

Rethinking Tokenized Graph Transformers for Node Classification

Jinsong Chen, Chenyang Li, Gaichao Li, John E. Hopcroft, Kun He*



Hopcroft Center on Computing Science,
School of Computer Science and Technology,
Huazhong University of Science and Technology, China
(*Corresponding author, brooklet60@hust.edu.cn)

Nov. 2025

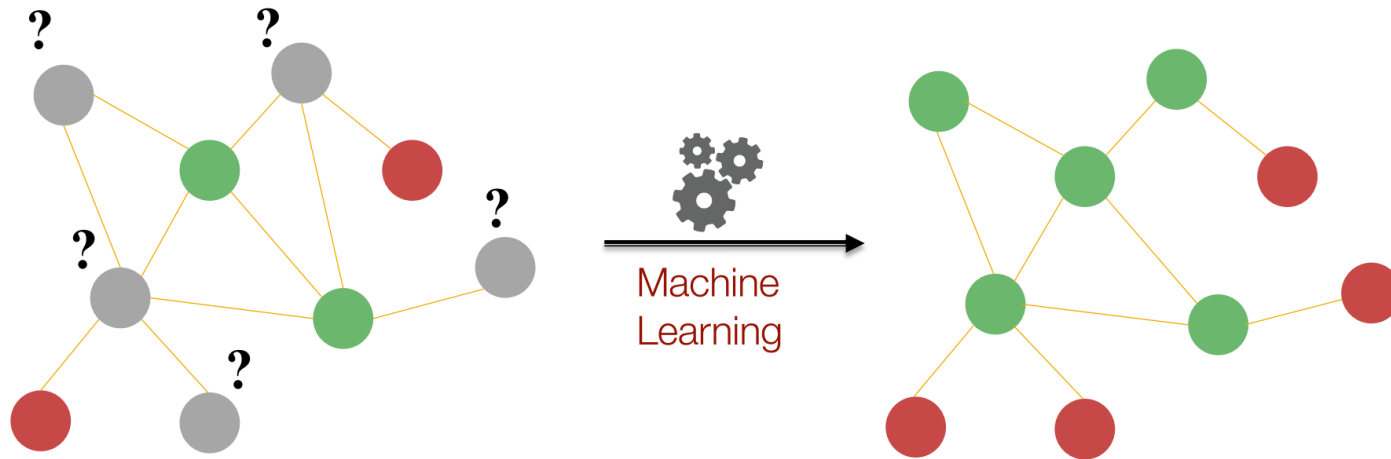
➤ Background



- **Node Classification**

An attributed graph $G = (V, E)$, the adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, the feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, and the label matrix $\mathbf{Y} \in \mathbb{R}^{n \times c}$.

Given a labeled node set V_L , predict the labels of other nodes in $V - V_L$.



➤ Background



- **Graph Transformers for node classification**

Leveraging the **Transformer layer** to learn the node representations.

Two main categories of existing GTs:

- **Entire graph-based GTs:**

Requiring the entire graph as the model input. Performing attention calculation on all node pairs.

Involving many irrelevant nodes and introducing high training cost.

- **Tokenized GTs:**

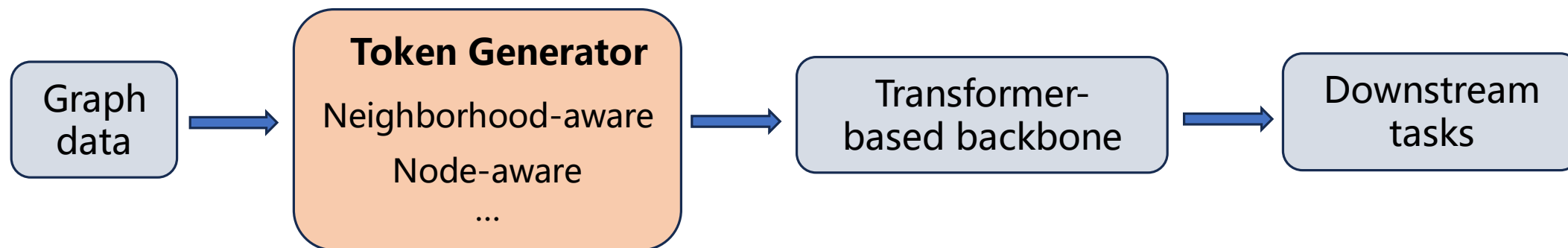
Transforming the input graph into token sequences for feeding Transformer to learn node representations.

Focusing on necessary graph information carried by tokens and requiring low training cost.

➤ Background



- **Tokenized Graph Transformers**



Neighborhood and node are two important elements in existing token generator.

Compared to neighborhood-aware tokens, node-aware tokens are more flexible to preserve various graph information.

➤ Background



- **Node-aware token generator**

- **Step 1: Measuring the similarity of nodes.**

Develop a function, such as cosine similarity and random walk-based strategies to calculate the similarity of each node pair.

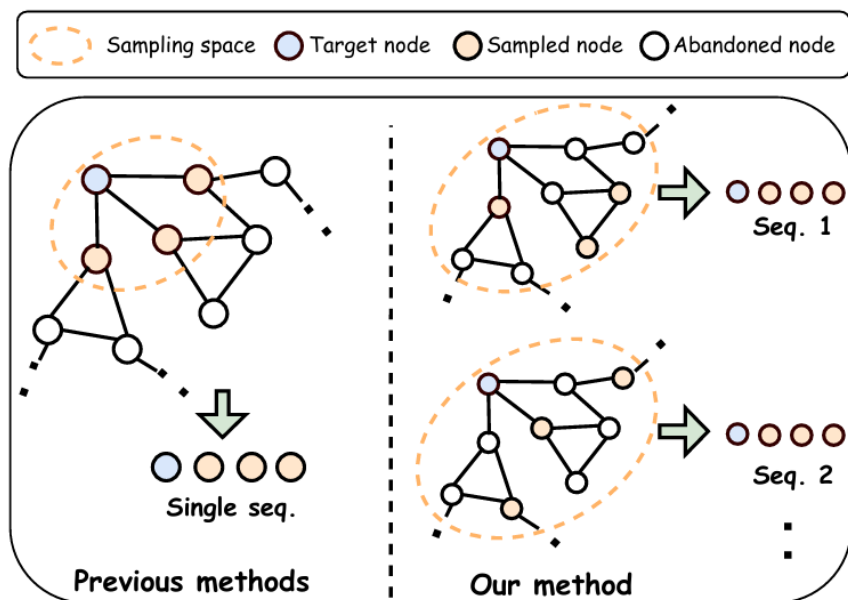
- **Step 2: Node sampling**

Apply top- k sampling strategy to sampling nodes with high similarity as tokens to construct the token sequence.

> Motivation



- **Rethinking Tokenized Graph Transformers**



Using the Top-K sampling is equivalent to selecting only first-order neighbors on the constructed K-NN graph.

How to fully utilize neighbors beyond first-order on constructed K-NN graph?

Figure 1: The toy example of token generation on the constructed k -NN graph. Previous methods only focus on 1-hop neighborhood to construct a single token sequence. While our method can flexibly select tokens from multi-hop neighborhoods to generate diverse token sequences.

- **Key idea**

Using **token swapping** to enlarge sampling space and construct multiple token sequences for model training.

- **Main steps**

- Token Swapping:

Swap semantically relevant tokens in different token sequences to generate new sequences.

- Multi-sequence learning:

Learn node representations from multiple token sequences by Transformer.

- Center Alignment Loss:

Introduce auxiliary alignment loss function for constraining model training.

> SwapGT

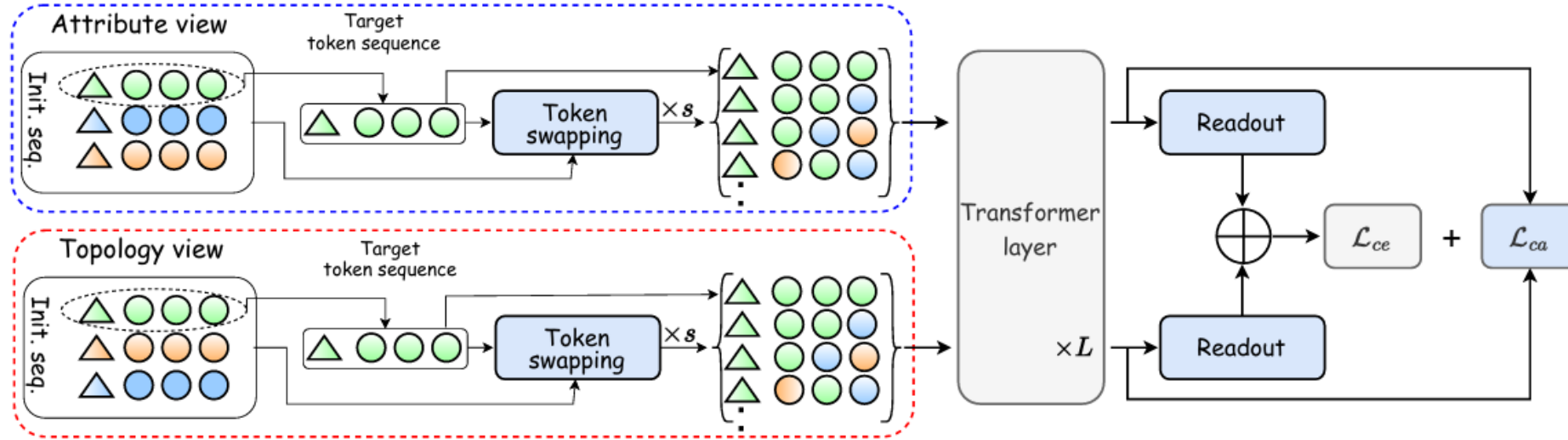


Figure 2: The overall framework of SwapGT. First, we generate the initial token sequences from both the attribute view and topology view. Then, we utilize the proposed token swapping operation to generate new token sequences for each target node. These generated token sequences are then fed into a Transformer-based backbone to learn node representations and generate predicted labels. Additionally, a center alignment loss is adopted to further constrain the representations extracted from different token sequences.

- **Token Sampling**

- Calculating node similarity matrix $\mathbf{S} \in \mathbb{R}^{n \times n}$

$$\mathbf{S} = \frac{\mathbf{X}^{in} \cdot \mathbf{X}^{inT}}{|\mathbf{X}^{in}| |\mathbf{X}^{in}|}, \mathbf{X}^{in} \in \mathbb{R}^{n \times d} \text{ represents the arbitrary node features}$$

$\mathbf{X}^{in} = \mathbf{X}$ for attribute feature view, $\mathbf{X}^{in} = \hat{\mathbf{A}}^k \mathbf{X}$ for topology feature view

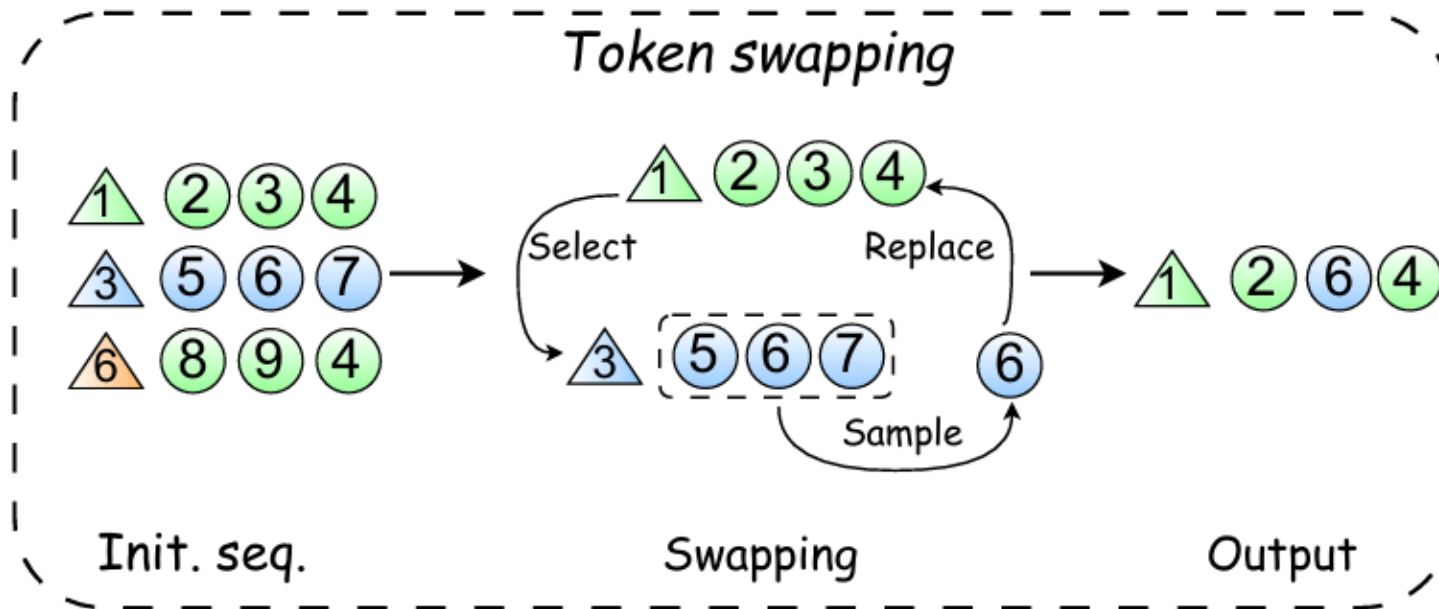
- Sampling initial token set

$$N_i = \{v_j | v_j \in \text{Top}(\mathbf{S}_i, k)\}$$

- **Token Swapping**

Randomly select tokens from node v_i 's token set N_i , and perform swapping operation to form a new token set N'_i .

$$N'_i = \{\delta(N_j) | v_j \in N_i\}$$



Swapping **once** can extend sampling space on **K-NN graph** from 1-order to **2-order**.

- **Token Swapping**

Algorithm 1 The Token Swapping Algorithm

Input: Sampled token set of all nodes $N \in \mathbb{R}^{n \times k}$; Target node v_i ; Probability p ; Swapping times t

Output: The new token set $N'_i \in \mathbb{R}^{1 \times k}$ of v_i

```
1: Initialize  $N'_i = N_i$ ;  
2: for  $t_0 = 1$  to  $t$  do  
3:   Initialize  $N^{new} = \{\}$ ;  
4:   for  $v_j \in N'_i$  do  
5:     if  $random(0, 1) > p$  then  
6:        $N^{new} = N^{new} \cup \{v_j\}$ ;  
7:     else  
8:        $v^{new} = \zeta(N_j)$ ;  
9:        $N^{new} = N^{new} \cup \{v^{new}\}$ ;  
10:    end if  
11:  end for  
12:   $N'_i = N^{new}$ ;  
13: end for  
14: return  $N'_i$ 
```

Totally, swap t times to expand the sampling space to **(t+1)** order.

- **Token sequence construction**

For node v_i , one token set N_i can generate one token sequence \mathbf{Z}_i^{in} as:

$$\mathbf{Z}_i^{in} = [\mathbf{X}_i, \mathbf{X}_{N_{i,0}}, \dots, \mathbf{X}_{N_{i,k-1}}]$$

By performing token swapping Algorithm for s times, we have $1+s$ token sets for v_i . We combine all token sequences to obtain the mode input $\mathbf{Z}_i \in \mathbb{R}^{(1+s) \times (1+k) \times d}$

For attribute and topology, we have \mathbf{Z}_i^A and \mathbf{Z}_i^T respectively.

- **Multi-sequence learning**

- Transformer-based backbone:

$$\mathbf{z}_i'^{A,(l)} = \text{MSA} \left(\mathbf{z}_i^{A,(l-1)} \right) + \mathbf{z}_i^{A,(l-1)}$$

$$\mathbf{z}_i^{A,(l)} = \text{FFN} \left(\mathbf{z}_i'^{A,(l)} \right) + \mathbf{z}_i'^{A,(l)}$$

- Readout function based on multi token sequences:

$$\mathbf{z}_i^{A,F} = \mathbf{z}_0^{A,i} \parallel \left(\frac{1}{s} \sum_{j=1}^s \mathbf{z}_j^{A,i} \right)$$

- **Multi-sequence learning**

- Fusing node representations in attribute and topology view:

$$\mathbf{Z}_i^F = \alpha \cdot \mathbf{Z}_i^{A,F} + (1 - \alpha) \cdot \mathbf{Z}_i^{T,F}$$

- Predicting labels of nodes:

$$\hat{\mathbf{Y}}_i = \text{MLP}(\mathbf{Z}_i^F)$$

$$\mathcal{L}_{ce} = - \sum_{i \in V_L} \mathbf{Y}_i \ln \hat{\mathbf{Y}}_i$$

- **Multi-sequence Center Alignment Loss**

- Center representation of multi sequences:

$$\mathbf{z}_c^i = \frac{1}{s+1} \sum_{j=0}^s \mathbf{z}_j^i$$

- Center alignment loss:

$$\text{CAL}(\mathbf{z}^i, \mathbf{z}_c^i) = 1 - \frac{1}{s+1} \sum_{j=0}^s \text{Cosine}(\mathbf{z}_j^i, \mathbf{z}_c^i)$$

- Overall loss function:

$$\mathcal{L}_{ca} = \text{CAL}(\mathbf{z}^{A,i}, \mathbf{z}_c^{A,i}) + \text{CAL}(\mathbf{z}^{T,i}, \mathbf{z}_c^{T,i})$$

$$\mathcal{L} = \mathcal{L}_{ce} + \lambda \cdot \mathcal{L}_{ca}$$

➤ Experiments



- **Datasets**

Table 3: Statistics of datasets, ranked by the homophily level.

| Dataset | # nodes | # edges | # features | # labels | $\mathcal{H} \downarrow$ |
|-------------|---------|---------|------------|----------|--------------------------|
| Photo | 7,650 | 238,163 | 745 | 8 | 0.83 |
| ACM | 3,025 | 1,3128 | 1,870 | 3 | 0.82 |
| Computer | 13,752 | 491,722 | 767 | 10 | 0.78 |
| Citeseer | 3,327 | 4,552 | 3,703 | 6 | 0.74 |
| WikiCS | 11,701 | 216,123 | 300 | 10 | 0.66 |
| BlogCatalog | 5,196 | 171,743 | 8,189 | 6 | 0.40 |
| UAI2010 | 3,067 | 28,311 | 4,973 | 19 | 0.36 |
| Flickr | 7,575 | 239,738 | 12,047 | 9 | 0.24 |

We use two strategies to split datasets:

Dense splitting: 50%/25%/25% and Sparse splitting: 2.5%/2.5%/95%

- Performance comparison (dense splitting)

Table 1: Comparison of all models in terms of mean accuracy \pm stdev (%) under dense splitting. The best results appear in **bold**. The second results appear in underline.

| Dataset \mathcal{H} | Photo 0.83 | ACM 0.82 | Computer 0.78 | Citeseer 0.74 | WikiCS 0.66 | BlogCatalog 0.40 | UAI2010 0.36 | Flickr 0.24 |
|--------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| SGC | 93.74 \pm 0.07 | 93.24 \pm 0.49 | 88.90 \pm 0.11 | 76.81 \pm 0.26 | 76.67 \pm 0.19 | 72.61 \pm 0.07 | 69.87 \pm 0.17 | 47.48 \pm 0.40 |
| APPNP | 94.98 \pm 0.41 | 93.00 \pm 0.55 | <u>91.31\pm0.29</u> | 77.52 \pm 0.22 | 81.96 \pm 0.14 | 94.77 \pm 0.19 | 77.41 \pm 0.47 | 84.66 \pm 0.31 |
| GPRGNN | 94.57 \pm 0.44 | 93.42 \pm 0.20 | 90.15 \pm 0.34 | 77.59 \pm 0.36 | 82.43 \pm 0.29 | 94.36 \pm 0.29 | 76.94 \pm 0.64 | 85.91 \pm 0.51 |
| FAGCN | 94.06 \pm 0.03 | 93.37 \pm 0.24 | 83.17 \pm 1.81 | 76.19 \pm 0.62 | 79.89 \pm 0.93 | 79.92 \pm 4.39 | 72.17 \pm 1.57 | 82.03 \pm 0.40 |
| BM-GCN | 95.10 \pm 0.20 | 93.68 \pm 0.34 | 91.28 \pm 0.96 | 77.91 \pm 0.58 | 83.90 \pm 0.41 | 94.85 \pm 0.42 | 77.39 \pm 1.13 | 83.97 \pm 0.87 |
| ACM-GCN | 94.56 \pm 0.21 | 93.04 \pm 1.28 | 85.19 \pm 2.26 | 77.62 \pm 0.81 | <u>83.95\pm0.41</u> | 94.53 \pm 0.53 | 76.87 \pm 1.42 | 83.85 \pm 0.73 |
| NAGphormer | <u>95.47\pm0.29</u> | 93.32 \pm 0.30 | 90.79 \pm 0.45 | 77.68 \pm 0.73 | 83.61 \pm 0.28 | 94.42 \pm 0.63 | 76.36 \pm 1.12 | 86.85 \pm 0.85 |
| SGFormer | <u>92.93\pm0.12</u> | 93.79 \pm 0.34 | 81.86 \pm 3.82 | 77.86 \pm 0.76 | 79.65 \pm 0.31 | 94.33 \pm 0.19 | 57.98 \pm 3.95 | 61.05 \pm 0.68 |
| Specformer | 95.22 \pm 0.13 | 93.63 \pm 1.94 | 85.47 \pm 1.44 | 77.96 \pm 0.89 | 83.74 \pm 0.62 | 94.21 \pm 0.23 | 73.06 \pm 0.77 | 86.55 \pm 0.40 |
| VCR-Graphormer | 95.38 \pm 0.51 | 93.11 \pm 0.79 | 90.47 \pm 0.58 | 77.21 \pm 0.65 | 80.82 \pm 0.72 | 94.19 \pm 0.17 | 76.08 \pm 0.52 | 85.96 \pm 0.55 |
| CoBFormer | 95.29 \pm 0.35 | 93.82 \pm 0.58 | 90.21 \pm 0.89 | 77.74 \pm 0.72 | 83.28 \pm 0.68 | 93.98 \pm 0.72 | 76.85 \pm 0.69 | 86.84 \pm 0.78 |
| PolyFormer | 95.45 \pm 0.21 | <u>94.27\pm0.44</u> | 90.87 \pm 0.74 | <u>78.03\pm0.86</u> | 83.79 \pm 0.75 | <u>95.08\pm0.43</u> | <u>77.92\pm0.82</u> | <u>87.01\pm0.57</u> |
| SwapGT | 95.92\pm0.18 | 94.98\pm0.41 | 91.73\pm0.72 | 78.49\pm0.95 | 84.52\pm0.63 | 95.93\pm0.56 | 79.06\pm0.73 | 87.56\pm0.61 |

SwapGT outperforms advanced GTs as well as representative GNNs on all datasets.

- Performance comparison (sparse splitting)

Table 2: Comparison of all models in terms of mean accuracy \pm stdev (%) under sparse splitting. The best results appear in **bold**. The second results appear in underline.

| Dataset \mathcal{H} | Photo 0.83 | ACM 0.82 | Computer 0.78 | Citeseer 0.74 | WikiCS 0.66 | BlogCatalog 0.40 | UAI2010 0.36 | Flickr 0.24 |
|--------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|----------------------------------|
| SGC | 91.90 \pm 0.35 | 89.57 \pm 0.28 | 86.79 \pm 0.19 | 66.41 \pm 0.59 | 74.99 \pm 0.19 | 71.23 \pm 0.06 | 51.61 \pm 0.41 | 39.43 \pm 0.50 |
| APNP | <u>92.24\pm0.28</u> | 89.91 \pm 0.89 | <u>87.64\pm0.39</u> | <u>66.70\pm0.11</u> | 77.42 \pm 0.31 | 81.76 \pm 0.38 | <u>61.65\pm0.71</u> | 71.39 \pm 0.62 |
| GPRGNN | 92.13 \pm 0.32 | 89.47 \pm 0.90 | 86.38 \pm 0.44 | 66.50 \pm 0.62 | 77.59 \pm 0.49 | <u>84.57\pm0.35</u> | <u>58.75\pm0.75</u> | <u>71.89\pm0.89</u> |
| FAGCN | 92.02 \pm 0.18 | 88.47 \pm 0.31 | 83.99 \pm 1.95 | 64.54 \pm 0.66 | 75.21 \pm 0.84 | 76.38 \pm 0.82 | 54.67 \pm 0.96 | 63.68 \pm 0.72 |
| BM-GCN | 91.19 \pm 0.39 | <u>90.11\pm0.60</u> | 86.14 \pm 0.51 | 66.11 \pm 0.47 | 77.39 \pm 0.37 | 84.05 \pm 0.54 | 57.51 \pm 1.14 | 60.82 \pm 0.76 |
| ACM-GCN | 91.71 \pm 0.64 | 89.68 \pm 0.45 | 86.64 \pm 0.59 | 64.85 \pm 1.19 | <u>77.68\pm0.57</u> | 77.17 \pm 1.34 | 56.05 \pm 2.11 | 64.58 \pm 1.53 |
| NAGphormer | 91.65 \pm 0.80 | 89.73 \pm 0.48 | 85.31 \pm 0.65 | 63.66 \pm 1.68 | 76.93 \pm 0.75 | 79.19 \pm 0.41 | 58.36 \pm 1.01 | 67.48 \pm 1.04 |
| SGFormer | 90.13 \pm 0.56 | 88.03 \pm 0.60 | 80.07 \pm 0.21 | 62.41 \pm 0.94 | 74.69 \pm 0.52 | 78.15 \pm 0.69 | 50.19 \pm 1.72 | 51.01 \pm 1.05 |
| Specformer | 90.57 \pm 0.55 | 88.20 \pm 1.05 | 85.55 \pm 0.63 | 62.64 \pm 1.54 | 75.24 \pm 0.71 | 79.75 \pm 1.29 | 57.42 \pm 1.06 | 56.94 \pm 1.48 |
| VCR-Graphormer | 91.39 \pm 0.75 | 86.81 \pm 0.84 | 85.06 \pm 0.64 | 57.61 \pm 0.60 | 72.81 \pm 1.44 | 74.90 \pm 1.18 | 56.43 \pm 1.10 | 50.93 \pm 1.12 |
| CoBFormer | 91.21 \pm 0.62 | 89.18 \pm 0.89 | 85.06 \pm 0.79 | 63.82 \pm 1.34 | 76.28 \pm 1.28 | 79.44 \pm 0.98 | 58.23 \pm 0.82 | 66.94 \pm 1.18 |
| PolyFormer | 91.52 \pm 0.78 | 89.83 \pm 0.62 | 85.75 \pm 0.78 | 64.77 \pm 1.27 | 75.12 \pm 1.16 | 81.02 \pm 0.81 | 58.89 \pm 0.77 | 67.85 \pm 1.43 |
| SwapGT | 92.93\pm0.26 | 90.92\pm0.69 | 88.14\pm0.52 | 69.91\pm1.02 | 78.11\pm0.83 | 88.11\pm0.58 | 63.96\pm1.09 | 72.16\pm1.19 |

SwapGT outperforms advanced GTs as well as representative GNNs on all datasets.

- Study on center alignment loss

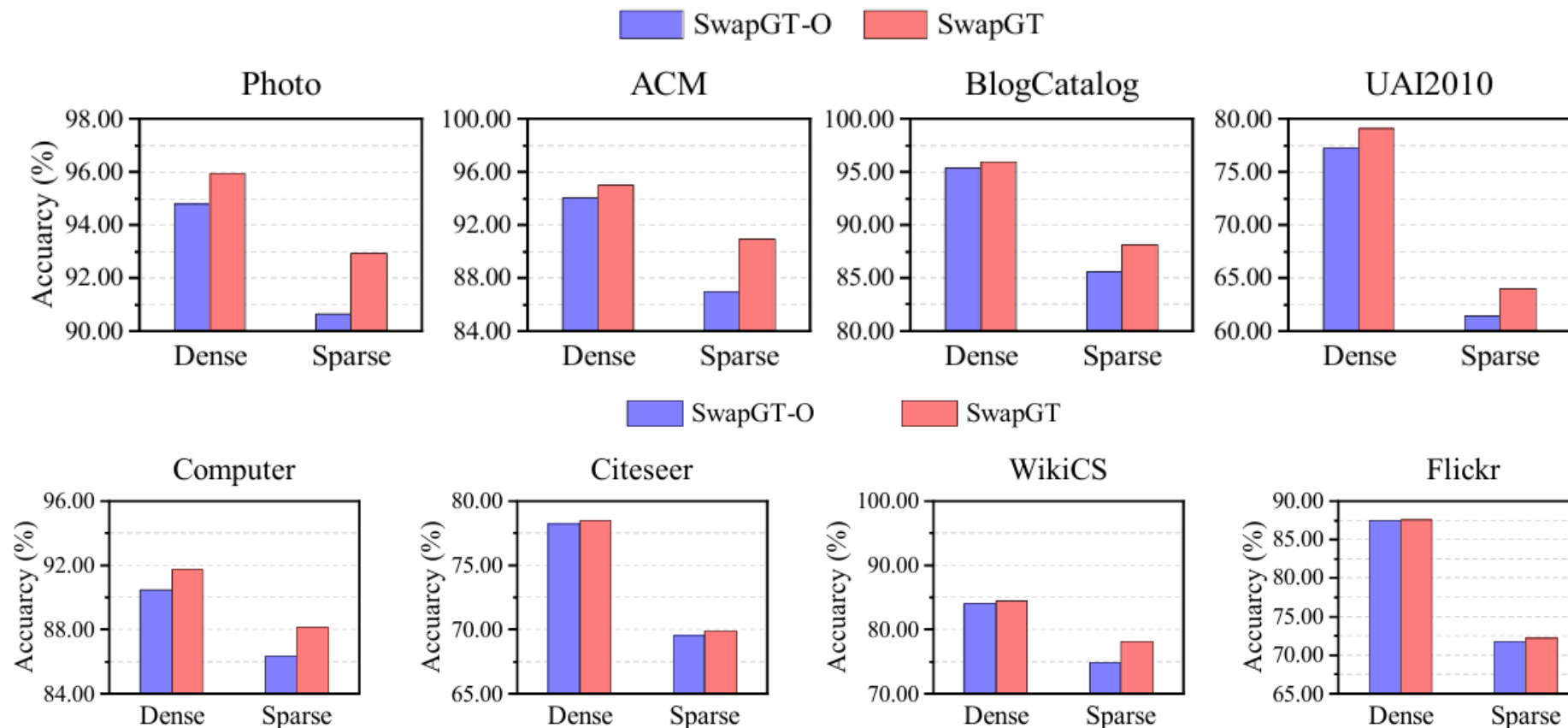


Figure 9: Performances of SwapGT with or without the center alignment loss.

- Study on token sequence generation

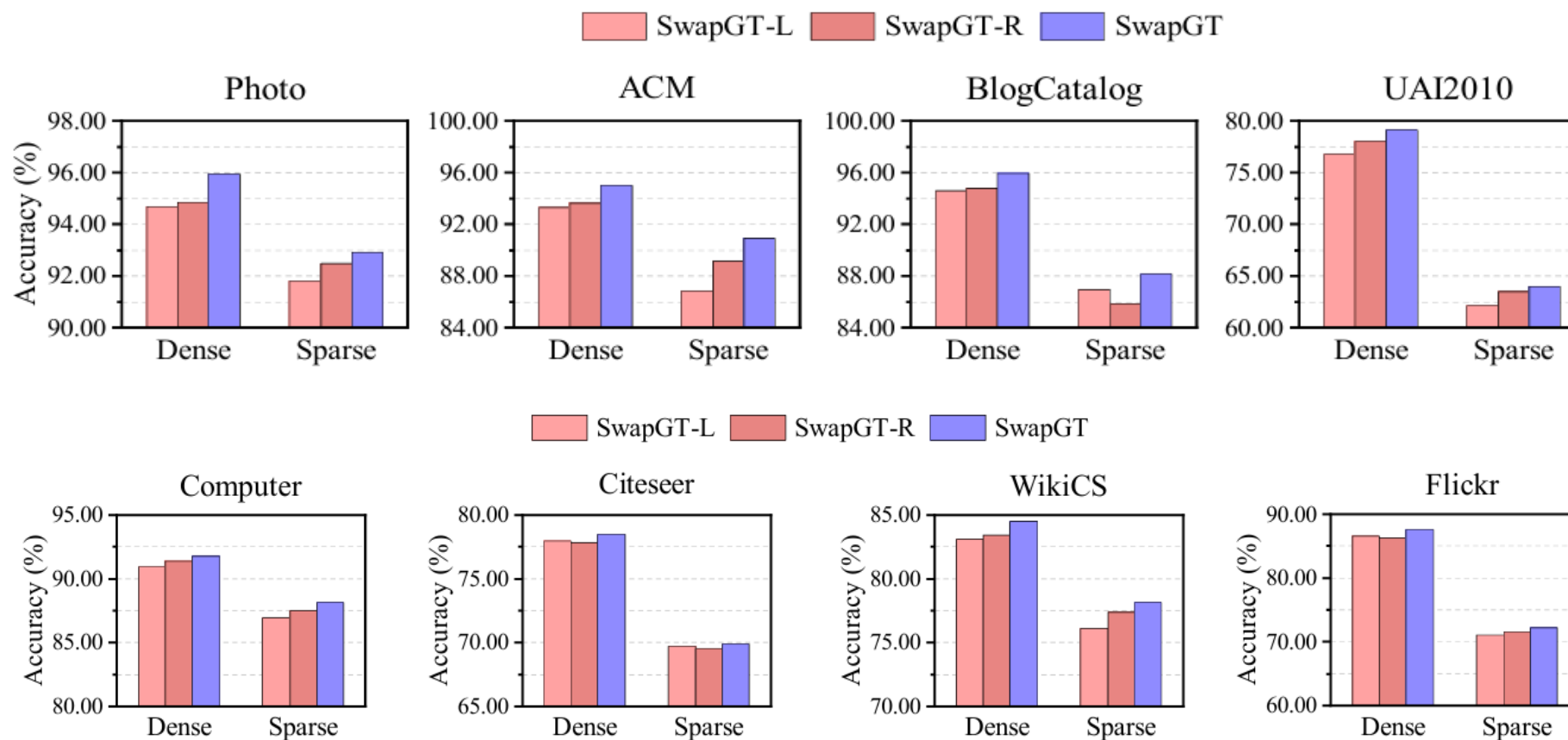


Figure 10: Performances of SwapGT with different token sequence generation strategies.

- Study on token swapping times t

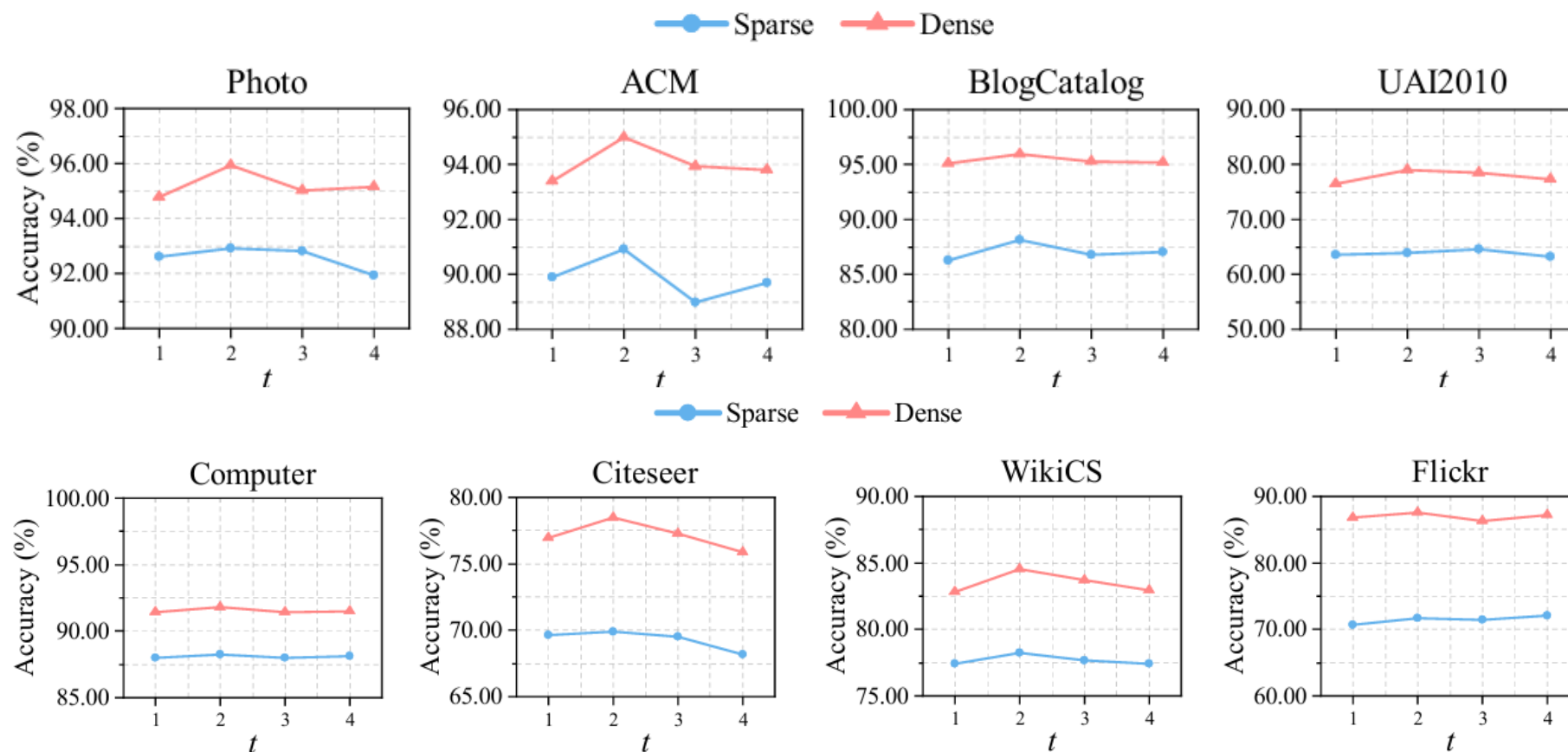


Figure 13: Performances of SwapGT with varying t .

➤ Future work



- **Rethinking SwapGT**

A potential limitation in SwapGT could be that SwapGT applies a **uniform swapping probability p** to all tokens in the sequence.

A **hierarchical probability** swapping framework may be a better solution, e.g., the **top $k/3$** nodes are assigned **higher** swapping probabilities, while others in top- k range receive lower probabilities.

This refinement could mitigate noise interference.

> End



Thanks for your attention!