



Carnegie Mellon University

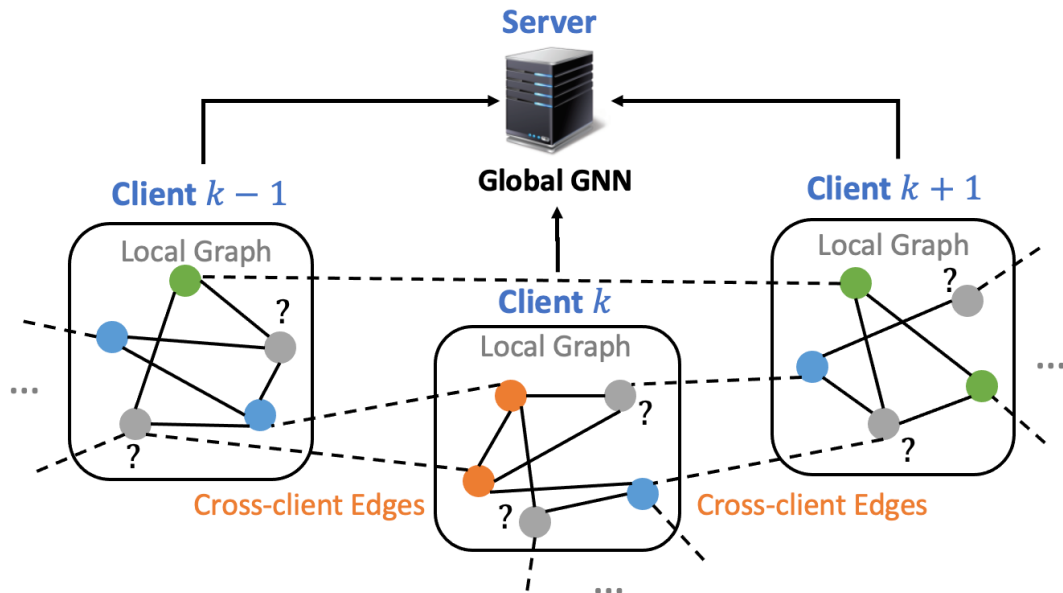
FedGCN: Convergence-Communication Tradeoffs in Federated Training of Graph Convolutional Networks

Presenter: Yuhang Yao

Author: Yuhang Yao, Weizhao Jin, Srivatsan Ravi, Carlee Joe-Wong

Federated Node Classification

- Nodes in a graph are **partitioned** across clients (e.g. private data across countries)
- **Cross-client edges** exist between nodes at different clients



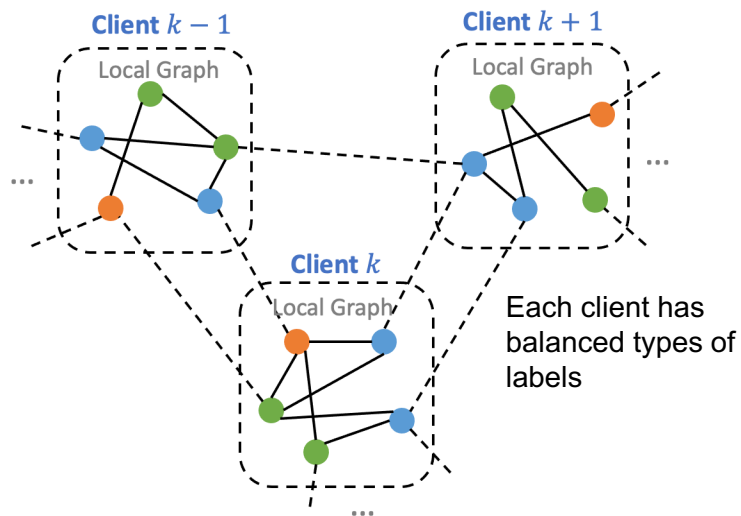
Each client knows

- Local graph structure
- Local node features
- Corresponding **cross-client edges**

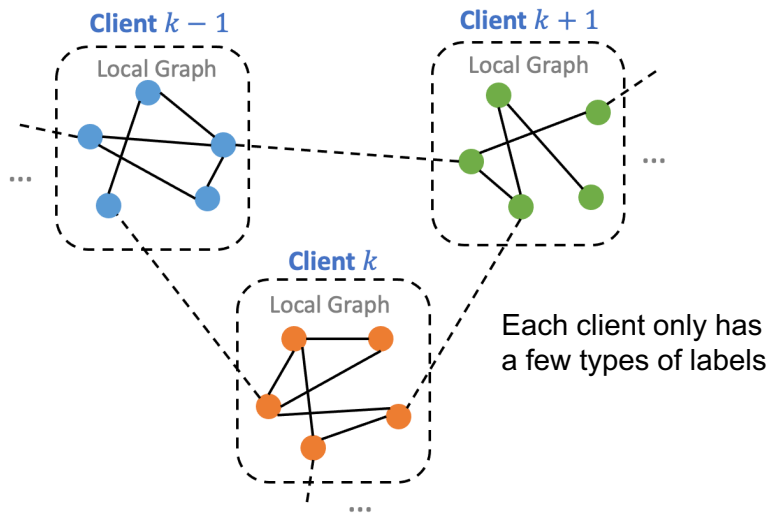
Node classification requires node features stored in other clients

Edges in Heterogeneous Data Distribution

IID (Independently Identically Distributed)



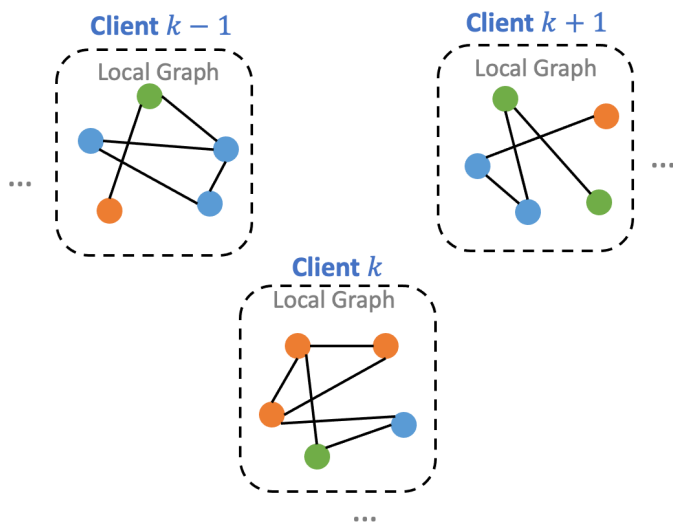
Non-IID



- Nodes connect more to nodes with the same label
- Non-IID may have fewer cross-client edges than IID
- More cross-client edges require more communication

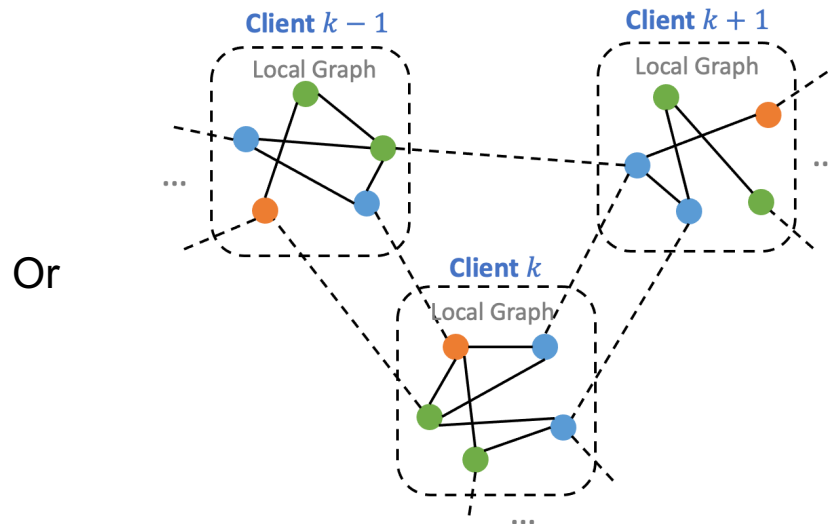
Limitation of Distributed Training

Ignore cross-client edges



Ignoring cross-client edges causes **information loss**

Send features and intermediate output at every round



Or

Sending features requires **huge communication cost**

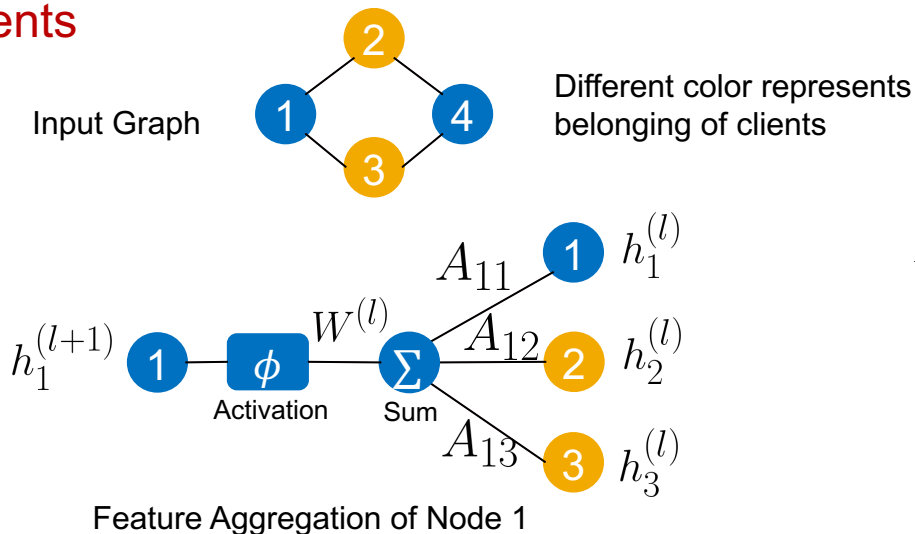
[1] He, Chaoyang, et al. "Fedgraphnn: A federated learning system and benchmark for graph neural networks." *arXiv preprint arXiv:2104.07145* (2021).
 [2] Wan, Cheng, et al. "BDS-GCN: Efficient full-graph training of graph convolutional nets with partition-parallelism and boundary sampling." (2020).

GCN in Federated Learning

In FL setting, nodes are stored in different clients

For each layer l

Node i in client $c(i)$ needs to aggregate information of nodes from $c(i)$ and **other clients**



$$\mathbf{h}_i^{(l+1)} = \phi \left(\sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{h}_j^{(l)} W_{c(i)}^{(l)} \right)$$

A_{ij} : Weight of connections between node i and node j

$h_i^{(l)}$: Output of node i at layer l

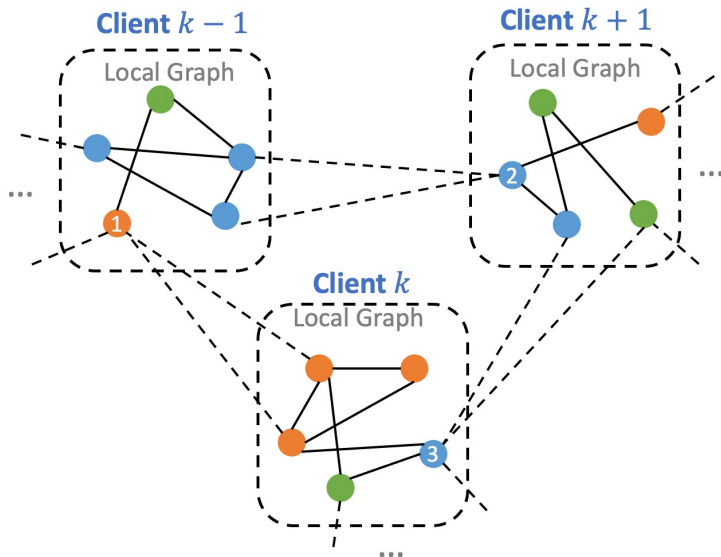
$c(i)$: index of the client that contains node i

$W_{c(i)}^{(l)}$: Parameters of GCN at layer l at client $c(i)$

$h_i^{(0)} = x_i$: Feature vector of node i at layer 0

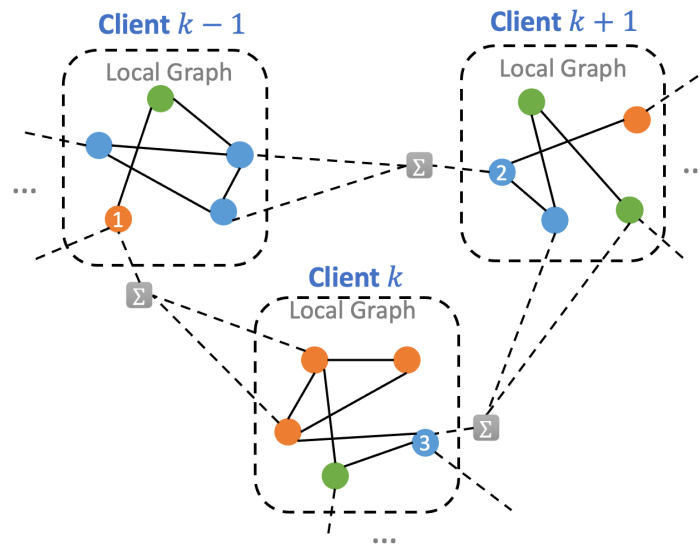
Feature Aggregation Instead of Sending Features

Send features and intermediate output at **every training round**



- High communication cost at every round

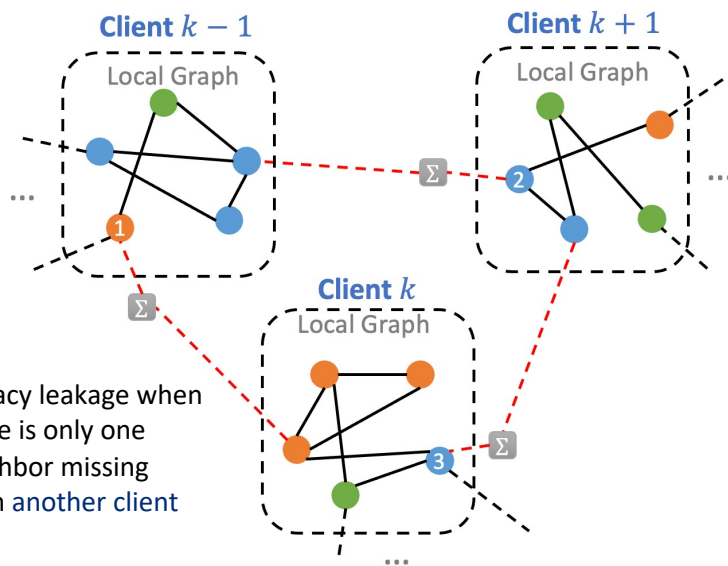
Send feature aggregations at **initial round**



- Same computation
- Much lower communication cost

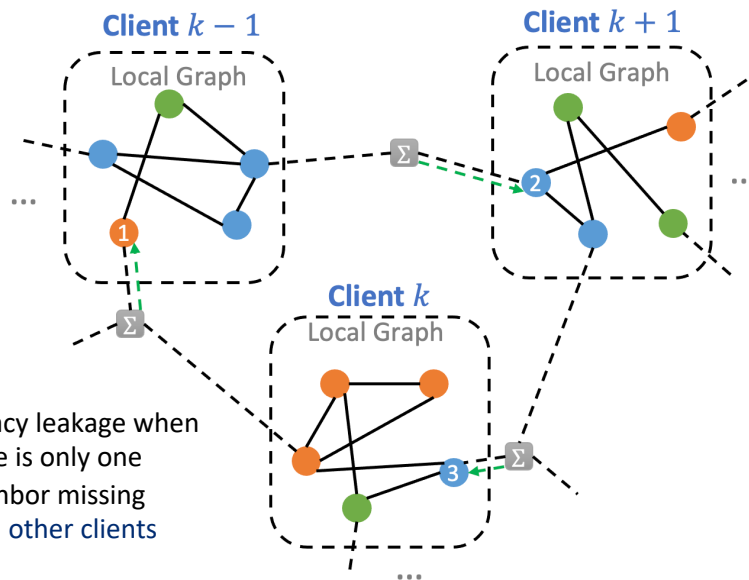
Server Aggregation Instead of Clients Aggregation

Clients Aggregation



$$\{\sum_{j \in \mathcal{N}_i} \mathbb{I}_z(c(j)) A_{ij} \mathbf{x}_j\}_{z \in [K]}$$

Server Aggregation



$$\sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{x}_j = \sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \mathbb{I}_k(c(j)) \cdot A_{ij} \mathbf{x}_j$$

Secure Neighbor Feature Aggregation

To guarantee privacy during the aggregation process of accumulated features, we leverage **Fully Homomorphic Encryption (FHE)**

$$\left[\left[\sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{x}_j \right] \right] = \sum_{k=1}^K \left[\left[\sum_{j \in \mathcal{N}_i} \mathbb{I}_k(c(j)) \cdot A_{ij} \mathbf{x}_j \right] \right]$$

1. All clients agree on and initialize a FHE keypair
2. Each client encrypts the local neighbor feature array and sends it to the server
3. Upon receiving all encrypted neighbor feature arrays from clients, the server performs secure neighbor feature aggregation

FedGCN with Three Types of Communication

- **No Communication(0-hop):** Use feature aggregation at the same client
- **1-hop Communication:** Communicate feature aggregation of 1-hop neighbors at all clients
- **2-hop Communication:** Communicate feature aggregation of 2-hop neighbors at all clients
- More hop means higher communication costs but with less information loss
- 2-hop communication does not have information loss for 2-layer GCN

Training Process of FedGCN

Communication at
initial round

Normal federated
training process

Algorithm 1 FedGCN Federated Training for Graph Convolutional Network

```

// Communication Round
for each client  $k \in [K]$  do in parallel
  Send  $\llbracket \{\sum_{j \in \mathcal{N}_i} \mathbb{1}_k(c(j)) \cdot A_{ij} \mathbf{x}_j\}_{i \in V_k} \rrbracket$  to the server
end
// Server Operation
for  $i \in V$  do in parallel
   $\llbracket \sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{x}_j \rrbracket = \sum_{k=1}^K \llbracket \sum_{j \in \mathcal{N}_i} \mathbb{1}_k(c(j)) \cdot A_{ij} \mathbf{x}_j \rrbracket$ 
end
for each client  $k \in [K]$  do in parallel
  if 1-hop then
    Receive  $\llbracket \{\sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{x}_j\}_{i \in V_k} \rrbracket$  and decrypt it
  end
  if 2-hop then
    Receive  $\llbracket \{\sum_{j \in \mathcal{N}_i} A_{ij} \mathbf{x}_j\}_{i \in \mathcal{N}_{V_k}} \rrbracket$  and decrypt it
  end
end
// Training Round
for  $t = 1, \dots, T$  do
  for each client  $k \in [K]$  do in parallel
    Receive  $\llbracket \mathbf{w}^{(t)} \rrbracket$  and decrypt it
    Set  $\mathbf{w}_k^{(t,1)} = \mathbf{w}^{(t)}$ 
    for  $e = 1, \dots, \tau$  do
      Set  $\mathbf{g}_{\mathbf{w}_k}^{(t,e)} = \nabla_{\mathbf{w}_k} f_k(\mathbf{w}_k^{(t,e)}; G_k)$ 
       $\mathbf{w}_k^{(t,e+1)} = \mathbf{w}_k^{(t,e)} - \eta \mathbf{g}_{\mathbf{w}_k}^{(t,e)}$  // Update Parameters
    end
     $\Delta \mathbf{w}_k^{(t,\tau)} = \mathbf{w}_k^{(t,\tau+1)} - \mathbf{w}_k^{(t,1)}$ 
    Send  $\llbracket \Delta \mathbf{w}_k^{(t,\tau)} \rrbracket$  to the server
  end
  // Server Operations
   $\llbracket \Delta \mathbf{w} \rrbracket = \frac{1}{K} \sum_{k=1}^K \llbracket \Delta \mathbf{w}_k^{(t,\tau)} \rrbracket$  // Difference6 Aggregation
   $\llbracket \mathbf{w}^{(t+1)} \rrbracket = \llbracket \mathbf{w}^{(t)} \rrbracket - \llbracket \Delta \mathbf{w} \rrbracket$  and broadcast to local clients // Update Global Models
end

```

Experiment Setups

Citation Network Datasets

Dataset	Nodes	Edges	Features	Classes
Cora	2,708	5,429	1,433	7
Citeseer	3,327	4,732	3,703	6
Pubmed	19,717	44,338	500	3
Ogbn-Arxiv	169,343	1,166,243	128	40

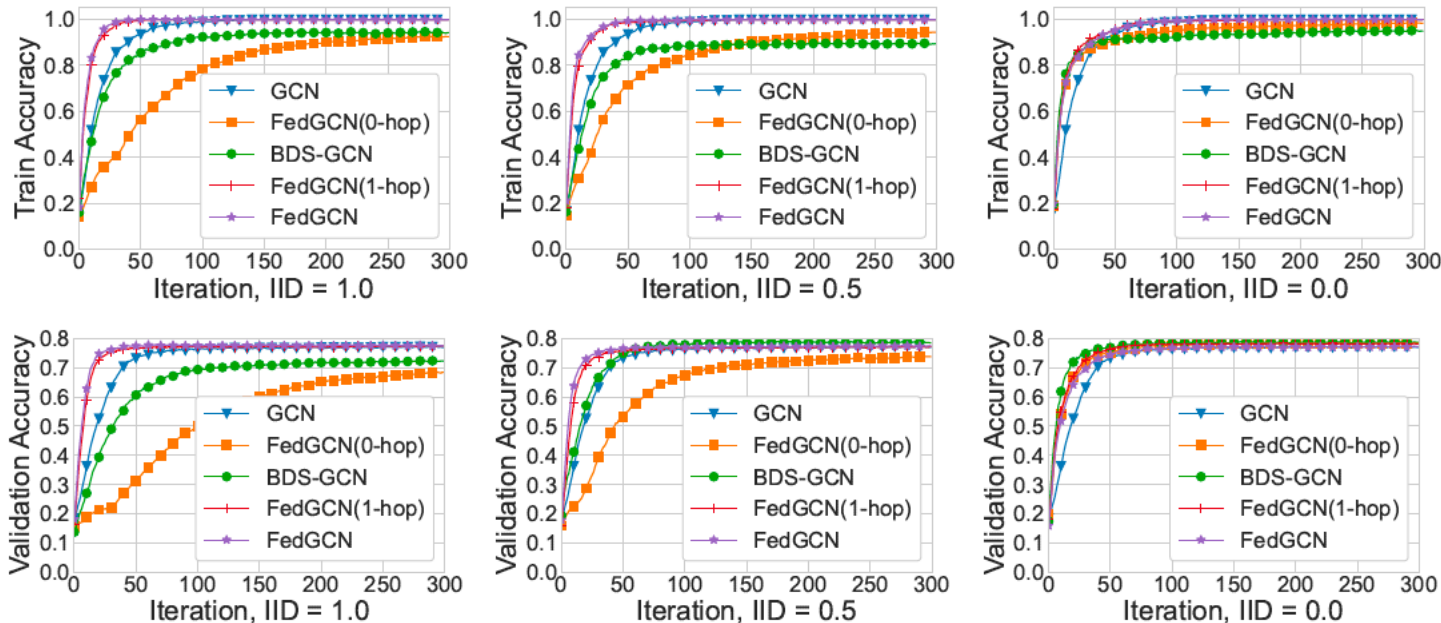
Compared methods

- Centralized GCN
- FedGCN(0-hop), FedGraphnn[1]
- BDS-GCN: Randomly samples cross-client edges
- FedSage+: Approximates 1-hop neighbors [2]
- FedGCN(1-hop)
- FedGCN(2-hop)

[1] He, Chaoyang, et al. "Fedgraphnn: A federated learning system and benchmark for graph neural networks." *arXiv* (2021).

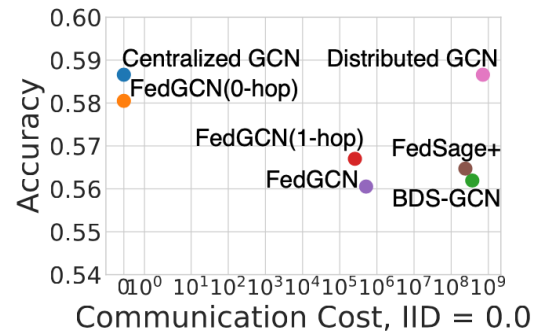
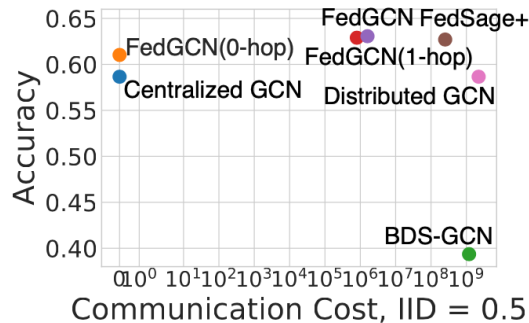
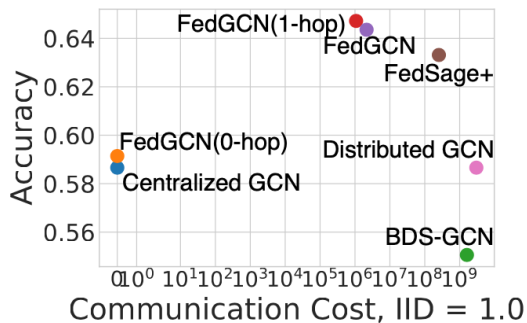
[2] Zhang, Ke, et al. "Subgraph federated learning with missing neighbor generation." *NeurIPS* (2021).

Training and Test Accuracy on Cora



- FedGCN converges much faster and has a higher test and training accuracy in all settings.
- Under the extreme non-i.i.d. setting, FedGCN (0-hop) has sufficient information to train a good model.

Test Accuracy vs Communication Cost on OGBN-ArXiv



- FedGCN (0-, 1-, and 2-hop) requires little communication with high accuracy
- FedGCN (0-hop) requires much less communication, but has lower accuracy due to information loss in the i.i.d. and partial-i.i.d. settings

Convergence Rate

0-hop

1-hop

2-hop

i.i.d. (SBM)

$$\left(1 - \frac{1}{K^4}\right) \frac{N^5}{K^5} \|B^4\|$$

$$\left(1 - \frac{1}{K^4} (1 + c_\alpha + c_\mu)^2\right) \frac{N^5}{K^5} \|B^4\|$$

$$\left(1 - \frac{1}{K^4} (1 + c_\alpha + c_\mu)^6\right) \frac{N^5}{K^5} \|B^4\|$$

In i.i.d., the 2-hop method has faster convergence rate

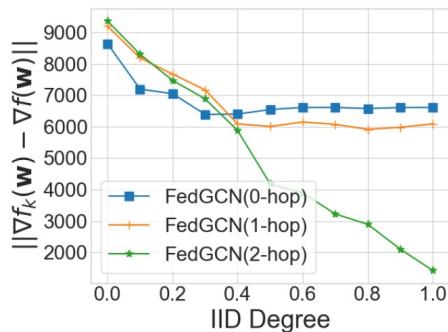
Non-i.i.d. (SBM)

$$\left(1 - \frac{1}{K^4}\right) \frac{N^5}{K^5} \|B^4\| + \sigma$$

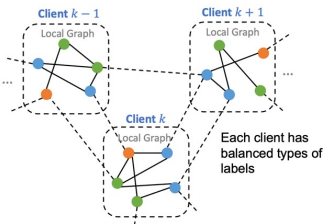
$$\left(1 - \frac{1}{K^4} (1 + c_\mu)^2\right) \frac{N^5}{K^5} \|B^4\| + \sigma$$

$$\left(1 - \frac{1}{K^4} (1 + c_\mu)^6\right) \frac{N^5}{K^5} \|B^4\| + \sigma$$

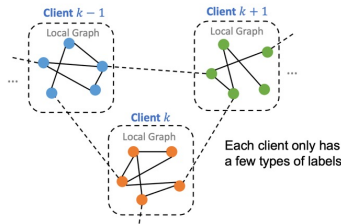
In extremely Non-i.i.d., the 1-hop and 2-hop methods help little.



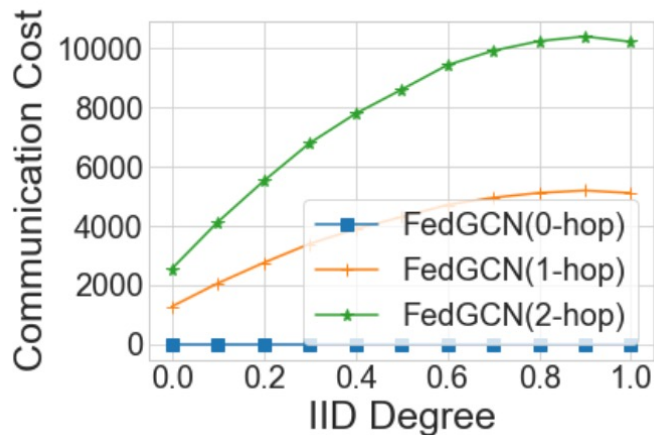
IID (Independently Identically Distributed)



Non-IID



Communication Cost



2-hop requires more communication but helps the convergence

Data Distribution	0-hop	1-hop	2-hop
Generic Graph	0	$\sum_{i \in V} c(\mathcal{N}_i) d + Nd$	$\sum_{i \in V} c(\mathcal{N}_i) d + \sum_{k=1}^K \mathcal{N}_{V_k} d$
Non-i.i.d. (SBM)	0	$(c_\mu + 2)Nd$	$2(c_\mu + 1)Nd$
Partial-i.i.d. (SBM)	0	$(c_\alpha p + c_\mu + 2)Nd$	$2(c_\alpha p + c_\mu + 1)Nd$
i.i.d. (SBM)	0	$(c_\alpha + c_\mu + 2)Nd$	$2(c_\alpha + c_\mu + 1)Nd$

Conclusion & Next Step

Conclusion

- Cross-client edges affect the model performance (convergence rate and test accuracy).
- Proposed FedGCN helps recover information on cross-client edges and only requires communication at the initial step
- Tradeoffs exist between convergence and communication under different data distributions.

Next Steps

- Large-scale experiments and more compared methods
- System deployment and library development