



# SE-GSL: A General and Effective Graph Structure Learning Framework through Structural Entropy Optimization

Dongcheng Zou  
Beihang University  
Beijing, China  
zoudongcheng@buaa.edu.cn

Hao Peng\*  
Beihang University  
Beijing, China  
penghao@buaa.edu.cn

Xiang Huang  
Beihang University  
Beijing, China  
huang.xiang@buaa.edu.cn

Renyu Yang  
Beihang University  
Beijing, China  
renyu.yang@buaa.edu.cn

Jianxin Li  
Beihang University  
Beijing, China  
lijx@buaa.edu.cn

Jia Wu  
Macquarie University  
Sydney, Australia  
jia.wu@mq.edu.au

Chunyang Liu  
Didi Chuxing  
Beijing, China  
liuchunyang@didiglobal.com

Philip S. Yu  
University of Illinois Chicago  
Chicago, USA  
psyu@uic.edu

WWW 2023

code: <https://github.com/RingBDStack/SE-GSL>

Reported by Zicong Dou

- i) robustness to system noises and heterophily graphs
- ii) model interpretability.

This paper proposes a general GSL framework, **SE-GSL**, through **structural entropy** and the graph hierarchy abstracted in the **encoding tree**.

## Graph and Community

$G = \{V, E, X\}$  denote a graph

$X \in \mathbb{R}^{n \times d}$  refers to the vertex attribute set.

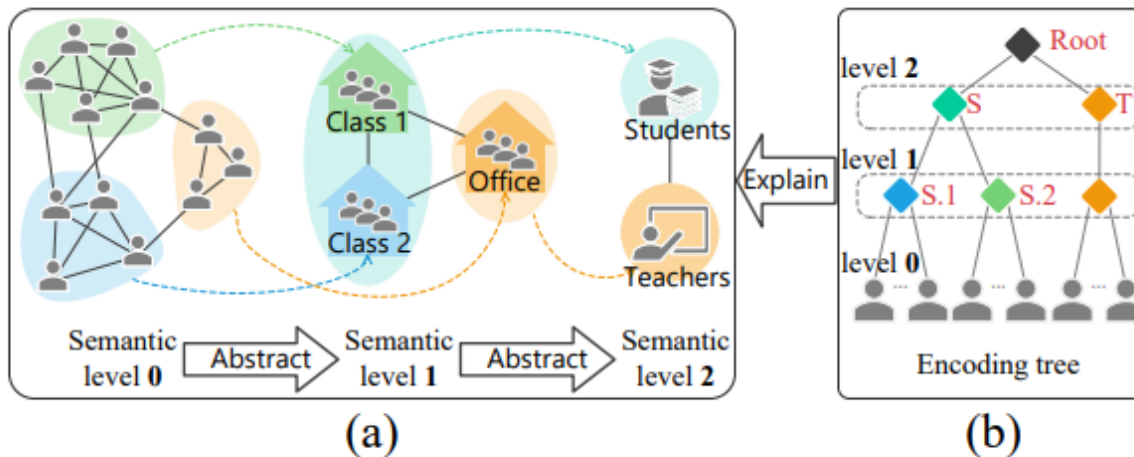
$A \in \mathbb{R}^{n \times n}$       $d(v_i) = \sum_j A_{ij}$

$D = \text{diag}(d(v_1), d(v_2), \dots, d(v_n))$

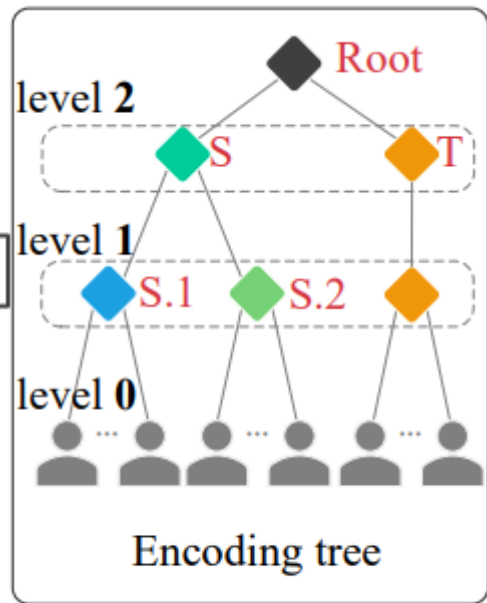
$\mathcal{P} = \{P_1, P_2, \dots, P_L\}$  is a partition of  $V$

## Graph Structure Learning (GSL)

$$\mathcal{L}_{gsl} = \mathcal{L}_{task}(Z, Y) + \alpha \mathcal{L}_{reg}(Z, G^*, G)$$



**Figure 1: An illustrative example of the hierarchical community (semantics) in a simple social network. (1) Vertices and edges represent the people and their interconnectivity (e.g., common locations, interests, occupations). There are different abstraction levels, and each community can be divided into sub-communities in a finer-grained manner (e.g., students are placed in different classrooms while teachers are allocated different offices). The lowest abstraction will come down to the individuals with own attributes, and the highest abstraction is the social network system. (b) An encoding tree is a natural form to represent and interpret such a multi-level hierarchy.**



(b)

*K*-level encoding tree.

## One-dimensional Structural Entropy

$$H^1(G) = - \sum_{v \in V} \frac{d_v}{\text{vol}(G)} \log_2 \frac{d_v}{\text{vol}(G)}, \quad (1)$$

where  $d_v$  is the degree of vertex  $v$ , and  $\text{vol}(G)$  is the sum of the degrees of all vertices in  $G$ .

## Encoding Tree

the encoding tree  $\mathcal{T}$  of graph  $G = (V, E)$

The root node  $\lambda$  in  $\mathcal{T}$  has a label  $T_\lambda = V$

Each non-root node  $\alpha$  has a label  $T_\alpha \subset V$ .

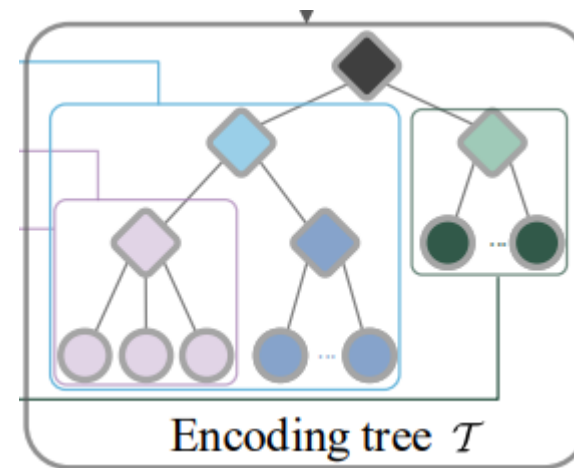
non-root node  $\alpha$ 's parent node in  $\mathcal{T}$  is denoted as  $\alpha^-$

its  $i$ -th children node  $\alpha^{(i)}$

non-leaf node  $\alpha$  the number of children  $\alpha$  is  $N$

all vertex subset  $T_{\alpha^{(i)}}$  form a partition of  $T_\alpha$

$$T_\alpha = \bigcup_{i=1}^N T_{\alpha^{(i)}} \text{ and } \bigcap_{i=1}^N T_{\alpha^{(i)}} = \emptyset$$



## High-dimensional Structural Entropy

$$H^K(G) = \min_{\forall \mathcal{T}: \text{height}(\mathcal{T}) \leq K} \{H^\mathcal{T}(G)\}, \quad (2)$$

$$H^\mathcal{T}(G) = \sum_{\alpha \in \mathcal{T}, \alpha \neq \lambda} H^\mathcal{T}(G; \alpha) = - \sum_{\alpha \in \mathcal{T}, \alpha \neq \lambda} \frac{g_\alpha}{\text{vol}(G)} \log_2 \frac{\mathcal{V}_\alpha}{\mathcal{V}_{\alpha^-}}, \quad (3)$$

where  $g_\alpha$  is the sum weights of the cut edge set  $[T_\alpha, T_\alpha/T_\lambda]$ , i.e., all edges connecting vertices inside  $T_\alpha$  with vertices outside  $T_\alpha$ .  $\mathcal{V}_\alpha$  is the sum of degrees of all vertices in  $T_\alpha$ .  $H^\mathcal{T}(G; \alpha)$  is the structural entropy of node  $\alpha$  and  $H^\mathcal{T}(G)$  is the structural entropy of  $\mathcal{T}$ .  $H^K(G)$  is the  $K$ -dimensional structural entropy, with the optimal encoding tree of  $K$ -level.

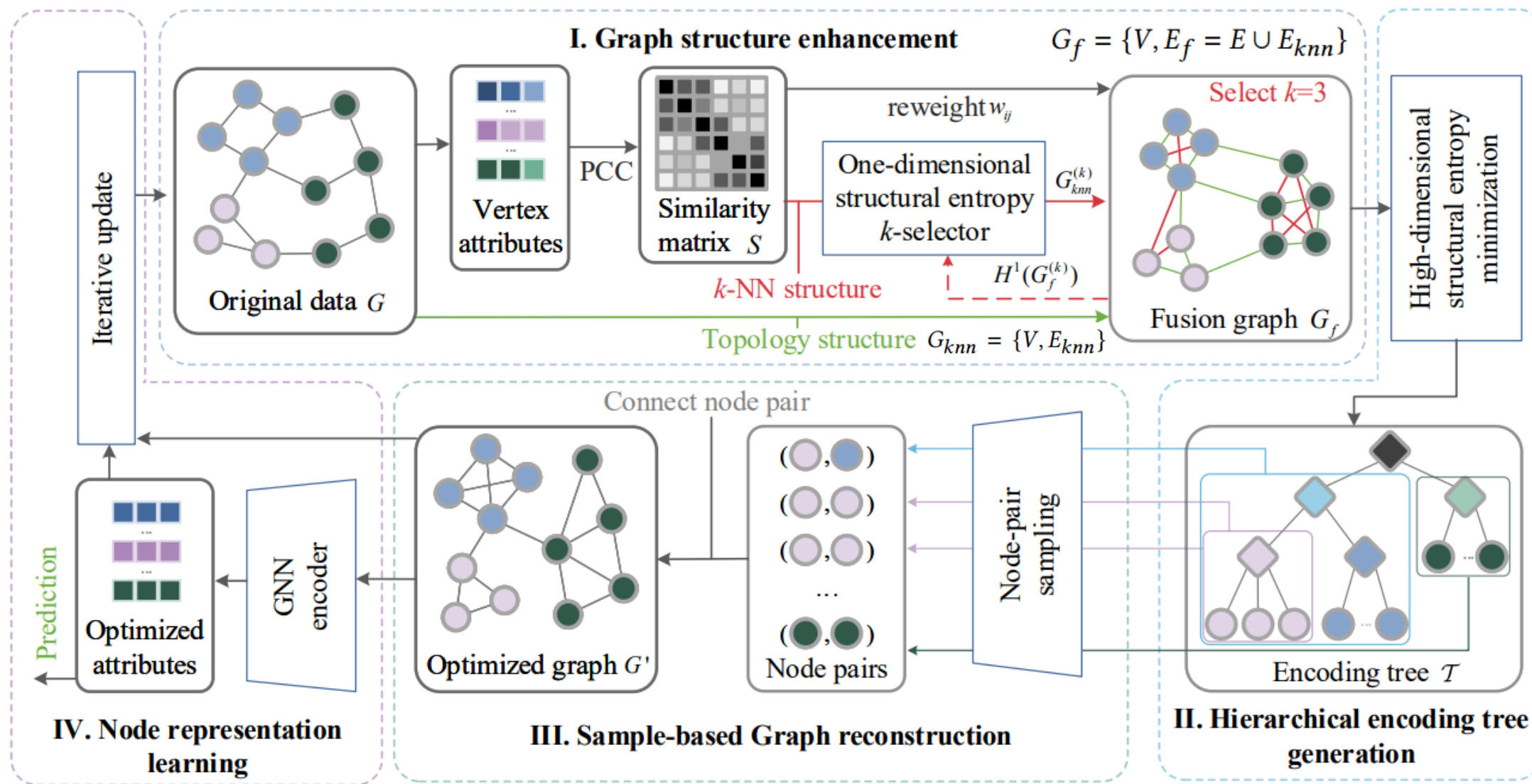
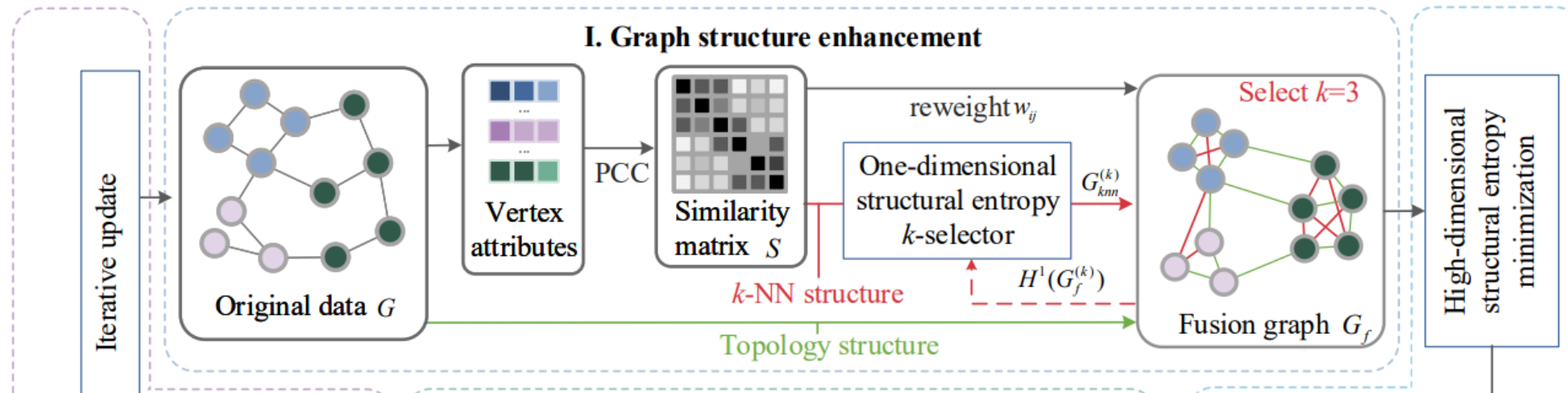


Figure 2: The overall architecture of SE-GSL.



## Graph Structure Enhancement

$$S_{ij} = \text{PCC}(x_i, x_j) = \frac{E((x_i - u_i)(x_j - u_j))}{\sigma_i \sigma_j}, \quad (4)$$

where  $x_i$  and  $x_j \in \mathbb{R}^{1 \times d}$  are the attribute vectors of vertices  $i$  and  $j$ , respectively.  $u_i$  and  $\sigma_i$  denote the mean value and variance of  $x_i$ , and  $E(\cdot)$  is the dot product function. Based on  $S$ , we can intrinsically construct the  $k$ -NN graph  $G_{knn} = \{V, E_{knn}\}$  where each edge in  $E_{knn}$  represents a vertex and its  $k$  nearest neighbors (e.g., the edges in red in Fig 2). We fuse  $G_{knn}$  and the original  $G$  to  $G_f = \{V, E_f = E \cup E_{knn}\}$ .

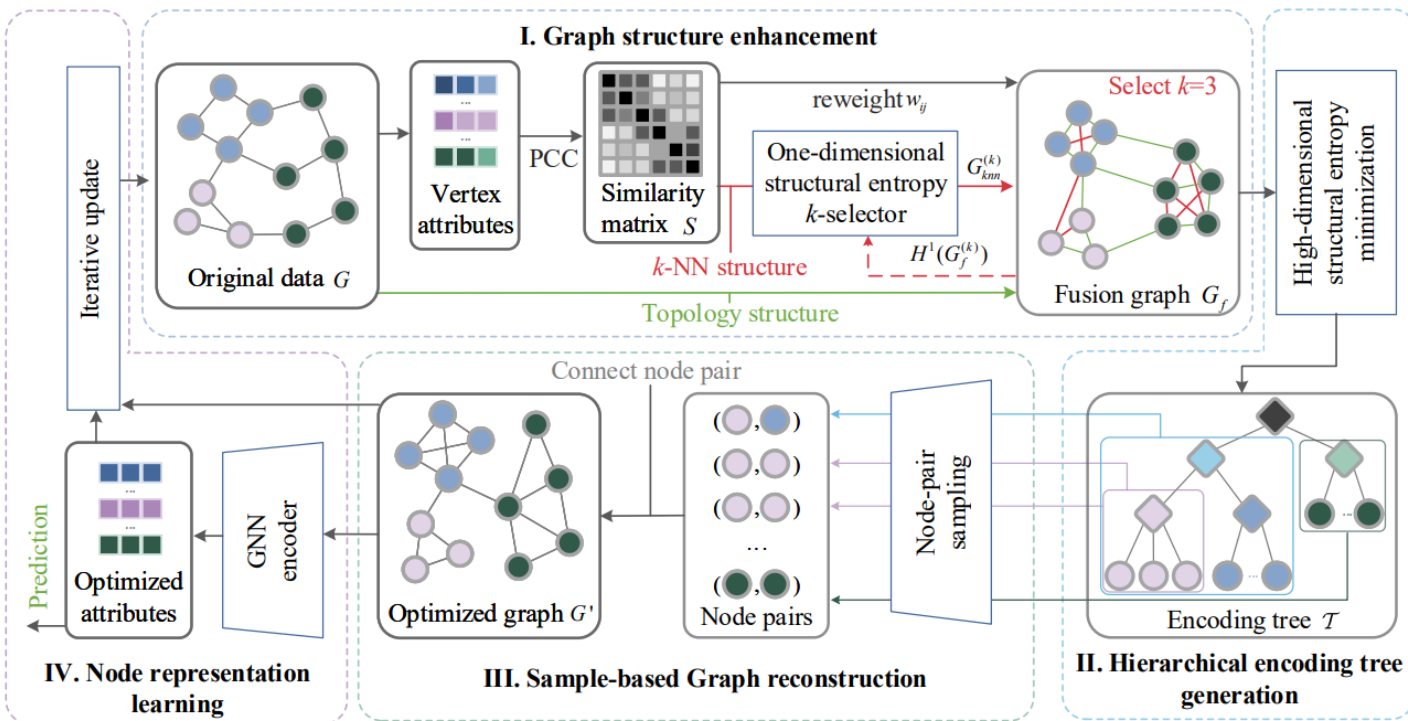
maximize the one-dimensional structural entropy  $H^1(G_f)$  to guide the selection of  $k$  for larger encoding information. In practice, we gradually increase the integer parameter  $k$ , generate the corresponding  $G_f^{(k)}$  and compute  $H^1(G_f^{(k)})$ .

$k$  reaches a threshold  $k_m$

the edge  $e_{ij}$  between  $v_i$  and  $v_j$  is reweighted as:

$$\omega_{ij} = S_{ij} + M, \quad M = \frac{1}{2|V|} \cdot \frac{1}{|E|} \sum_{1 < i, j < n} S_{ij}, \quad (5)$$

where  $M$  is a modification factor that amplifies the trivial edge weights and thus makes the  $k$ -selector more sensitive to noises.

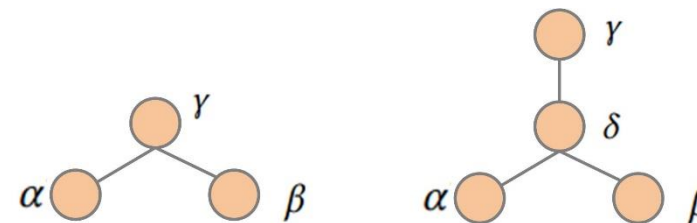


## Hierarchical Encoding Tree Generation

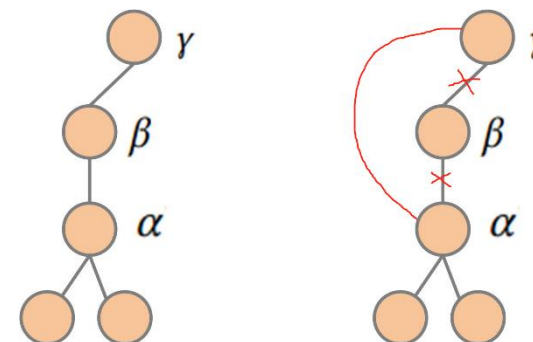
$$\mathcal{T}^* = \arg \min_{\forall \mathcal{T}: \text{height}(\mathcal{T}) \leq K} (H^{\mathcal{T}}(G)). \quad (6)$$

To address this optimization problem, we design a greedy-based heuristic algorithm to approximate  $H^K(G)$ . To assist the greedy heuristic, we define two basic operators:

**Definition 1. Combining operator:** Given an encoding tree  $\mathcal{T}$  for  $G = (V, E)$ , let  $\alpha$  and  $\beta$  be two nodes in  $\mathcal{T}$  sharing the same parent  $\gamma$ . The combining operator  $\text{CB}_{\mathcal{T}}(\alpha, \beta)$  updates the encoding tree as:  $\gamma \leftarrow \delta^-; \delta \leftarrow \alpha^-; \delta \leftarrow \beta^-$ . A new node  $\delta$  is inserted between  $\gamma$  and its children  $\alpha, \beta$ .



**Definition 2. Lifting operator:** Given an encoding tree  $\mathcal{T}$  for  $G = (V, E)$ , let  $\alpha, \beta$  and  $\gamma$  be the nodes in  $\mathcal{T}$ , satisfying  $\beta^- = \gamma$  and  $\alpha^- = \beta$ . The lifting operator  $\text{LF}_{\mathcal{T}}(\alpha, \beta)$  updates the encoding tree as:  $\gamma \leftarrow \alpha^-$ ; IF  $:T_{\beta} = \emptyset$ , THEN  $:\text{drop}(\beta)$ . The subtree rooted at  $\alpha$  is lifted by placing itself as  $\gamma$ 's child. If no more children exist after lifting,  $\beta$  will be deleted from  $\mathcal{T}$ .



## Sample-based Graph Reconstruction

**Definition 3. Structural entropy of deduction:** Let  $\mathcal{T}$  be an encoding tree of  $G$ . We define the structural entropy of the deduction from non-leaf node  $\lambda$  to its descendant  $\alpha$  as:

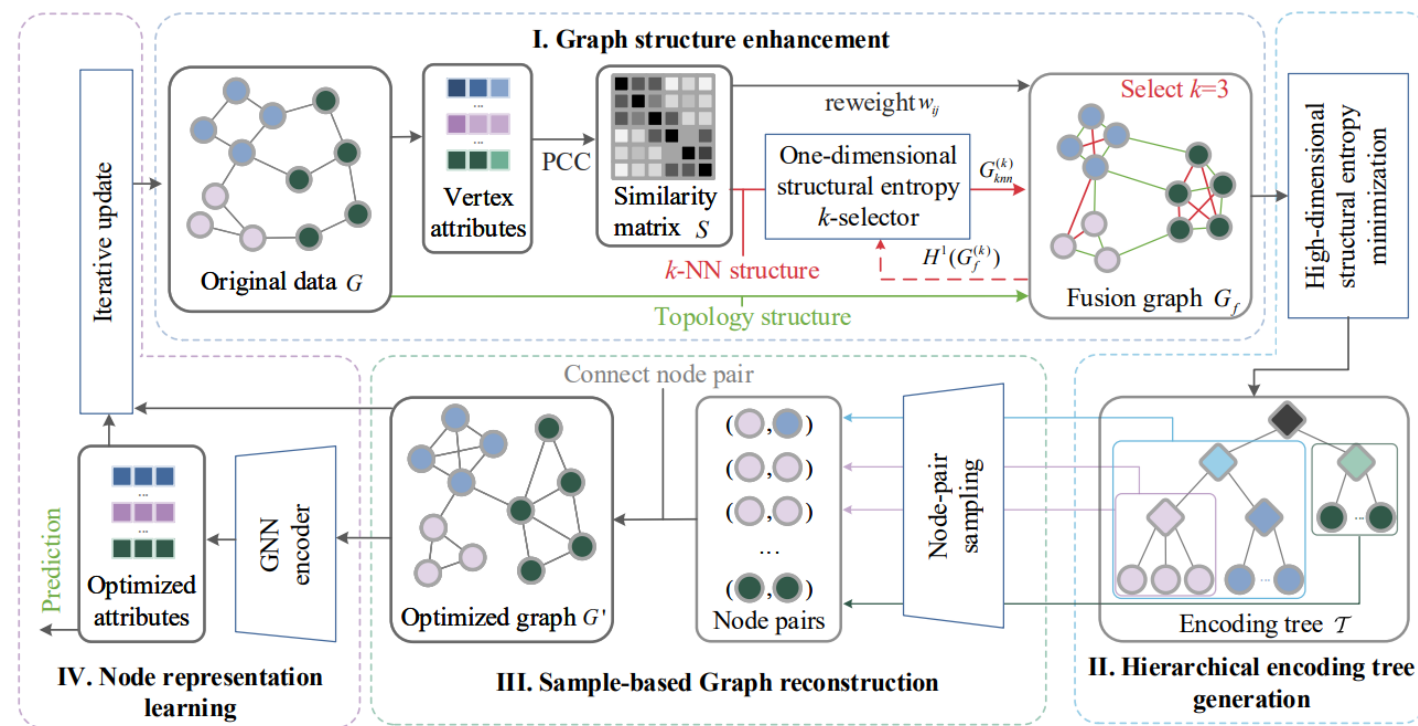
$$H^{\mathcal{T}}(G; (\lambda, \alpha]) = \sum_{\beta, T_{\alpha} \subseteq T_{\beta} \subset T_{\lambda}} H^{\mathcal{T}}(G; \beta). \quad (7)$$

Specifically, for a given  $\mathcal{T}$ , assume the node  $\delta$  has a set of child nodes  $\{\delta^{(1)}, \delta^{(2)}, \dots, \delta^{(n)}\}$ . The probability of the child  $\delta^{(i)}$  is defined as:  $P(\delta^{(i)}) = \sigma_{\delta}(H^{\mathcal{T}}(G_f; (\lambda, \delta^{(i)}]))$ , where  $\lambda$  is the root of  $\mathcal{T}$  and  $\sigma_{\delta}(\cdot)$  represents a distribution function. Take softmax function as an example, the probability of  $\delta^{(i)}$  can be calculated as:

$$P(\delta^{(i)}) = \frac{\exp(H^{\mathcal{T}}(G_f; (\lambda, \delta^{(i)}]))}{\sum_{j=1}^n \exp(H^{\mathcal{T}}(G_f; (\lambda, \delta^{(j)}]))}. \quad (8)$$

(3) After recursively performing step (2), we sample two leaf nodes  $\delta_1$  and  $\delta_2$ , while adding the edge connecting vertex  $v_1 = T_{\delta_1}$  and  $v_2 = T_{\delta_2}$  into the edge set  $E'$  of graph  $G'$ .

For each subtree rooted at  $\delta$ , we conduct independent samplings for  $\theta \times n$  times, where  $n$  is the number of  $\delta$ 's children,  $\theta$  is a hyperparameter that positively correlated with the density of reconstructed graph.



(1) For the encoding tree (or subtree) with root node  $\delta$ , two different child nodes  $\delta^{(i)}$  and  $\delta^{(j)}$  are selected by sampling according to  $P(\delta^{(i)})$  and  $P(\delta^{(j)})$ . Let  $\delta_1 \leftarrow \delta^{(i)}$  and  $\delta_2 \leftarrow \delta^{(j)}$

(2) If  $\delta_1$  is a non-leaf node, we perform sampling once on the subtree rooted at  $\delta_1$  to get  $\delta_1^{(i)}$ , then update  $\delta_1 \leftarrow \delta_1^{(i)}$ . The same is operated on  $\delta_2$ .

**Table 1: Classification Accuracy (%) comparison, with improvement range of SE-GSL against the baselines.** The best results are bolded and the second-best are underlined. Green denotes the outperformance percentage, while yellow denotes underperformance.

Method	Cora	Citeseer	Pubmed	PT	TW	Actor	Cornell	Texas	Wisconsin
GCN	87.26 $\pm$ 0.63	76.22 $\pm$ 0.71	87.46 $\pm$ 0.12	67.62 $\pm$ 0.21	62.46 $\pm$ 1.94	27.65 $\pm$ 0.55	49.19 $\pm$ 1.80	57.30 $\pm$ 2.86	48.57 $\pm$ 4.08
GAT	87.52 $\pm$ 0.69	76.04 $\pm$ 0.78	86.61 $\pm$ 0.15	68.76 $\pm$ 1.01	61.68 $\pm$ 1.20	27.77 $\pm$ 0.59	57.09 $\pm$ 6.32	58.10 $\pm$ 4.14	51.34 $\pm$ 4.78
GCNII	87.57 $\pm$ 0.87	75.47 $\pm$ 1.01	<u>88.64<math>\pm</math>0.23</u>	68.93 $\pm$ 0.93	65.17 $\pm$ 0.47	30.66 $\pm$ 0.66	58.76 $\pm$ 7.11	55.36 $\pm$ 6.45	51.96 $\pm$ 4.36
Grand	<b>87.93<math>\pm</math>0.71</b>	77.59 $\pm$ 0.85	86.14 $\pm$ 0.98	69.80 $\pm$ 0.75	<u>66.79<math>\pm</math>0.22</u>	29.80 $\pm$ 0.60	57.21 $\pm$ 2.48	56.56 $\pm$ 1.53	52.94 $\pm$ 3.36
Mixhop	85.71 $\pm$ 0.85	75.94 $\pm$ 1.00	87.31 $\pm$ 0.44	69.48 $\pm$ 0.30	66.34 $\pm$ 0.22	33.72 $\pm$ 0.76	64.47 $\pm$ 4.78	63.16 $\pm$ 6.28	72.12 $\pm$ 3.34
Droptedge	86.32 $\pm$ 1.09	76.12 $\pm$ 1.32	87.58 $\pm$ 0.34	68.49 $\pm$ 0.91	65.24 $\pm$ 1.45	30.10 $\pm$ 0.71	58.94 $\pm$ 5.95	59.20 $\pm$ 5.43	60.45 $\pm$ 4.48
Geom-GCN-P	84.93	75.14	88.09	-	-	31.63	60.81	67.57	64.12
Geom-GCN-S	85.27	74.71	84.75	-	-	30.30	55.68	59.73	56.67
GDC	87.17 $\pm$ 0.36	76.13 $\pm$ 0.53	88.08 $\pm$ 0.16	66.14 $\pm$ 0.54	64.14 $\pm$ 1.40	28.74 $\pm$ 0.76	59.46 $\pm$ 4.35	56.42 $\pm$ 3.99	48.30 $\pm$ 4.29
GEN	<u>87.84<math>\pm</math>0.69</u>	<b>78.77<math>\pm</math>0.88</b>	86.13 $\pm$ 0.41	<u>71.62<math>\pm</math>0.78</u>	65.16 $\pm$ 0.77	<b>36.69<math>\pm</math>1.02</b>	65.57 $\pm$ 6.74	73.38 $\pm$ 6.65	54.90 $\pm$ 4.73
H <sub>2</sub> GCN-2	87.81 $\pm$ 1.35	76.88 $\pm$ 1.77	<b>89.59<math>\pm</math>0.33</b>	68.15 $\pm$ 0.30	63.33 $\pm$ 0.77	35.62 $\pm$ 1.30	<b>82.16<math>\pm</math>6.00</b>	<u>82.16<math>\pm</math>5.28</u>	<u>85.88<math>\pm</math>4.22</u>
SE-GSL	<b>87.93<math>\pm</math>1.24</b>	<u>77.63<math>\pm</math>1.65</u>	88.16 $\pm$ 0.76	<b>71.91<math>\pm</math>0.66</b>	<b>66.99<math>\pm</math>0.93</b>	<u>36.34<math>\pm</math>2.07</u>	<u>75.21<math>\pm</math>5.54</u>	<b>82.49<math>\pm</math>4.80</b>	<b>86.27<math>\pm</math>4.32</b>
Improvement	0.00~3.00	-1.14~2.92	-1.43~3.41	0.29~5.77	0.20~5.31	-0.35~8.69	-6.95~26.02	0.33~27.13	0.39~37.97





**Table 2: Classification accuracy(%) of SE-GSL and corresponding backbones. Wisc. is short for Wisconsin.**

Method	Actor	TW	Texas	Wisc.	Improvement
SE-GSL <sub>GCN</sub>	35.03	66.88	75.68	79.61	↑ 5.20~31.04
SE-GSL <sub>SAGE</sub>	36.20	66.92	<b>82.49</b>	<b>86.27</b>	↑ 0.25~6.79
SE-GSL <sub>GAT</sub>	32.46	63.57	74.59	78.82	↑ 4.69~27.48
SE-GSL <sub>APPNP</sub>	<b>36.34</b>	<b>66.99</b>	81.28	83.14	↑ 2.01~12.16

**Table 3: The  $k$  selection for each iteration in structural optimization. Bolds represent the  $k$  selection when the accuracy reaches maximum.**

Iteration	1	2	3	4	5	6	7	8	9
Cora	22	<b>22</b>	19	22	21	22	20	21	20
Actor	23	15	15	15	14	15	14	<b>14</b>	15
TW	50	16	16	<b>17</b>	15	17	27	16	16
Wisconsin	21	16	<b>11</b>	16	14	13	16	13	11
Texas	21	13	13	<b>13</b>	13	10	14	10	14

**Table 4: Glossary of Notations.**

Notation	Description
$G; A; S$	Graph; Adjacency matrix; Similarity matrix.
$v; e; x$	Vertex; Edge; Vertex attribute.
$V; E; X$	Vertex set; Edge set; Attribute set.
$ V ;  E $	The number of vertices and edges.
$\mathcal{P}; P_i$	The partition of $V$ ; A community.
$D; d(v_i)$	The degree matrix; The degree of vertex $v_i$ .
$e_{ij}$	The edge between $v_i$ and $v_j$ .
$w_{ij}$	The weight of edge $e_{ij}$ .
$vol(G)$	The volume of graph $G$ , i.e., degree sum in $G$ .
$G_{knn}^{(k)}$	The $k$ -NN graph with parameter $k$ .
$G_f$	Fusion graph.
$G_f^{(k)}$	The fusion graph with parameter $k$ .
$\mathcal{T}$	Encoding tree.
$\mathcal{T}^*$	The optimal encoding tree.
$\lambda$	The root node of an encoding tree.
$\alpha$	A non-root node of an encoding tree.
$\alpha^-$	The parent node of $\alpha$ .
$\alpha^{(i)}$	the $i$ -th child of $\alpha$ .
$T_\lambda$	The label of $\lambda$ , i.e., node set $V$ .
$T_\alpha$	The label of $\alpha$ , i.e., a subset of $V$ .
$\mathcal{V}_\alpha$	Volume of graph $G$ .
$g_a$	the sum weights of cut edge set $[T_\alpha, T_\alpha/T_\lambda]$ .
$N(\mathcal{T})$	The number of non-root node in $\mathcal{T}$ .
$H^{\mathcal{T}}(G)$	Structural entropy of $G$ under $\mathcal{T}$ .
$H^K(G)$	$K$ -dimensional structural entropy.
$H^1(G)$	One-dimensional structural entropy.
$H^{\mathcal{T}}(G; \alpha)$	Structural entropy of node $\alpha$ in $\mathcal{T}$ .
$H^{\mathcal{T}}(G; (\lambda, \alpha])$	Structural entropy of a deduction from $\lambda$ to $\alpha$ .

**Table 5: Statistics of benchmark datasets.**

Dataset	Nodes	Edges	Classes	Features	homophily
Cora	2708	5429	7	1433	0.83
Citeseer	3327	4732	6	3703	0.71
Pubmed	19717	44338	3	500	0.79
PT	1912	31299	2	3169	0.59
TW	2772	63462	2	3169	0.55
Actor	7600	33544	5	931	0.24
Cornell	183	295	5	1703	0.30
Texas	183	309	5	1703	0.11
Wisconsin	251	499	5	1703	0.21

**Table 6: Comparison of training time(hr.) of achieving the best performance based on GPU.**

Method	Cora	Citeseer	Pubmed	PT	TW	Actor	Cornell	Texas	Wisconsin
SE-GSL <sub>GCN</sub>	0.071	0.213	4.574	0.178	0.374	1.482	0.006	0.008	0.009
SE-GSL <sub>SAGE</sub>	0.074	0.076	4.852	0.169	0.214	0.817	0.006	0.007	0.009
SE-GSL <sub>GAT</sub>	0.071	0.180	4.602	0.172	0.329	1.273	0.006	0.008	0.009
SE-GSL <sub>APPNP</sub>	0.073	0.215	4.854	0.138	0.379	1.367	0.010	0.011	0.013

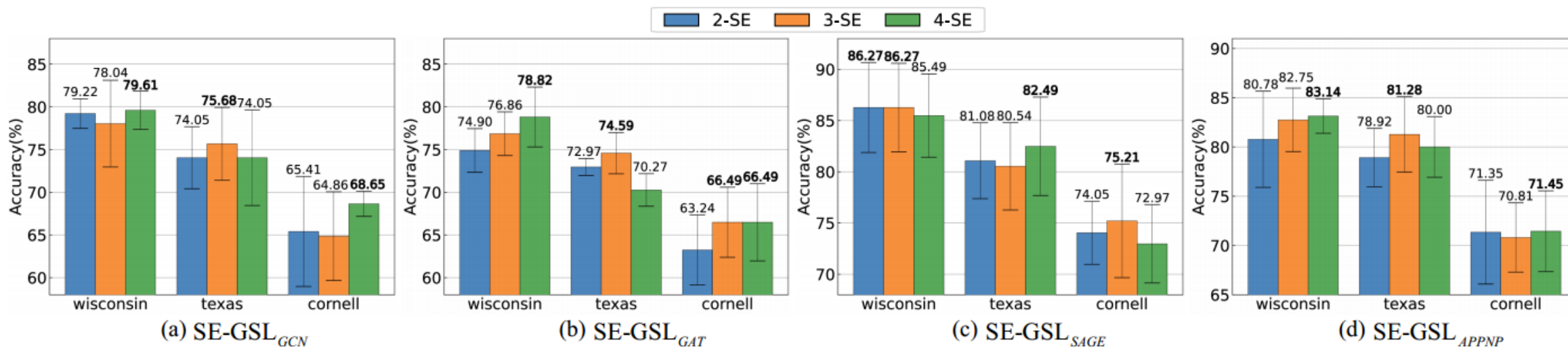
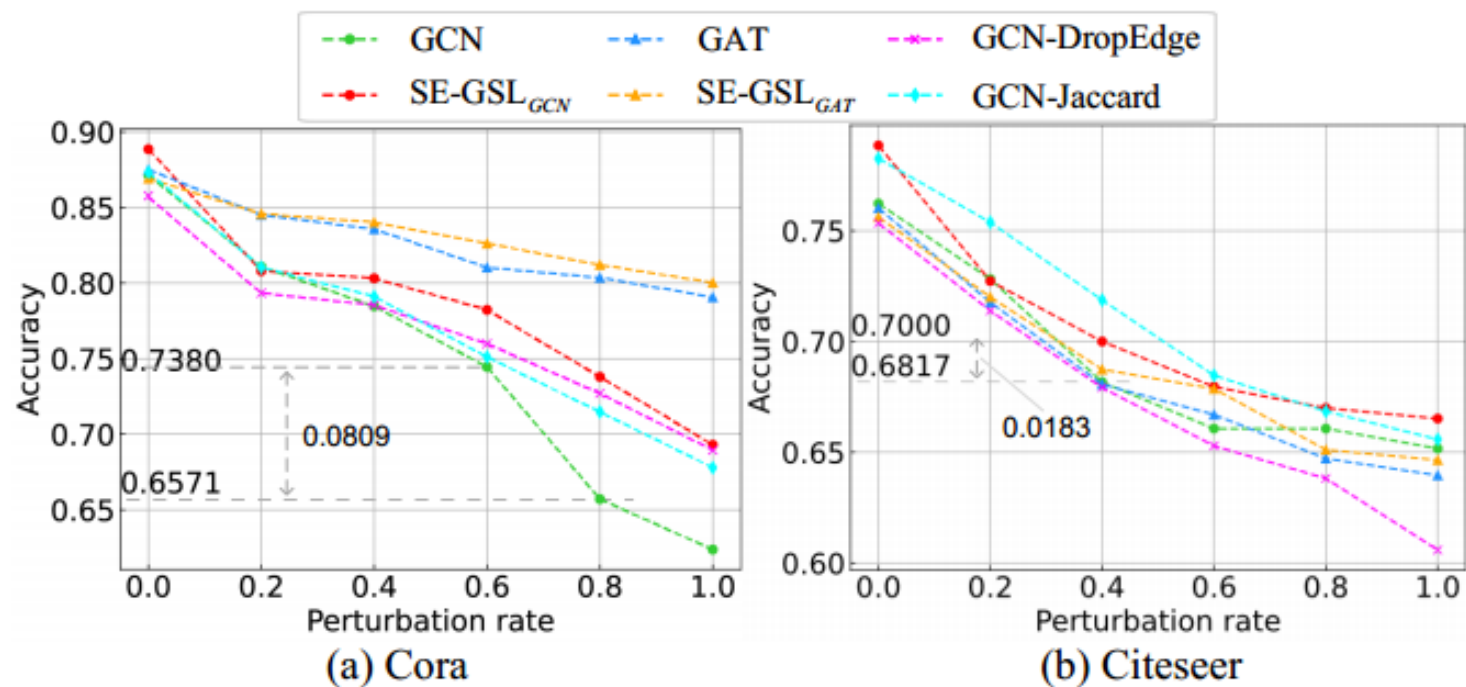
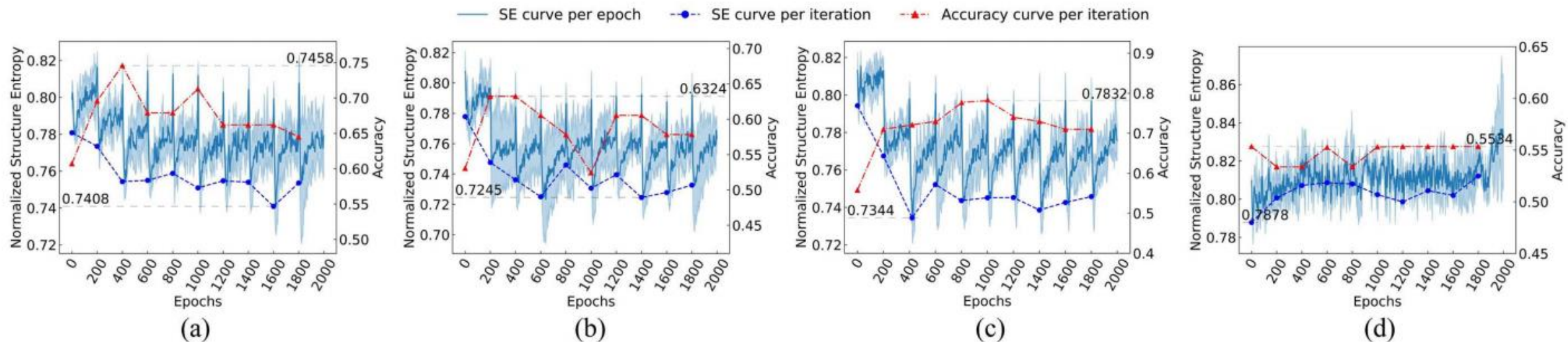


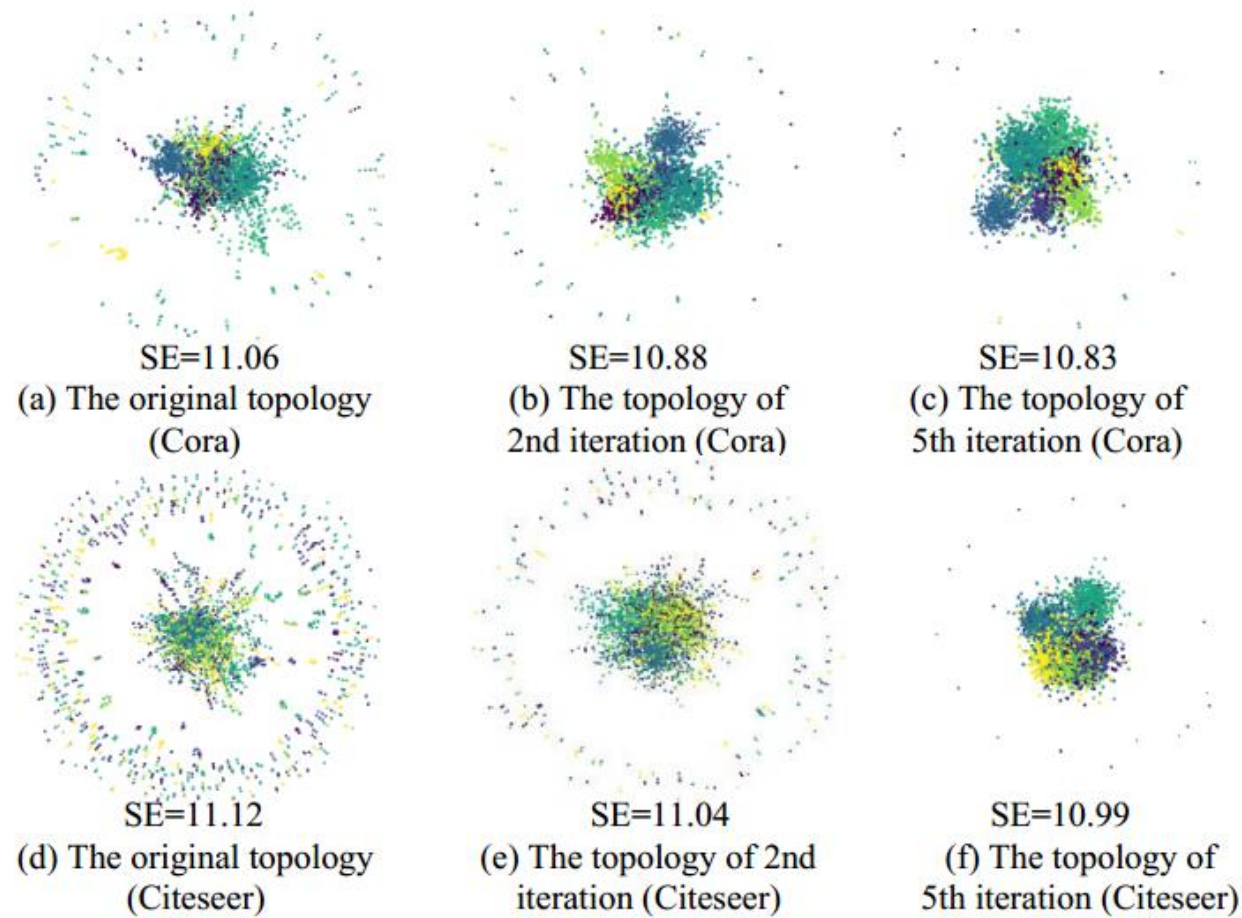
Figure 3: Results of SE-GSL with different encoding tree heights.



**Figure 4: Robustness of SE-GSL against random noises.**



**Figure 5: The normalized structural entropy changes during the training of SE-GSL<sub>GAT</sub> with 2-dimensional structural entropy on (a) Texas, (b) Cornell, and (c) Wisconsin. The structure is iterated every 200 epochs. By comparison, (d) shows the entropy changes on Wisconsin without the graph reconstruction strategy.**



**Figure 6: The visualized evolution of the graph structure on Cora (a,b,c) and Citeseer (d,e,f). The corresponding Structural Entropy (SE) is also shown.**



# Thanks