

***EGonc: Energy – based Open – Set Node Classification
with substitute Unknowns***

Qin Zhang¹, Zelin Shi¹, Shirui Pan², Junyang Chen¹, Huisi Wu¹, and Xiaojun Chen^{1*}

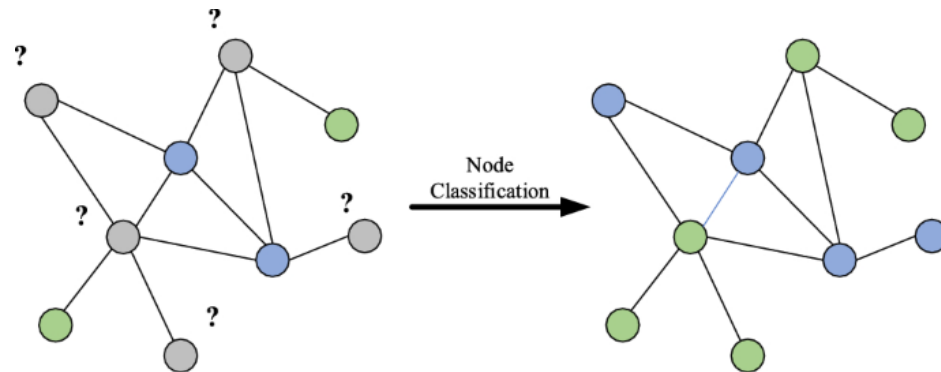
***Shenzhen University¹
Griffith University²***

NeurIPS(2024)

Introduction

■ Background

- Node classification is the task of predicting the labels of unlabeled nodes in a graph.
- SOTA methods based on graph neural networks achieve excellent performance when all labels are available during training.
- But in real-life, models are often applied on data with new classes, which can lead to massive misclassification and thus significantly degrade performance.
- Hence, developing open-set classification methods is crucial to resolve this issue.



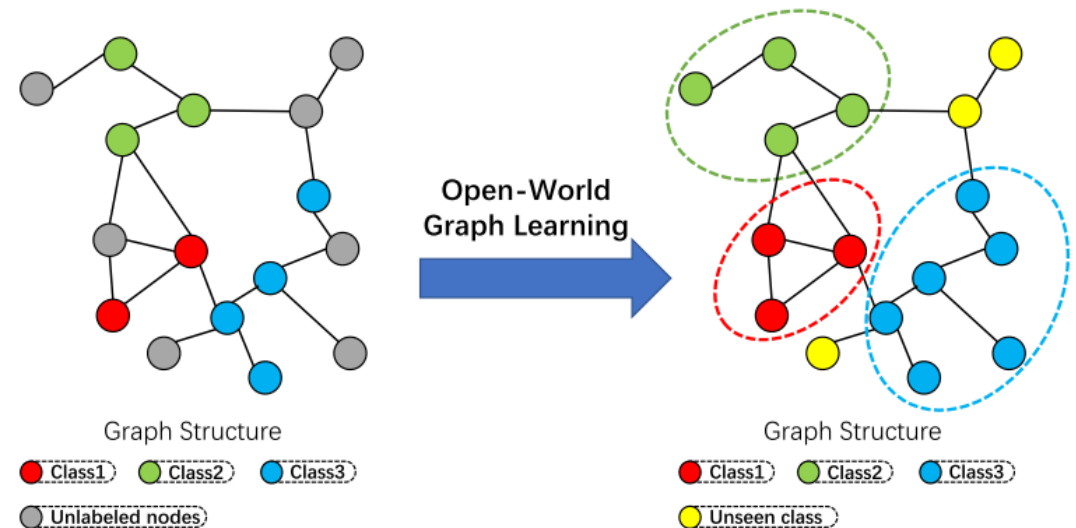
Introduction

■ Motivation

Major challenges in Graph Learning:

mostly based on closed-world assumptions, lacking generalization ability

- Restricting the category space to remain consistent between the training and the testing stages.
- Most methods based on the open-world assumption adhere to a transductive setting.
- Most open-set node classification methods on graphs are based on discriminative or generative models, lacking new approaches.



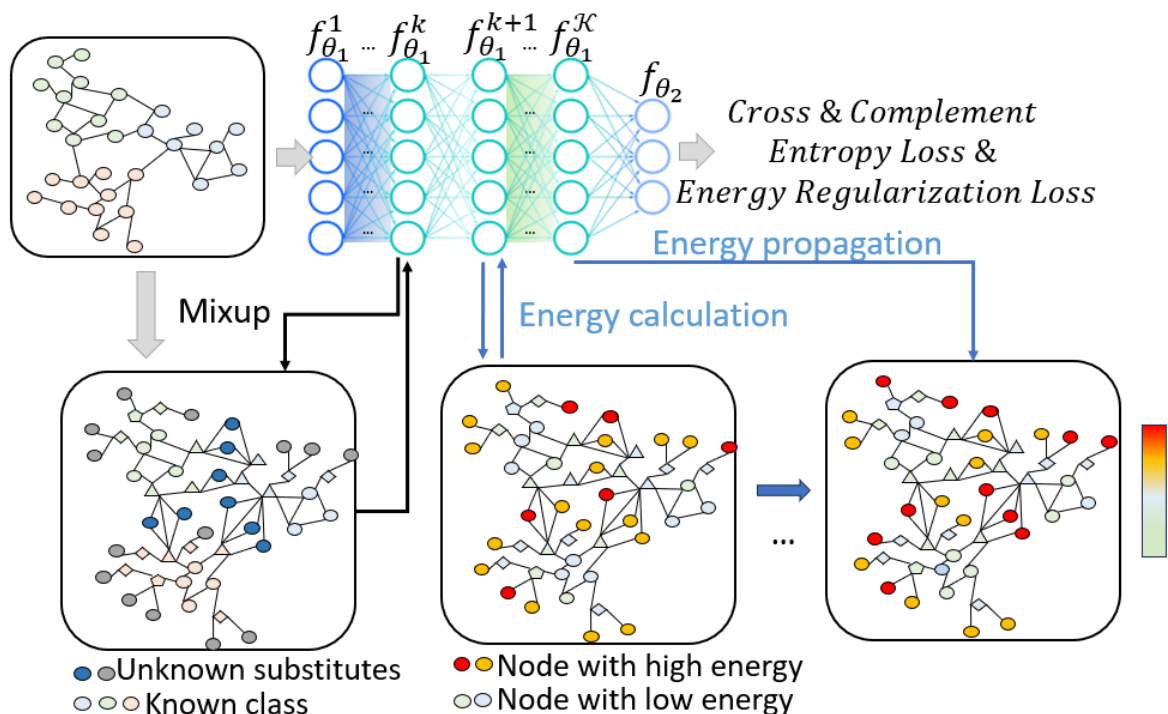
Introduction

■ Contribution

- A novel method, *EGonc*, for open-set node classification is proposed by redefining the open-world graph learning paradigm based on the energy model and elaborate unknown-substitute generation.
- *EGonc* has nice theoretical properties that guarantee an overall distinguishable margin between the detection scores for IND and OOD samples.
- No open-set data (samples of unknown classes or any side information of unknown classes) is required during training and validation.
- *EGonc* is agnostic to specific GNN architecture and demonstrates robust generalization capabilities.

Model

■ Overview of *EGonc* model.



Our model is mainly consists of three components:

- **Substitute Unknowns Generation**
 - An effective way for generating Substitute unknown nodes
 - ✓ Inter-Class Unknown substitute
 - ✓ External Unknown substitute
- **Energy Propagation**
 - An bridge between the energy function and an open-set classifier
- **Open-Set Classifier Learning**
 - An learning module to guarantee the classification of known classes and the rejection of the unknown class.

Mode I

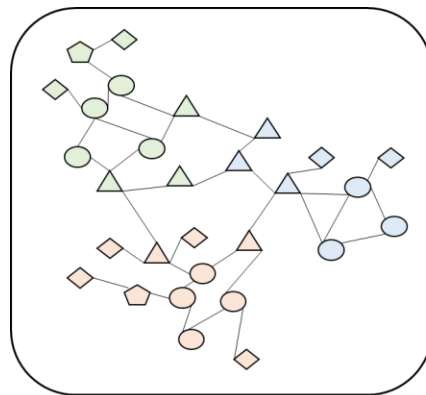
■ Inter-Class Unknown substitutes

$$\begin{cases} \tilde{x}_i = \alpha h_i^k(\theta_1; x_i, A) + (1 - \alpha)h_j^k(\theta_1; x_j, A) \\ \tilde{y}_i = C + 1 \end{cases}$$

■ External Unknown Substitutes

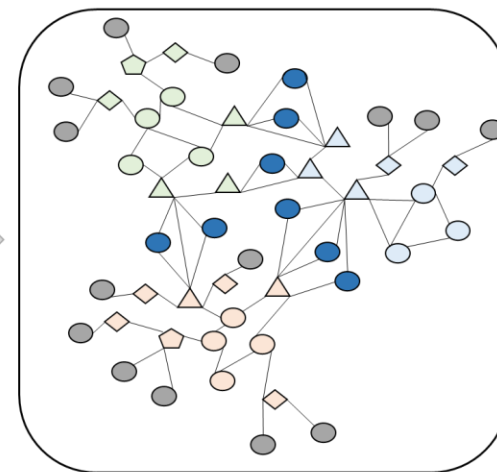
$$h_{(c)}^k = \frac{1}{|X^c|} \sum_{x_i \in X^c} h_i^k(\theta_1; x_i, A), c = 1, \dots, C$$

$$\begin{cases} \tilde{x}_i = \beta h_i^k(\theta_1; x_i, A) + (-\gamma h_{(y_i)}^k) \\ \tilde{y}_i = C + 1 \end{cases}$$



- Known class 1
- Known class 2
- Known class 3

Mixup



- Inter-class unknown proxies
- External unknown proxies

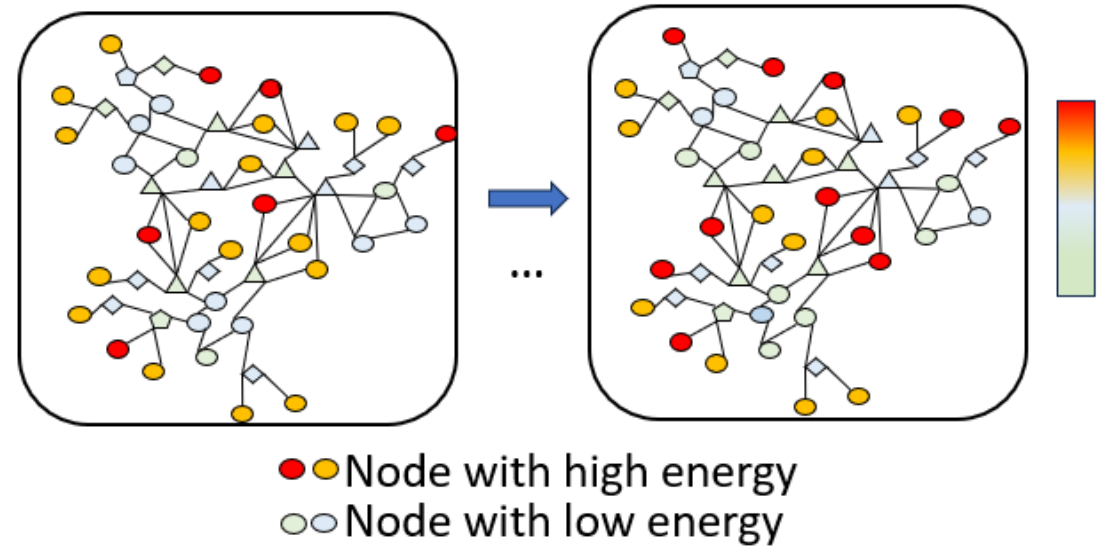
Mode I

■ Energy Propagation

$$E^{(k)} = kE^{(k-1)} + (1 - k)D^{-1}\hat{A}E^{(k-1)}$$

$$\hat{A} = D^{-1/2}AD^{-1/2}$$

$$E^{(k)} = [E_i^{(k)}]$$



Model

Algorithm 1 EGonc: open-set node classification

Require: $G = (V, E, X)$: a graph with links and features;

$\mathcal{D}_{tr} = \{G, Y\}$: train set with labeled nodes;

$X_{te} = S \cup U$: test set where S are the known classes appeared in training and U are the unknown classes;

Ensure: $f(X_{te} \rightarrow \mathcal{Y}), \mathcal{Y} \in \{1, \dots, C, unknown\}$.

1: Obtain the inter-class node pairs $\{(x_i, y_i), (x_j, y_j)\} \in \mathcal{D}_{tr}$ s.t. $y_i \neq y_j \& a_{ij} = 1$

2: Obtain the peripheral nodes that are leaf nodes and low confident nodes.

3: **while** not convergence **do**

4: For the first $m = 1, \dots, k$ layer:

$$h_i^m = f^m(\theta_1; h_i^{m-1}, h_j^{m-1}, j \in \mathcal{N}_i), \forall x_i \in \mathcal{D}_{tr}$$

5: At the k -th layer:

Create unknown substitutes X_{sub} using Eq. (8) & (9)

Augment the substitutes to known class samples:

$$\overline{\mathcal{D}}_{tr} = \mathcal{D}_{tr} \cup (X_{sub}, Y_{C+1})$$

6: For the $m = k + 1, \dots, k_1 - 1$ layers:

$$h_i^m = f^m(\theta_1; h_i^{m-1}, h_j^{m-1}, j \in \mathcal{N}_i), \forall x_i \in \overline{\mathcal{D}}_{tr}$$

7: For the $m = k_1, \dots, K$ layers:

$$h_i^m = f^m(\theta_1; h_i^{m-1}, h_j^{m-1}, j \in \mathcal{N}_i), \forall x_i \in \overline{\mathcal{D}}_{tr}$$

$$E_i^m = f^m(E_i^{m-1}, E_j^{m-1}, j \in \mathcal{N}_i), \forall x_i \in \overline{\mathcal{D}}_{tr}$$

8: For open-set classifier layer:

Obtain cross entropy loss as Eq. (11)

Obtain complement entropy loss as Eq. (12)

Obtain energy regularization loss as Eq. (13)

9: Back-propagate loss gradient using Eq. (14) and update weights

10: **if** early stopping condition is satisfied **then**

11: Terminate

12: **end if**

13: **end while**

$$\checkmark l_1 = \sum_{(x_i, y_i) \in \mathcal{D}_{tr}} l_{CrE}(\hat{y}_i, y_i) + \lambda_1 \sum_{x_i \in X_{sub}} l_{CrE}(\hat{y}_i, C + 1)$$

$$\checkmark l_1 \text{ use } l_{CrE}(\hat{y}_i, y_i) = -y_i \log \hat{y}_i \text{ to maximize data likelihood}$$

$$\checkmark l_2 = \sum_{(x_i, y_i) \in \mathcal{D}_{tr}} l_{CrE}(\hat{y}_i / y_i, C + 1) + \sum_{x_i \in X_{sub}} l_{CoE}(\hat{y}_i, y_i)$$

$$\checkmark l_{CoE} = - \sum_{c=1, c \neq y_i}^{C+1} \frac{\hat{y}_{i,c}}{1 - \hat{y}_{i,y_i}} \log \frac{\hat{y}_{i,c}}{1 - \hat{y}_{i,y_i}}$$

$\checkmark l_2$ use l_{CoE} and to l_{CrE} to eliminate the effects of complement classes

$$\checkmark l_3 = k_1 \left(\sum_{(x_i, y_i) \in \mathcal{D}_{tr}} \sigma(E_{ind}(x_i)) + \sum_{x_j \in X_{sub}} \sigma(E_{ood}(x_j)) \right)$$

$$\checkmark + k_2 \left(\sum_{(x_i, y_i) \in \mathcal{D}_{tr}} \sigma(E_{ind}(x_i))^2 + \sum_{x_j \in X_{sub}} \sigma(E_{ood}(x_j))^2 \right)$$

\checkmark And the loss function is $l_{total} = l_1 + \lambda_2 l_2 + \lambda_3 l_3$

Experiment

Table 5: Statistics of the experimental datasets.

Dataset	Nodes	Edges	Features	Labels
Cora	2708	5429	1433	7
Citeseer	3312	4732	3703	6
DBLP	17716	105734	1639	4
PubMed	19717	44325	500	3
Ogbn_arxiv	169343	1166243	128	40

We select five real-world datasets.

- ✓ Cora
- ✓ Citeseer
- ✓ DBLP
- ✓ PubMed
- ✓ Ogbn_arxiv

Table 1: Near open-set classification on five citation network datasets with one unknown class ($u=1$) in the *inductive learning setting*. Numbers reported are all percentage (%).

Methods	Cora		Citeseer		DBLP		Pubmed		Ogbn_arxiv		Average	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GCN_soft	70.6	67.6	44.6	38.9	63.8	59.2	28.9	29.9	49.8	17.5	51.5	42.6
GCN_sig	69.2	64.7	45.3	44.5	63.5	58.7	28.9	29.8	48.8	9.5	51.1	41.4
GCN_soft_ τ	73.6	73.8	57.3	54.5	65.0	62.4	49.7	48.6	47.3	20.6	58.6	52.0
GCN_sig_ τ	79.7	80.1	62.1	54.6	69.2	68.2	45.1	46.0	46.0	8.3	60.4	51.4
Openmax	74.6	75.1	56.2	54.5	67.2	67.2	49.1	48.7	45.5	16.3	58.5	52.4
DOC	77.8	78.1	66.0	56.7	69.9	69.2	45.6	46.2	46.7	20.7	61.2	52.2
PROSER	83.2	83.7	73.7	63.6	71.7	72.6	71.0	58.4	53.0	<u>31.1</u>	70.5	61.9
OpenWGL	78.1	78.9	64.1	60.8	71.4	72.2	65.3	63.4	45.4	20.7	64.9	60.2
GNNSAFE	79.6	81.0	69.8	60.3	72.5	74.1	70.1	66.8	51.2	24.2	68.6	61.3
\mathcal{G}^2Pxy	<u>84.3</u>	<u>84.8</u>	<u>75.5</u>	<u>71.0</u>	<u>77.3</u>	<u>79.0</u>	<u>73.7</u>	<u>70.2</u>	<u>62.7</u>	33.0	<u>74.7</u>	<u>67.6</u>
EGonc	84.5	84.9	75.8	71.5	79.1	80.8	80.2	75.5	63.0	33.0	76.5	69.1

- ✓ Our method consistently **outperforms** baseline methods for all datasets.
- ✓ Specifically, Our method is **better** than GNNSAFE, g^2pxy and OpenWGL in the inductive learning setting, which are the state-of-the-art method.

Experiment

1. As shown in Table 8, our proposed method consistently **outperforms** the baselines in terms of Acc and F1 on different datasets in the transductive learning setting.
2. As shown in Table 3, when compared under far open-set classification setting, our model consistently outperforms them **in all metrics**.

Table 8: Near open-set classification on five citation network datasets with one unknown class ($u=1$) under in the *transductive learning setting*. Numbers reported are all percentage (%).

Methods	Cora		Citeseer		DBLP		Pubmed		Ogbn_arxiv		Average	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GCN_soft	70.8	68.2	44.7	38.9	62.9	57.0	29.2	29.7	50.2	18.4	51.6	42.4
GCN_sig	68.8	64.5	44.6	40.1	63.4	59.2	29.0	29.5	46.8	8.4	50.5	40.3
GCN_soft_ τ	78.1	78.9	67.3	57.0	67.3	67.7	68.9	27.2	49.6	19.0	66.2	50.0
GCN_sig_ τ	78.3	78.5	65.4	55.3	71.4	71.5	69.0	27.2	45.9	7.7	66.0	48.0
Openmax	77.2	76.9	57.5	56.7	69.0	70.6	55.0	52.1	49.2	18.9	61.6	55.0
DOC	77.3	77.9	65.1	55.3	71.7	72.0	68.4	34.2	49.9	19.4	66.5	51.8
PROSER	84.7	83.6	74.3	66.6	75.3	71.6	72.8	60.8	55.0	30.7	72.4	62.7
OpenWGL	83.3	83.5	70.0	65.4	74.3	74.2	71.2	68.0	46.0	20.0	69.0	62.2
GNNSAFE	80.7	81.9	73.1	62.2	74.2	75.8	73.5	69.9	52.8	24.1	70.8	62.8
\mathcal{G}^2Pxy	<u>90.7</u>	<u>89.7</u>	<u>76.3</u>	<u>71.8</u>	<u>77.5</u>	<u>79.5</u>	<u>78.0</u>	<u>73.4</u>	<u>63.7</u>	<u>31.4</u>	<u>77.2</u>	<u>69.2</u>
EGonc	91.2	90.4	77.2	72.9	79.4	80.7	86.5	80.5	63.8	31.6	79.6	71.2

Table 3: Accuracy and macro-F1 for far open-set classification on benchmark datasets. Numbers reported are all percentage (%).

Methods	Co_Ci		Ci_DB		DB_Pub		Average	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GCN_soft	43.0	58.9	38.4	42.5	41.9	53.7	41.1	51.7
GCN_sig	41.6	57.5	36.3	42.1	41.6	45.2	39.8	48.3
GCN_soft_ τ	81.2	77.6	86.2	71.1	85.0	75.6	84.1	74.8
GCN_sig_ τ	69.4	51.8	68.7	48.0	79.8	69.1	72.6	56.3
Openmax	56.2	55.1	69.6	60.3	69.6	58.7	65.1	58.0
DOC	69.4	57.8	75.5	62.3	78.0	70.7	74.3	63.6
PROSER	78.5	79.1	81.5	66.4	78.6	69.0	79.5	71.5
OpenWGL	80.6	76.7	44.6	11.9	84.6	70.7	69.9	53.1
GNNSAFE	79.3	79.9	80.9	65.9	80.0	65.0	80.1	70.3
\mathcal{G}^2Pxy	<u>81.3</u>	<u>80.5</u>	<u>87.5</u>	<u>74.4</u>	<u>86.5</u>	<u>72.3</u>	<u>85.1</u>	<u>75.7</u>
EGonc	81.7	81.0	88.1	75.2	87.2	<u>72.8</u>	85.7	76.3

Experiment

Table 2: Accuracy and macro-F1 scores of EGonc and its variants with respect to different losses and generation strategies.

Components				Cora		Citeseer		DBLP		Pubmed		O_arxiv		Average	
l_1	l_2	l_3		Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
✓				84.2	84.7	75.2	69.0	76.5	77.7	70.1	47.3	61.9	34.1	73.6	62.6
✓	✓			84.3	84.8	75.5	71.0	77.3	79.0	73.7	70.2	62.7	33.0	74.7	67.6
✓	✓	✓		84.5	84.9	75.8	71.5	79.1	80.8	80.2	75.5	63.0	33.0	76.5	69.1
X_{far}	X_{rand}	X_{int}	X_{ext}	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
				82.7	83.2	73.5	69.6	69.5	71.3	70.4	67.2	60.1	30.0	71.2	64.3
✓				83.7	84.0	75.5	66.9	72.3	72.7	71.8	68.5	62.3	29.3	73.1	64.3
	✓			81.3	82.2	74.6	63.7	71.2	71.5	70.0	66.9	61.9	32.3	71.8	63.3
		✓		84.2	84.7	75.3	70.8	75.3	76.9	73.4	68.7	62.3	31.4	74.1	66.5
			✓	84.1	84.6	75.4	70.9	75.5	74.8	71.4	66.9	61.5	29.5	73.6	65.3
✓	✓			84.0	84.4	75.7	71.2	72.0	71.7	73.0	69.1	61.9	32.0	73.3	65.7
		✓	✓	84.5	84.9	75.8	71.5	79.1	80.8	80.2	75.5	63.0	33.0	76.5	69.1

- As shown in Table 2, we compare variants of *EGonc* with respect to the generative strategy and different losses to demonstrate its effect.

Experiment

Table 4: Accuracy and macro-F1 scores of open-set classification methods with different backbone neural network. Numbers reported are all percentage (%).

Methods	Cora		Citeseer		Dblp		PubMed		Ogbn_arxiv	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
GCN_soft_ τ	73.6	73.8	57.3	54.5	65.0	62.4	49.7	48.6	47.3	20.6
GCN_DOC	77.8	78.1	66.0	56.7	69.9	69.2	45.6	46.2	46.7	<u>20.7</u>
GCN_Openmax	74.6	75.1	56.2	54.5	67.2	67.2	49.1	48.7	45.5	16.3
GCN_ \mathcal{G}^2Pxy	<u>84.3</u>	<u>84.8</u>	<u>75.5</u>	<u>71.0</u>	<u>77.3</u>	<u>79.0</u>	<u>73.7</u>	<u>70.2</u>	<u>62.7</u>	33.0
GCN_EGonc	84.5	84.9	75.8	71.5	79.1	80.8	80.2	75.5	63.0	33.0
GAT_soft_ τ	71.6	69.2	58.9	51.1	65.4	66.6	43.2	43.7	49.1	16.7
GAT_DOC	71.1	72.6	62.4	59.5	64.2	61.8	42.1	42.9	48.3	16.2
GAT_Openmax	66.3	63.4	48.6	48.9	62.5	56.9	48.6	47.0	32.2	8.4
GAT_ \mathcal{G}^2Pxy	<u>80.4</u>	<u>81.0</u>	<u>75.2</u>	<u>70.9</u>	<u>72.9</u>	<u>73.7</u>	<u>71.7</u>	<u>47.0</u>	<u>53.7</u>	<u>22.6</u>
GAT_EGonc	80.8	81.3	75.3	71.0	73.1	74.0	74.3	63.6	56.1	24.5
Graphsage_soft_ τ	72.7	72.9	63.5	51.2	64.3	64.0	46.6	46.9	51.5	16.0
Graphsage_DOC	76.0	<u>75.4</u>	63.6	59.9	68.9	72.2	44.6	45.7	49.5	14.7
Graphsage_Openmax	71.1	70.6	47.9	48.7	62.3	56.9	44.4	45.1	43.2	8.0
Graphsage_ \mathcal{G}^2Pxy	<u>87.2</u>	87.3	<u>78.6</u>	<u>76.9</u>	<u>74.4</u>	<u>74.7</u>	<u>72.8</u>	<u>64.9</u>	<u>62.8</u>	<u>36.5</u>
Graphsage_EGonc	87.3	87.3	79.5	77.4	78.0	79.6	73.0	65.0	63.4	38.4

1. As shown in Table 4, the proposed model *EGonc* is **agnostic to specific GNN architecture** and demonstrates robust generalization capabilities.

Conclusion

- In this paper, we propose a novel energy-based generative open-set node classification method, *EGonc*, by estimating the underlying density of the training data to decide whether a given input is close to the IND data.
- Two kinds of substitute unknowns are generated to mimic the distribution of real open-set samples.
- Under constraint of cross entropy loss, complement entropy loss, and energy regularization loss, *EGonc* achieves superior effectiveness for unknown class detection and known class classification, which is validated by experiments on benchmark graph datasets.
- Moreover, *EGonc* also has good generalization since it has no specific requirement on the GNN architecture.

***EGonc: Energy – based Open – Set Node Classification
with substitute Unknowns***



Thank you!

2110276101@email.szu.edu.cn