (B_W B_D C_\sigma)^{L-1} \left\| D_{local}^{(L+1),m} - D_{full}^{(L+1),m} \right\|_F \\ & \leq \left(B_W B_D^I C_\sigma B_D + B_W^2 B_D^2 B_D^I L_\sigma B_H^i + B_W^2 B_D^2 B_D^I L_\sigma B_\Delta^r \right)^{L-1} \\ & + (B_W B_D C_\sigma)^{L-1} L_l \left\| H_{local}^{(L),m} - H_{full}^{(L),m} \right\|_F \\ & \leq \left(B_W B_D^I C_\sigma B_D + B_W^2 B_D^2 B_D^I L_\sigma B_H^i + B_W^2 B_D^2 B_D^I L_\sigma B_\Delta^r \right)^{L-1} \\ & +$$

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

E.3 ERRORS OF STOCHASTIC GRADIENTS

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

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

where
$$B_{\Delta G}^{l} = \max_{1 \leq l \leq L} \left((B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{H}^{l}B_{\Delta P}^{l} + B_{P}B_{H}^{l}C_{\sigma}B_{\Delta D}^{l} \right. \\ \left. + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{H}^{l}B_{\Delta P}^{l} + B_{P}B_{H}^{l}C_{\sigma}B_{\Delta D}^{l} \right. \\ \left. + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right), \qquad (12)$$

$$B_{\Delta G}^{l} = \max_{1 \leq l \leq L} \left((B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right), \qquad (13)$$

$$H_{1412}^{l+12} \qquad B_{\Delta G}^{l} = \max_{1 \leq l \leq L} \left((B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right), \qquad (13)$$

$$H_{1413}^{l+144} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right) \qquad (13)$$

$$H_{1414}^{l+141} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right) \qquad (13)$$

$$H_{1415}^{l+141} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right) \qquad (13)$$

$$H_{1414}^{l+141} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right) \qquad (13)$$

$$H_{1414}^{l+141} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{P}B_{\Delta H}^{l} \right) \qquad (13)$$

$$H_{1414}^{l+141} \qquad + (B_{D}^{l}C_{\sigma} + B_{P}B_{H}^{l}B_{D}^{l}L_{\sigma}B_{W}) B_{D}^{l}B_{D}^{l}L_{\sigma}B_{W} \right) B_{P}B_{\Delta H}^{l}$$

$$H_{1414}^{l}B_{D}^{l}C_{\sigma}B_{D}^{l}B_{D}^{l}C_{\sigma}B_{D}^{l}B_{D}^{l}C_{\sigma}B_{D}^{l}C_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}B_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D}^{l}B_{D}^{l}B_{D}^{l}C_{D}^{l}B_{D$$

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

1456

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:

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

$$\| \widetilde{G}_{full}^{(l),m} - G_{full}^{(l),m} \|_{F}$$

$$= \| [\widetilde{P}_{local}^{(l),m} \widetilde{H}_{local}^{(l-1),m} + \widetilde{P}_{remote}^{(l),m} \widetilde{H}_{remote}^{(l-1),m}]^{\top} \widetilde{D}_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right)$$

$$- [P_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(Z_{full}^{(l),m} \right) \|_{F}$$

$$\leq \| [\widetilde{P}_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m}]^{\top} \widetilde{D}_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) - [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} \widetilde{D}_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) \|_{F}$$

$$+ \| [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} \widetilde{D}_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) - [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) \|_{F}$$

$$+ \| [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) - [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) \|_{F}$$

$$+ \| [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) - [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) \|_{F}$$

$$+ \| [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) - [P_{full}^{(l),m} H_{full}^{(l-1),m}]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \right) \|_{F}$$

$$+ \| P_{full}^{(l),m} H_{full}^{(l-1),m} - P_{full}^{(l),m} H_{full}^{(l-1),m} \|_{F} + B_{F} B_{H}^{l} C_{\sigma} \| \widetilde{D}_{full}^{(l),m} - D_{full}^{(l),m} \|_{F}$$

$$+ \| B_{F} B_{H}^{l} B_{D}^{l} \| \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} - P_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} \|_{F}$$

$$+ \| B_{F} B_{H}^{l} B_{D}^{l} \| \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} - P_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} \|_{F}$$

$$+ \| B_{F} B_{H}^{l} B_{D}^{l} \| \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} - P_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} \|_{F}$$

$$+ \| B_{F} B_{H}^{l} B_{D}^{l} \| \nabla \sigma \left(\widetilde{Z}_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} - P_{full}^{(l),m} \widetilde{H}_{full}^{(l-1),m} \|_{F}$$

$$+ \| B_{F} B_{H}^{l} B_{D}^{l} \| \nabla \sigma \left($$

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

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

Lemma E.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:

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

where

$$B_{\Delta G}^{r} = \max_{1 \le l \le L} \left(\left(B_{D}^{l} C_{\sigma} + B_{P} B_{H}^{f} B_{D}^{f} L_{\sigma} B_{W} \right) B_{P} B_{\Delta H}^{r} + B_{P} B_{H}^{f} C_{\sigma} B_{\Delta D}^{r} + \left(B_{D}^{l} C_{\sigma} + B_{P} B_{H}^{f} B_{D}^{f} L_{\sigma} B_{W} \right) B_{H}^{f} B_{P} \right)$$

$$(14)$$

Proof.

$$\begin{aligned} & \| \mathbf{G}_{local}^{(l),m} - \mathbf{G}_{full}^{(l),m} \|_{F} \\ & = \| \left[\mathbf{P}_{local}^{(l),m} H_{local}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \|_{F} \\ & \leq \| \left[\mathbf{P}_{local}^{(l),m} H_{local}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \|_{F} \\ & + \| \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \|_{F} \\ & + \| \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) \|_{F} \\ & + \| \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{local}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{local}^{(l),m} \right) - \left[\mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \right]^{\top} D_{full}^{(l),m} \circ \nabla \sigma \left(\mathbf{Z}_{full}^{(l),m} \right) \|_{F} \\ & + \| \mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \|_{F} + \mathbf{P}_{F}_{B}_{H}^{H} C_{\sigma} \| \mathbf{D}_{local}^{(l),m} - \mathbf{D}_{full}^{(l),m} \|_{F} \\ & + \| \mathbf{P}_{F}_{H}^{H}_{B}_{D}^{L} \mathbf{D}_{F}_{B}_{W} \| \mathbf{P}_{local}^{(l),m} H_{full}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} H_{full}^{(l-1),m} \|_{F} \\ & + \| \mathbf{P}_{F}_{H}^{H}_{B}_{D}^{L} \mathbf{D}_{F}_{B}_{W} \| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \|_{F} \\ & + \| \mathbf{P}_{F}_{H}^{H}_{B}_{D}^{L} \mathbf{D}_{F}_{B}_{W} \| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \|_{F} \\ & + \| \mathbf{P}_{F}_{B}_{H}^{H}_{B}_{D}^{L} \mathbf{D}_{F}_{B}_{W} \| \mathbf{P}_{local}^{(l),m} \mathbf{H}_{local}^{(l-1),m} - \mathbf{P}_{full}^{(l),m} \mathbf{H}_{full}^{(l-1),m} \|_{F} \\ & + \| \mathbf{P}_{F}_{B}_{H}^{H}_{B}_{D}^{L}_{B}_{B}_{W} \| \mathbf{P}_{B}_{B}^{L}_{B}_{B}_{W} \| \mathbf{P}_{hall}^{(l-1),m} - \mathbf{P}_{full}^{(l-1),m} \|_{F} \\$$

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

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

E.4 Main Proof of Theorem 5.6

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

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

Proof.

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

$$\overset{(a)}{\leq} \left\langle \nabla \mathcal{L}\left(\boldsymbol{\theta}_{t}\right), \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \right\rangle + \frac{L_{F}}{2} \left\| \boldsymbol{\theta}_{t+1} - \boldsymbol{\theta}_{t} \right\|^{2}$$

$$\stackrel{(b)}{=} -\alpha \left\langle \nabla \mathcal{L} \left(\boldsymbol{\theta}_{t} \right), \frac{1}{M} \sum_{\boldsymbol{\delta}, \boldsymbol{\delta}} \nabla \widetilde{F}^{m} \left(\boldsymbol{\theta}_{t}^{m} \right) \right\rangle + \frac{L_{F}}{2} \alpha^{2} \left\| \frac{1}{M} \sum_{\boldsymbol{\delta}, \boldsymbol{\delta}} \nabla \widetilde{F}^{m} \left(\boldsymbol{\theta}_{t}^{m} \right) \right\|^{2}$$

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

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

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

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

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

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

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

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

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

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

1620
$$\stackrel{(a)}{\leq} LB_{\Delta G}^r + LB_{\Delta G}^l,$$
 (16)

where (a) follows from Lemmas E.6 and E.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

$$\left\| \nabla F^m \left(\boldsymbol{\theta}_t^m \right) - \nabla \widetilde{F}^m \left(\boldsymbol{\theta}_t^m \right) \right\| = \left\| \nabla F_{full}^m \left(\boldsymbol{\theta}_t^m \right) - \nabla \widetilde{F}_{full}^m \left(\boldsymbol{\theta}_t^m \right) \right\| \overset{(a)}{\leq} L B_{\Delta G}^f,$$

where (a) uses Lemma E.6.

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

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

Choosing $T = n_t I$ yields

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

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

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

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