# On Which Nodes Does GCN Fail? Enhancing GCN From the Node Perspective

Jincheng Huang[1], Jialie Shen[2], Xiaoshuang Shi[1]*, Xiaofeng Zhu[1]*

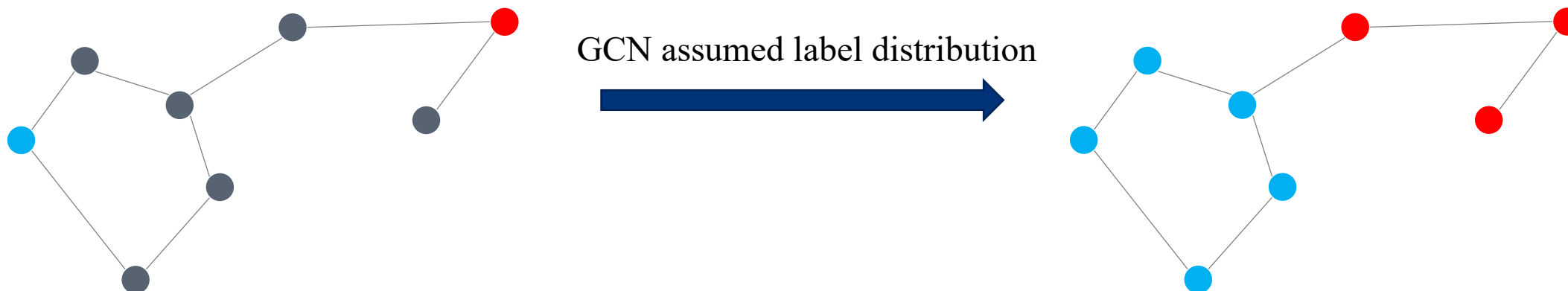1.University of Electronic Science and Technology of China 2. City, University of London, London

# Background

| Model | Operation | Result |
|---|---|---|
| Graph Convolutional Networks (GCNs) | Feature Smoothing | Connected nodes have similar features |



GCNs excel at handling graph-structured data, with most methods relying on their feature smoothing operations.

# Background

What kind of graph data does GCN expect?



GCN assumed label distribution

GCNs assume that Connect nodes are highly likely to share the same labels. (i.e., label smoothness assumption)(Zhang et al., 2021)

> Question: Is the label distribution obtained by GCN feature smoothing consistent with the label smoothness assumption?
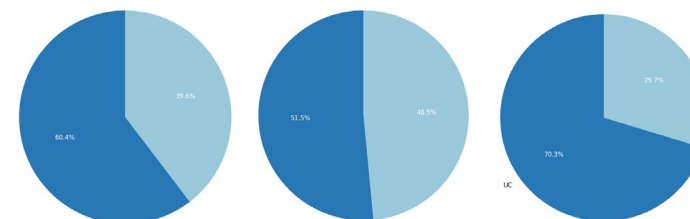
## Theorem 1

For nodes with unknown labels in the graph, the upper bound of the GCN's generalization ability reaches optimal if the true labels of these nodes are equal to the labels generated by the LPA.

● Theorem 1 establishes the link between the output of LPA and the expected label distribution of GCN (i.e., label smoothness assumption)

## Label-Feature Smoothing Alignment Algorithm

1. GCN feature smoothing: $\mathbf{Y}_{fs} = \widehat{\mathbf{A}}^L MLP(\mathbf{X})$      2. GCN label smoothness assumption: $\mathbf{Y}_{lp} = \widehat{\mathbf{A}}^L \mathbf{Y}$

$$\mathbf{V}_{OOC} = \{\mathbf{V}_i | argmax(\mathbf{Y}_{fs,i}) \neq argmax(\mathbf{Y}_{lp,i}), i \in [n]\} \quad \mathbf{V}_{UC} = \mathbf{V} - \mathbf{V}_{OOC}$$

**Answer the Question:** There is a fairly significant proportion of nodes that are not consistent.
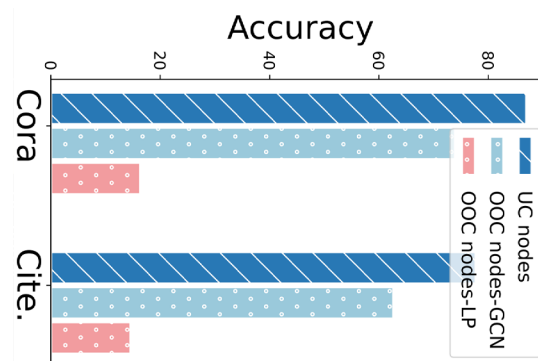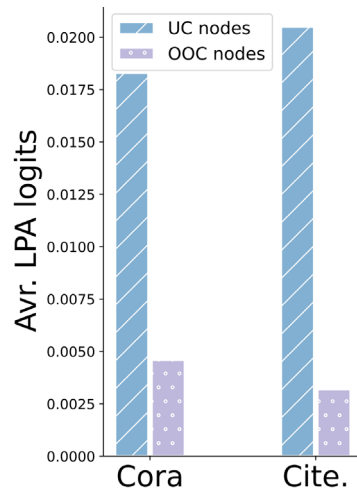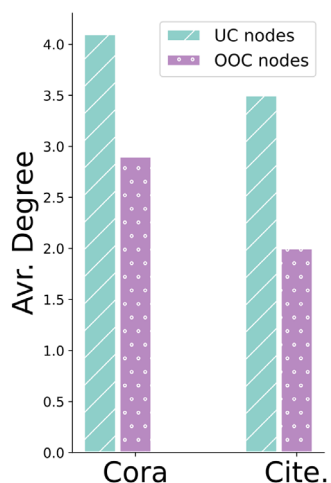
Cora      Citeseer      Pubmed

**Unlabeled Nodes**

## UC Nodes

**Nodes** that achieve label smoothing assumptions using GCN feature smoothing operations are under the control of GCN.



Accuracy

## OOC Nodes

**Nodes** affected by GCN's feature smoothing operation conflict with the label smoothness assumption, making it difficult to correct representation under the GCN framework.



**Character of OOC nodes.**

(i)  Nodes with few neighbors (left figure).

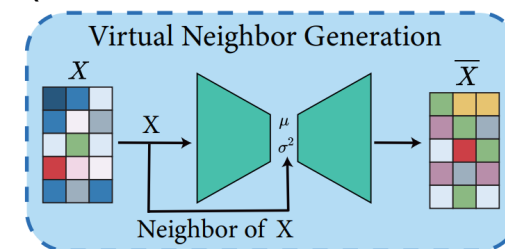(ii)  Nodes away from labeled nodes (right figure).

## For Nodes with Few Neighbors

### Virtual Neighbor Generation

Use $X_v(v \in V)$ as a condition, and to learn the neighbor distribution of $X_u(u \in N_v)$ (Liu et al., 2022, Sohn et al., 2015).

$$\mathcal{L}_{ELBO} = -KL(q(\mathbf{z} \mid \mathbf{X}_u \mathbf{X}_v) \| p(\mathbf{z} \mid \mathbf{X}_v))$$
$$+ \mathbb{E}_{q(\mathbf{z} \mid \mathbf{X}_u, \mathbf{x}_v)}(p(\mathbf{X}_u \mid \mathbf{X}_v, \mathbf{z}))$$

This process allows us to obtain the node v's virtual neighbor feature vector $\overline{\mathbf{X}}_v$.



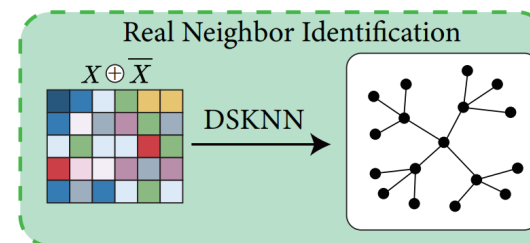Virtual Neighbor Generation

### Potential Real Neighbor Identification

Virtual nodes contain only first-order information and can't affect message passing. We posit potential non-directly connected neighbors can augment message passing for OOC nodes if:
- They are in the same subspace.
- Their neighbors are in the same subspace.

$$\min_{\mathbf{S}} \sum_{i,j=0}^{n} \left( -s_{i,j}\mathbf{X}_i^T\mathbf{X}_j - s_{i,j}\overline{\mathbf{X}}_i^T\overline{\mathbf{X}}_j + s_{i,j}^2 \right)$$

$$s_{i,j} = \frac{1}{2}\left( \mathbf{X}_i^T\mathbf{X}_j + \overline{\mathbf{X}}_i^T\overline{\mathbf{X}}_j \right) = \frac{1}{2}\left( \mathbf{X}_i \oplus \overline{\mathbf{X}}_i \right)^T \left( \mathbf{X}_j \oplus \overline{\mathbf{X}}_j \right)$$



Real Neighbor Identification

## Nodes away from labeled nodes

### Theorem 2

Given an undirected graph G(V, E) has n nodes and m edges. Assuming there are q nodes in the graph with labels selected uniformly at random. The occurrence probability of nodes that are not affected by labels with a two-layer GCN is equal to
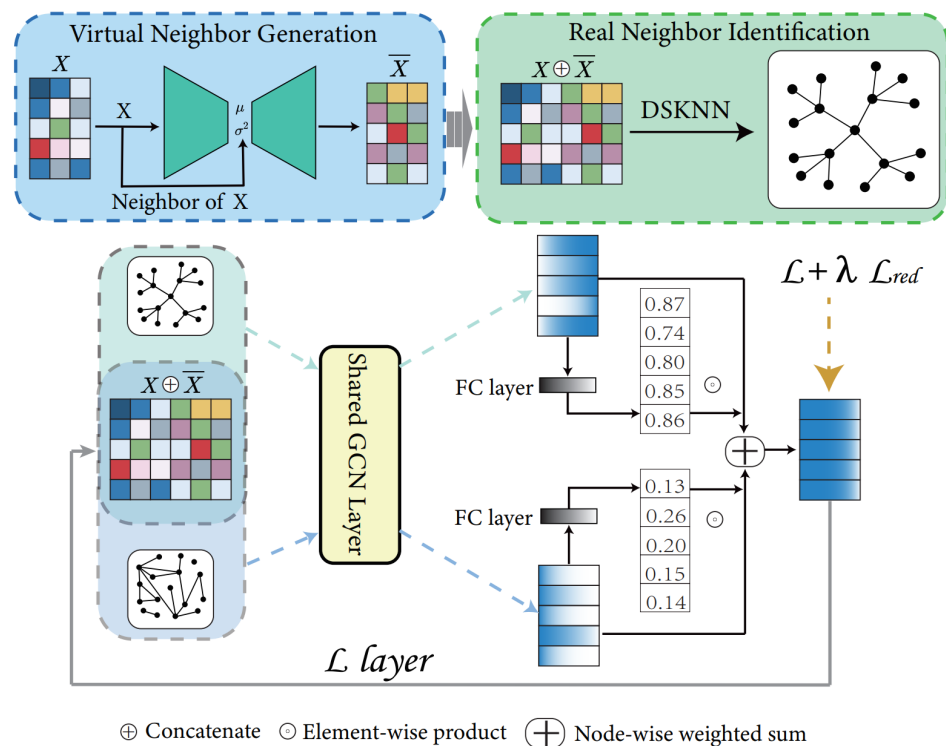
$$(1 - \frac{q}{n})(1 - \frac{q}{n-1})\prod_{i=1}^{q}(1 - \frac{2m}{n(n-1)-2i})\prod_{i=q}^{2q}(1 - \frac{2(m-1)}{n(n-1)-2i})$$

- Theorem 2 tells us the occurrence probability of unaffected by labeled nodes is negatively correlated with the number of labels and total edges.

- DSKNN-graph
  - 1. Can reduce the probability of OOC nodes.
  - 2. Allowing flexible addition or removal of edges.

| Number of UC nodes | Cora | Citeseer | Pubmed |
|---|---|---|---|
| Original Graph | 633 | 485 | 708 |
| DSKNN Graph | 660 | 711 | 720 |
| Combine Graph | 833 | 840 | 879 |
| Improve Ratio | 31.6% | 73.2% | 24.2% |

**Solution**: we just need to make sure that the number of edges in constructing the DSKNN graph is much larger than the average degree of the original graph.

## Overall Architecture



Virtual Neighbor Generation

Real Neighbor Identification

Shared GCN Layer

FC layer

$\mathcal{L} + \lambda \, \mathcal{L}_{red}$

$\mathcal{L}$ layer

⊕ Concatenate  ⊙ Element-wise product  ⊕ Node-wise weighted sum

▶ Concatenate virtual neighbors'' feature as input feature:

$$\overline{\mathcal{X}} = \mathbf{X} \oplus \overline{\mathbf{X}}$$

▶ Propagating the features on the original graph and the DSKNN graph :

$$\mathbf{H}_{ori}^{(l)} = \widehat{\mathbf{A}}\mathbf{H}^{(l-1)}\mathbf{W}^{(l-1)}, \mathbf{H}_{ds}^{(l)} = \widehat{\mathbf{S}}\mathbf{H}^{(l-1)}\mathbf{W}^{(l-1)}$$

▶ Adaptive node-level assembling:

$$\mathbf{H}^{(l)} = \mathrm{diag}\left(\boldsymbol{\lambda}_0^{(l)}\right)\mathbf{H}_{ori}^{(l)} + \mathrm{diag}\left(\boldsymbol{\lambda}_1^{(l)}\right)\mathbf{H}_{ds}^{(l)}, \quad \boldsymbol{\lambda}_0^{(l)} + \boldsymbol{\lambda}_1^{(l)} = \mathbf{1}$$
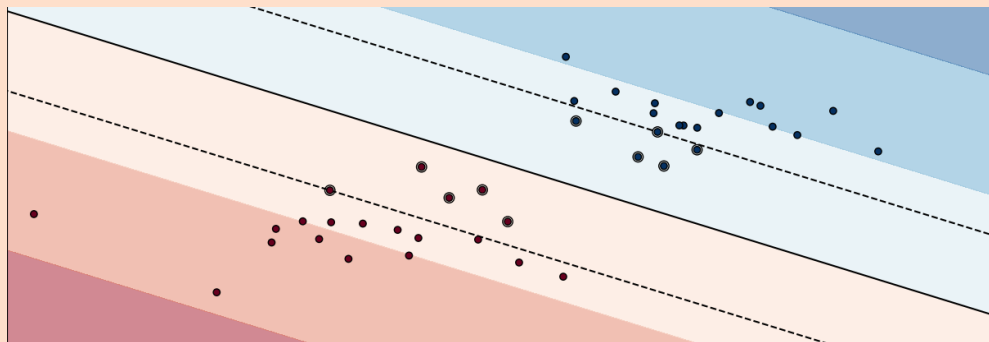
$$\boldsymbol{\lambda}_0^{(l)} = \sigma\left(FC_0^{(l)}\left(\mathbf{H}_{ori}^{(l)}\right)\right), \boldsymbol{\lambda}_1^{(l)} = \sigma\left(FC_1^{(l)}\left(\mathbf{H}_{ds}^{(l)}\right)\right)$$

$$\left[\boldsymbol{\lambda}_0^{(l)}, \boldsymbol{\lambda}_1^{(l)}\right] = \frac{\left[\boldsymbol{\lambda}_0^{(l)}, \boldsymbol{\lambda}_1^{(l)}\right]}{\max\left(\left\|\left[\boldsymbol{\lambda}_0^{(l)}, \boldsymbol{\lambda}_1^{(l)}\right]\right\|_2, \epsilon\right)}$$

## Some Problem

### Fundamental Assumption in Semi-Supervised Learning

In semi-supervised learning, the classifier's decision boundary should avoid high-density regions of the data distribution.



### An accomplish way

Ensuring the classifier outputs low-entropy predictions on unlabeled data.



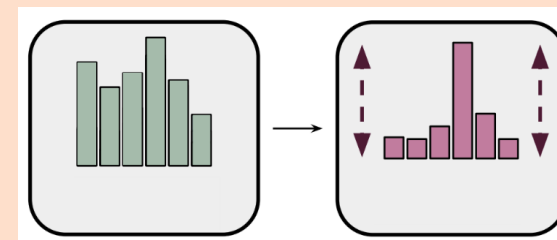**In Adaptive node-level assembling**

**Violation**

**Solution**

### Lemma 1

$$H(\lambda p_1 + (1-\lambda)p_2) \geq \lambda H(p_1) + (1-\lambda)H(p_2)$$

- When the output layer assembles the logits, the entropy will increase beyond a linear combination of the two view.

**Entropy Reduction Loss:**

$$\mathcal{L}_{red} = \frac{1}{c}\sum_{i=1}^{c}(\mathbf{y}_i - \mathbf{y}_i^{\frac{1}{\tau}} \Big/ \sum_{j}^{c} \mathbf{y}_j^{\frac{1}{\tau}})^2 + \mathbb{I}(||\mathbf{y}_{ori} - \mathbf{y}_{ds}||_2)$$

**Main Results**

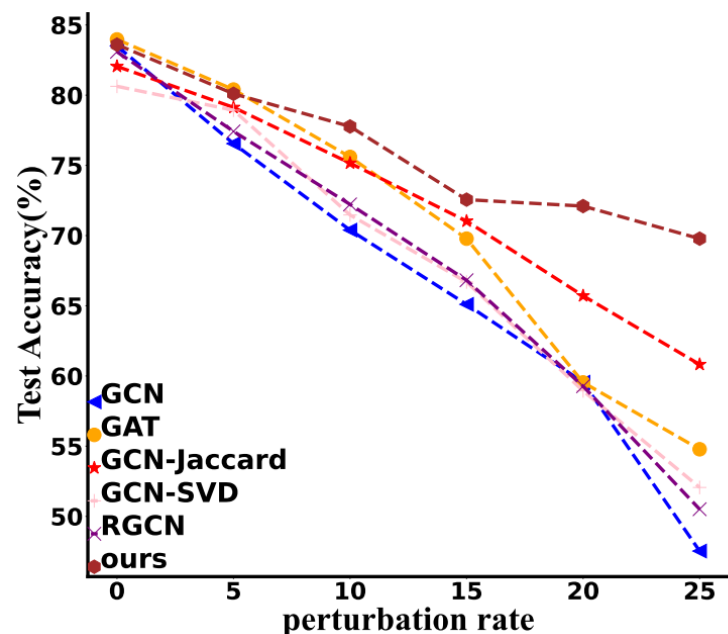| Datasets | Cora | Citeseer | Pubmed | Computers | Photo | Physics | CS |
|---|---|---|---|---|---|---|---|
| GCN | $81.5_{\pm0.82}$ | $70.9_{\pm0.71}$ | $79.0_{\pm0.52}$ | $82.6_{\pm2.43}$ | $91.2_{\pm1.21}$ | $92.8_{\pm1.00}$ | $91.1_{\pm0.52}$ |
| GAT | $83.0_{\pm0.41}$ | $71.1_{\pm0.51}$ | $79.1_{\pm0.44}$ | $78.0_{\pm19.0}$ | $85.7_{\pm20.3}$ | $92.5_{\pm0.94}$ | $90.5_{\pm0.61}$ |
| APPNP | $83.3_{\pm0.51}$ | $72.5_{\pm0.62}$ | $79.9_{\pm0.32}$ | $82.2_{\pm2.13}$ | $90.8_{\pm1.32}$ | $93.7_{\pm0.69}$ | $92.5_{\pm0.32}$ |
| GCN-LPA | $83.1_{\pm0.73}$ | $72.6_{\pm0.80}$ | $78.6_{\pm1.32}$ | $83.5_{\pm1.41}$ | $91.1_{\pm0.83}$ | $93.6_{\pm1.06}$ | $91.8_{\pm0.42}$ |
| DAGNN | $84.4_{\pm0.57}$ | $73.3_{\pm0.65}$ | $80.5_{\pm0.53}$ | $83.5_{\pm1.28}$ | $92.0_{\pm1.22}$ | $94.0_{\pm0.62}$ | $91.5_{\pm0.33}$ |
| $w$GCN | $83.1_{\pm0.31}$ | $73.9_{\pm0.46}$ | $80.8_{\pm0.25}$ | $83.6_{\pm0.86}$ | $92.4_{\pm0.18}$ | $92.8_{\pm0.23}$ | $89.3_{\pm0.14}$ |
| AERO-GNN | $83.9_{\pm0.51}$ | $73.2_{\pm0.68}$ | $80.6_{\pm0.55}$ | $83.3_{\pm0.72}$ | $91.1_{\pm0.83}$ | $93.3_{\pm0.65}$ | $92.0_{\pm0.71}$ |
| Ours | $\mathbf{84.8}_{\pm0.53}$ | $\mathbf{75.3}_{\pm0.41}$ | $\mathbf{81.7}_{\pm0.88}$ | $\mathbf{84.0}_{\pm1.25}$ | $\mathbf{92.9}_{\pm0.56}$ | $\mathbf{94.3}_{\pm0.25}$ | $\mathbf{93.4}_{\pm0.18}$ |

- In the node classification task, our proposed method outperformance the SOTA baseline.

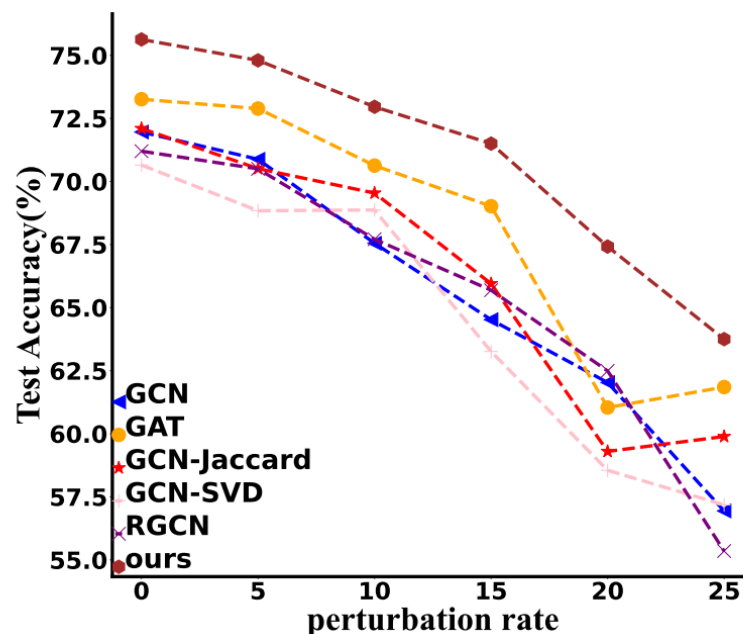| Datasets | Cora | | Citeseer | | Pubmed | | Computers | | Photo | | Physics | | CS | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | UC nodes | OOC nodes | UC nodes | OOC nodes | UC nodes | OOC nodes | UC nodes | OOC nodes | UC nodes | OOC nodes | UC nodes | OOC nodes | UC nodes | OOC nodes |
| GCN | $87.01_{\pm0.6}$ | $73.95_{\pm1.1}$ | $77.75_{\pm0.5}$ | $62.66_{\pm0.8}$ | $83.66_{\pm0.3}$ | $67.05_{\pm1.0}$ | $87.41_{\pm0.5}$ | $70.13_{\pm0.8}$ | $96.58_{\pm0.7}$ | $78.02_{\pm1.4}$ | $97.01_{\pm0.2}$ | $86.45_{\pm0.3}$ | $95.65_{\pm0.4}$ | $84.53_{\pm0.6}$ |
| APPNP | $87.47_{\pm0.5}$ | $76.21_{\pm1.3}$ | $78.21_{\pm0.6}$ | $67.59_{\pm0.9}$ | $84.36_{\pm0.5}$ | $67.96_{\pm1.1}$ | $87.23_{\pm0.9}$ | $69.84_{\pm2.6}$ | $95.98_{\pm0.8}$ | $78.13_{\pm1.5}$ | $97.13_{\pm0.5}$ | $89.25_{\pm0.9}$ | $95.31_{\pm0.2}$ | $87.01_{\pm0.5}$ |
| DAGNN | $87.80_{\pm0.5}$ | $78.52_{\pm1.5}$ | $78.33_{\pm0.7}$ | $68.27_{\pm0.93}$ | $84.48_{\pm0.8}$ | $68.32_{\pm0.7}$ | $88.21_{\pm0.7}$ | $71.97_{\pm1.5}$ | $95.36_{\pm0.8}$ | $80.64_{\pm1.2}$ | $97.16_{\pm0.5}$ | $89.98_{\pm0.7}$ | $94.75_{\pm0.2}$ | $87.53_{\pm0.7}$ |
| AERO-GNN | $87.74_{\pm0.3}$ | $77.38_{\pm0.8}$ | $78.14_{\pm0.8}$ | $68.78_{\pm1.0}$ | $85.38_{\pm0.3}$ | $69.79_{\pm1.1}$ | $88.56_{\pm0.8}$ | $71.72_{\pm1.3}$ | $96.34_{\pm0.6}$ | $77.65_{\pm1.0}$ | $97.03_{\pm0.4}$ | $88.65_{\pm0.9}$ | $95.89_{\pm0.6}$ | $86.01_{\pm1.1}$ |
| Ours | $87.70_{\pm0.5}$ | $\mathbf{79.26}_{\pm0.7}$ | $78.39_{\pm0.5}$ | $\mathbf{72.04}_{\pm1.5}$ | $85.54_{\pm0.3}$ | $\mathbf{73.16}_{\pm1.0}$ | $88.12_{\pm0.8}$ | $\mathbf{73.39}_{\pm1.5}$ | $95.57_{\pm0.5}$ | $\mathbf{82.75}_{\pm0.8}$ | $97.12_{\pm0.2}$ | $\mathbf{91.15}_{\pm0.3}$ | $95.53_{\pm0.1}$ | $\mathbf{89.51}_{\pm0.3}$ |

- Most methods (including ours) show similar effectiveness on UC nodes. The key factor differentiating their performance is their behavior on OOC nodes. Thus, research on GCNs should primarily focus on OOC nodes.

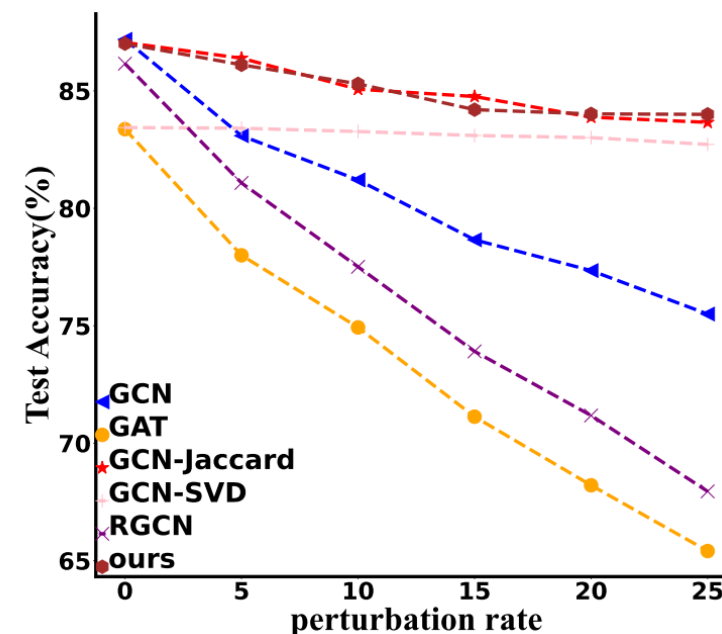- Our proposed method significantly improves the performance of GCNs on OOC nodes.

# Experiments

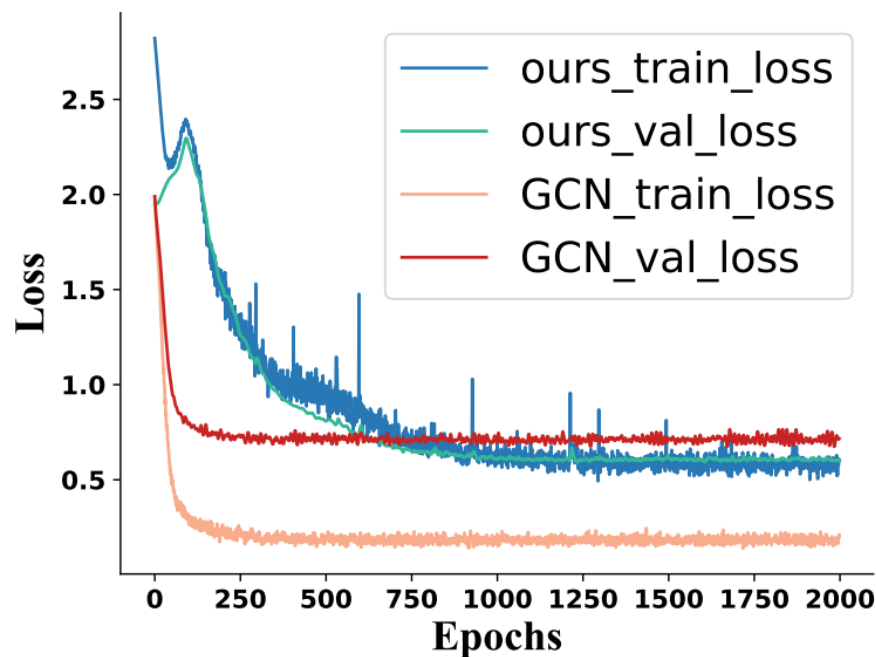**Adversarial Robustness-Metaattack**
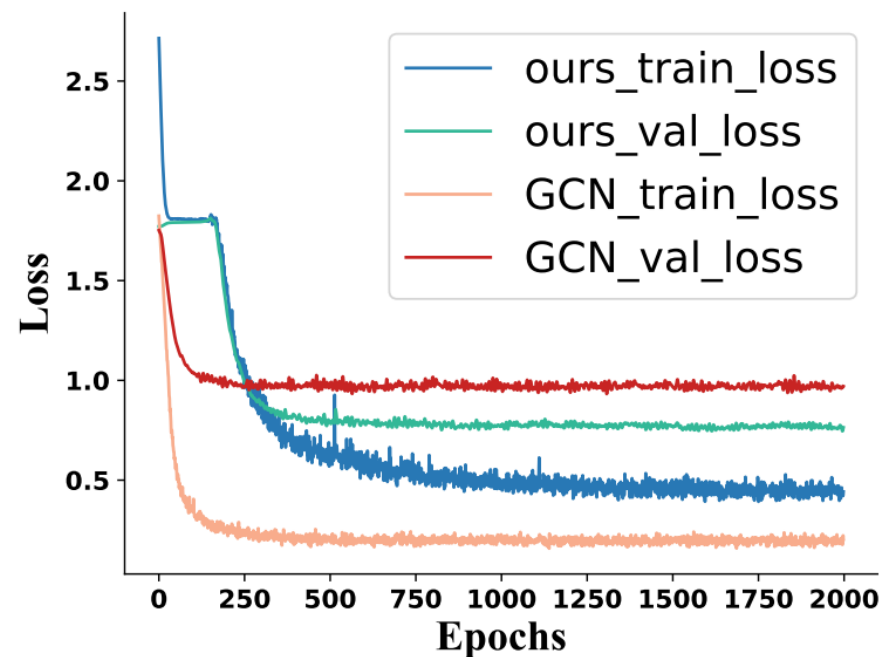


(a) Cora

(b) Citeseer

(c) Pubmed

- Our proposed method has strong adversarial robustness
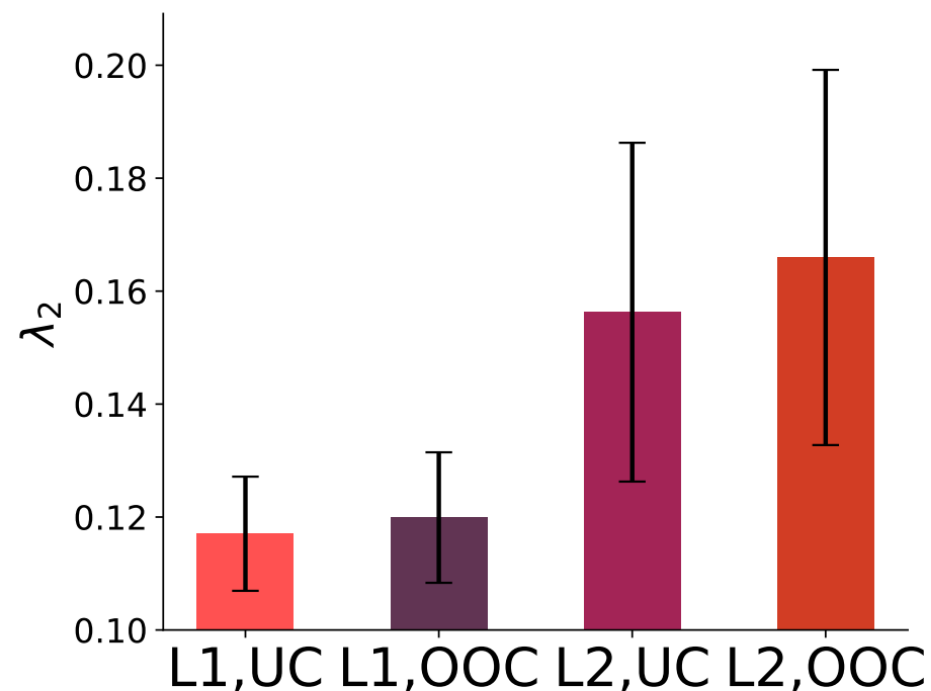
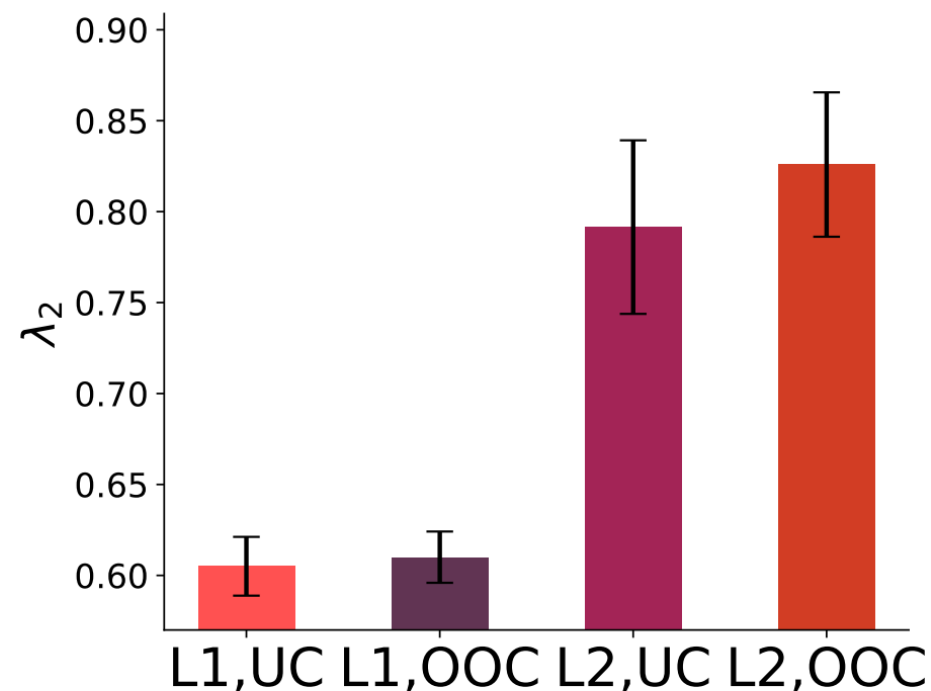**Analysis Generalization Ability**



(a) Cora

(b) Citeseer

- Our proposed improves the GCN's generalization ability.

**Analysis Adaptive Node-level Assembling**



(a) Cora



(b) Pubmed

- The OOC nodes have heavier average weights in the second layer of the DSKNN side compared to the UC nodes, suggesting greater benefit for OOC nodes from the DSKNN side.
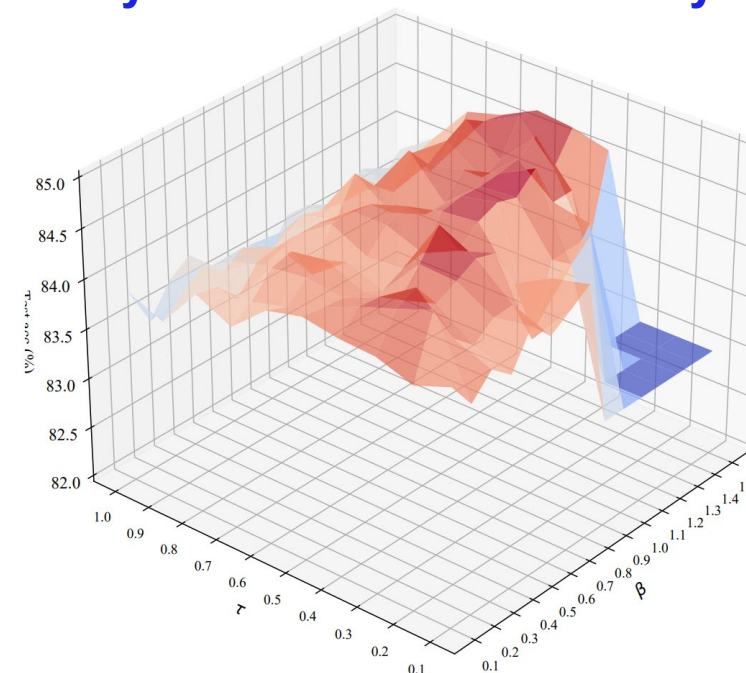- The weights learned by each layer are differentiated.

# Experiments

**Analysis Adaptive Node-level Assembling**

| Ablation | Cora | Citeseer | Pubmed |
|---|---|---|---|
| DaGCN | $84.8_{\pm0.53}$ | $75.3_{\pm0.41}$ | $81.7_{\pm0.88}$ |
| - w/o VNG | $84.2_{\pm0.96}$ | $74.5_{\pm0.66}$ | $81.2_{\pm1.00}$ |
| - w/o RNG | $83.6_{\pm0.46}$ | $73.6_{\pm0.37}$ | $80.7_{\pm0.62}$ |
| - w/o ERL | $84.0_{\pm0.56}$ | $73.8_{\pm0.72}$ | $81.3_{\pm0.75}$ |
| GCN | $81.5_{\pm0.82}$ | $70.9_{\pm0.71}$ | $79.0_{\pm0.52}$ |

**Analysis Parameter Sensitivity**



- All components are valid.
- The DSKNN-graph part played the biggest effect.

- the temperature parameter τ is significantly important, since when τ is in the interval [0.4, 0.8], the model performance maintains an excellent level. The DSKNN-graph part played the biggest effect.
- if we ensure that τ is in a suitable range, the selection of β is not sensitive.

# Summary

**Conclusion**

- vanilla GCN has been able to achieve high-quality representation learning on UC nodes. The advanced model should focus on improving OOC nodes to promote GCN.
- We provide algorithms for locating OOC nodes and provide directions and models to promote OOC nodes.

**Future Work**

- Optimize graph structure from the perspective of reducing OOC nodes and Generalization.