



# NAGphormer: A Tokenized Graph Transformer For Node Classification In Large Graphs

Jinsong Chen<sup>1,2,3,\*</sup>, Kaiyuan Gao<sup>2,3,\*</sup>, Gaichao Li<sup>1,2,3</sup>, Kun He<sup>2,3,†</sup>

Institute of Artificial Intelligence

Huazhong University of Science and Technology, School of Computer Science and Technology

Hopcroft Center on Computing Science

chenjinsong,im\_kai,gaichaolee,brooklet60}@hust.edu.cn

Code: <https://github.com/JHL-HUST/NAGphormer>.

— ICLR 2023

2023. 8. 24 • ChongQing



gesis  
Leibniz-Institut  
für Sozialwissenschaften



Reported by Jinyuan Zhang



- 1. Introduction**
- 2. Overview**
- 3. Methods**
- 4. Experiments**





# Introduction

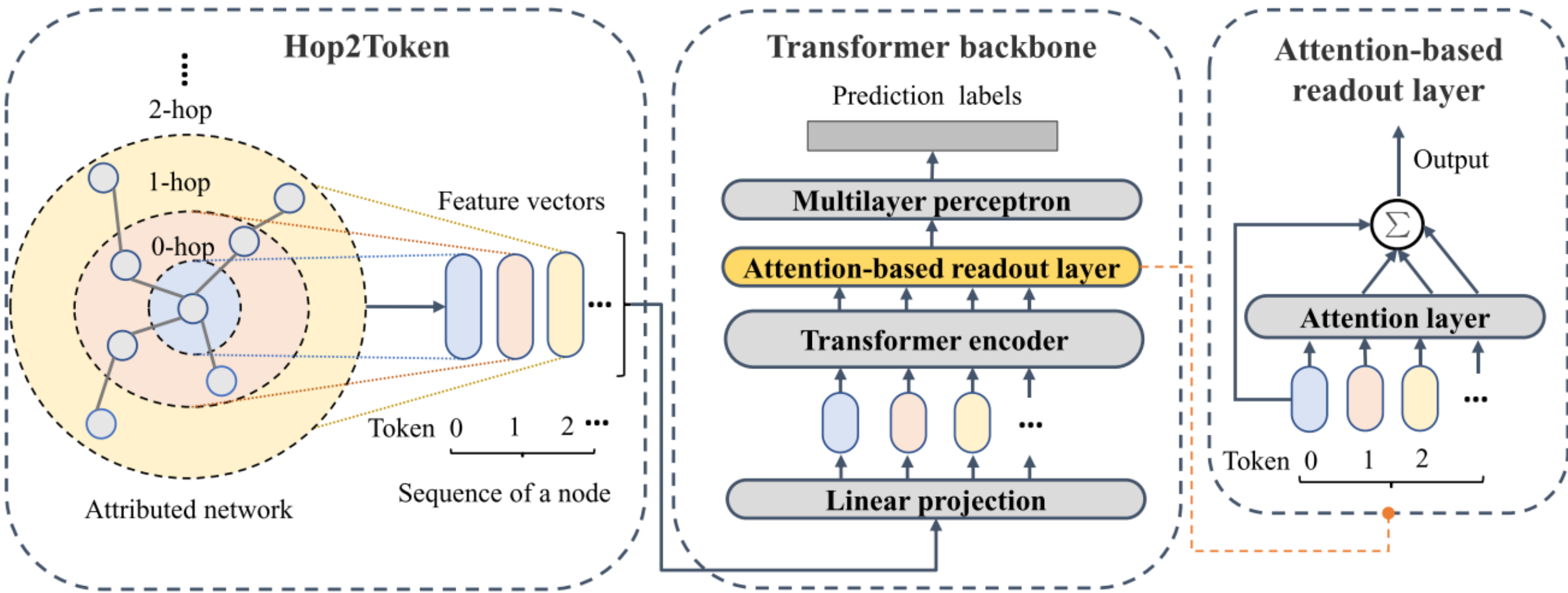
## QUESTION:

- 1、 Existing graph Transformers **hard to scale to large graphs**
- 2、 **over-smoothing** and **over-squashing** problems, the negative influence of these inherent limitations cannot be eliminated completely.

## WORK:

- 1、 **Hop2Token**, resulting in a sequence of token vectors that preserves neighborhood information for different hops, regard each node in the complex graph data as a sequence of tokens。
- 2、 **NAGphormer**, for the node classification task. In self-attention mechanism, it can learn more expressive node representations from the multi-hop neighborhoods
- 3、 develop an **attention-based readout function** to adaptively learn the importance of different-hop neighborhoods to further boost the model performance.

# Overview



Model framework of NAGphormer

# Method

## Graph Neural Network (GNN)

$$\mathbf{H}^{(l+1)} = \sigma(\hat{\mathbf{A}}\mathbf{H}^{(l)}\mathbf{W}^{(l)}), \quad (1)$$

## Decoupled Graph Convolutional Network (GCN)

$$\mathbf{Z} = \sum_{k=0}^K \beta_k \mathbf{H}^{(k)}, \mathbf{H}^{(k)} = \hat{\mathbf{A}}\mathbf{H}^{(k-1)}, \mathbf{H}^{(0)} = \mathbf{f}_\theta(\mathbf{X}),$$

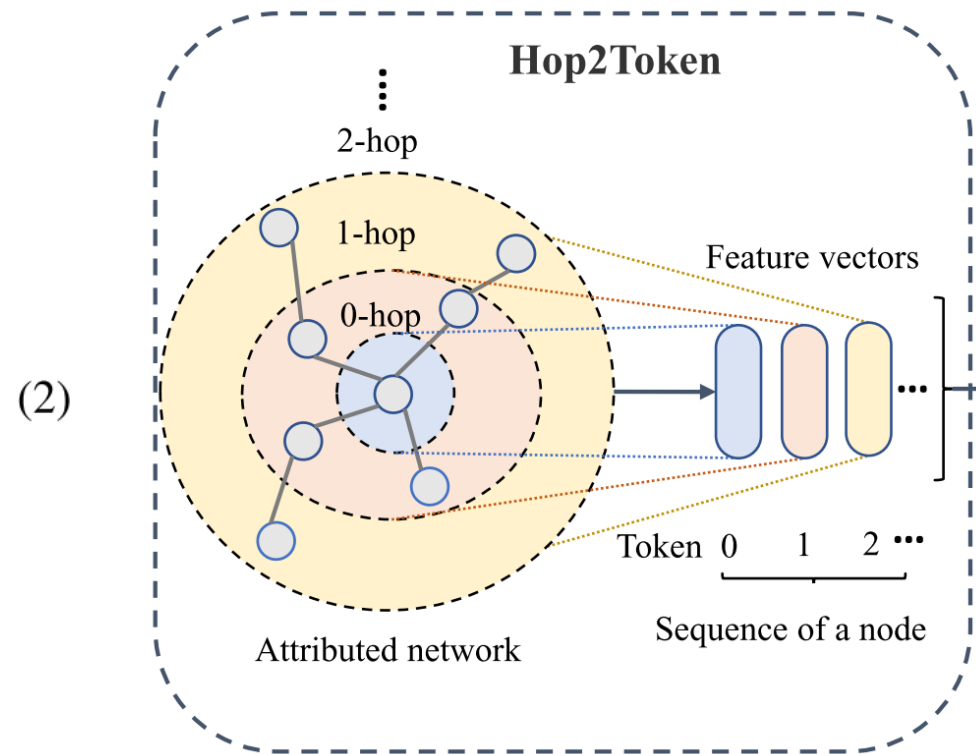
## Hop2Token

$$\mathbf{X}' = \mathbf{X} \parallel \mathbf{U}. \quad (10)$$

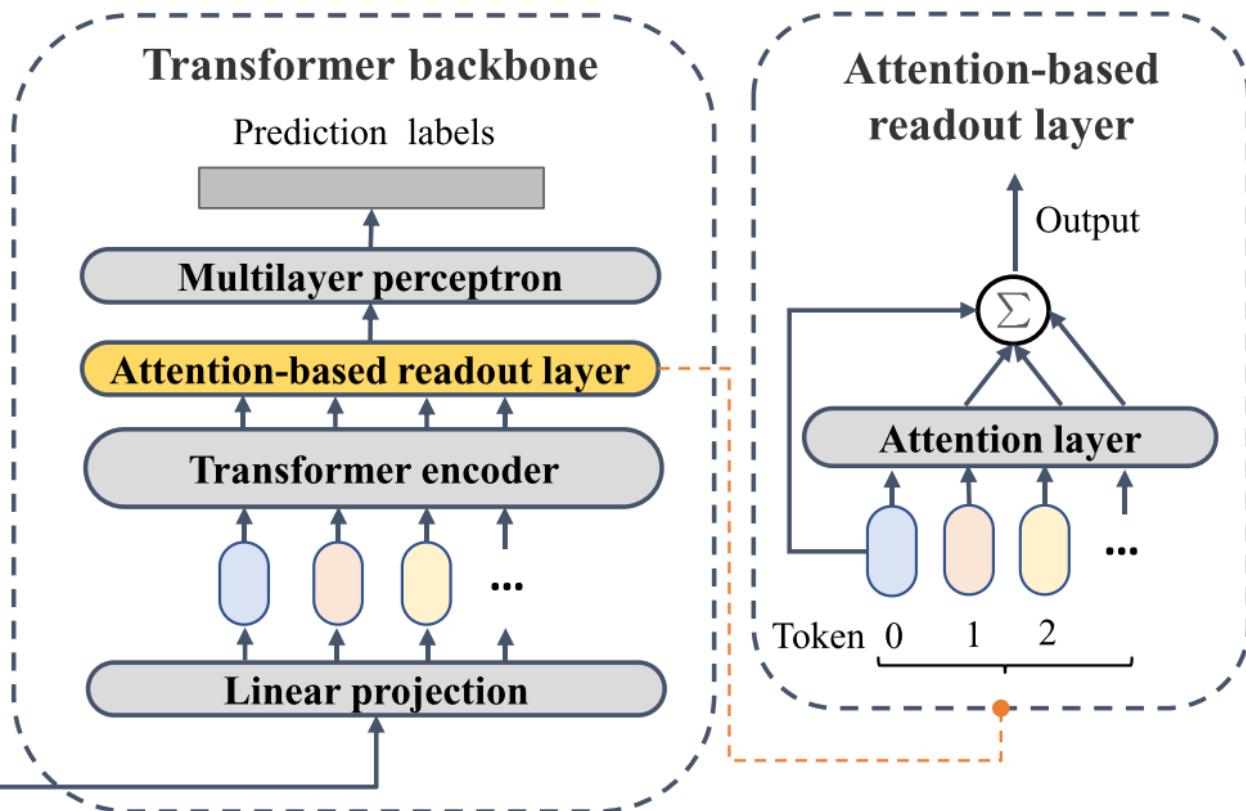
$$\mathbf{x}_v^k = \phi(\mathcal{N}^k(v)). \quad (5)$$

$$\mathbf{X}_k = \hat{\mathbf{A}}^k \mathbf{X}. \quad (6)$$

$$\mathcal{S}_v = (\mathbf{x}_v^0, \mathbf{x}_v^1, \dots, \mathbf{x}_v^K)$$



# Method



$$\mathbf{Z}_v^{(0)} = [\mathbf{x}_v^0 \mathbf{E}; \mathbf{x}_v^1 \mathbf{E}; \dots; \mathbf{x}_v^K \mathbf{E}], \quad (7)$$

$$\mathbf{Z}'_v^{(\ell)} = \text{MSA} \left( \text{LN} \left( \mathbf{Z}_v^{(\ell-1)} \right) \right) + \mathbf{Z}_v^{(\ell-1)}, \quad (8)$$

$$\mathbf{Z}_v^{(\ell)} = \text{FFN} \left( \text{LN} \left( \mathbf{Z}'_v^{(\ell)} \right) \right) + \mathbf{Z}'_v^{(\ell)}, \quad (9)$$

$$\alpha_k = \frac{\exp((\mathbf{Z}_0 \| \mathbf{Z}_k) \mathbf{W}_a^\top)}{\sum_{i=1}^K \exp((\mathbf{Z}_0 \| \mathbf{Z}_i) \mathbf{W}_a^\top)}, \quad (11)$$

$$\mathbf{Z}_{out} = \mathbf{Z}_0 + \sum_{k=1}^K \alpha_k \mathbf{Z}_k. \quad (12)$$

# Experiments

Method	Pubmed	CoraFull	Computer	Photo	CS	Physics
GCN	86.54 $\pm$ 0.12	61.76 $\pm$ 0.14	89.65 $\pm$ 0.52	92.70 $\pm$ 0.20	92.92 $\pm$ 0.12	96.18 $\pm$ 0.07
GAT	86.32 $\pm$ 0.16	64.47 $\pm$ 0.18	90.78 $\pm$ 0.13	93.87 $\pm$ 0.11	93.61 $\pm$ 0.14	96.17 $\pm$ 0.08
APPNP	88.43 $\pm$ 0.15	65.16 $\pm$ 0.28	90.18 $\pm$ 0.17	94.32 $\pm$ 0.14	94.49 $\pm$ 0.07	96.54 $\pm$ 0.07
GPRGNN	89.34 $\pm$ 0.25	67.12 $\pm$ 0.31	89.32 $\pm$ 0.29	94.49 $\pm$ 0.14	95.13 $\pm$ 0.09	96.85 $\pm$ 0.08
GraphSAINT	88.96 $\pm$ 0.16	67.85 $\pm$ 0.21	90.22 $\pm$ 0.15	91.72 $\pm$ 0.13	94.41 $\pm$ 0.09	96.43 $\pm$ 0.05
PPRGo	87.38 $\pm$ 0.11	63.54 $\pm$ 0.25	88.69 $\pm$ 0.21	93.61 $\pm$ 0.12	92.52 $\pm$ 0.15	95.51 $\pm$ 0.08
GRAND+	88.64 $\pm$ 0.09	71.37 $\pm$ 0.11	88.74 $\pm$ 0.11	94.75 $\pm$ 0.12	93.92 $\pm$ 0.08	96.47 $\pm$ 0.04
GT	88.79 $\pm$ 0.12	61.05 $\pm$ 0.38	91.18 $\pm$ 0.17	94.74 $\pm$ 0.13	94.64 $\pm$ 0.13	97.05 $\pm$ 0.05
Graphormer	OOM	OOM	OOM	92.74 $\pm$ 0.14	OOM	OOM
SAN	88.22 $\pm$ 0.15	59.01 $\pm$ 0.34	89.83 $\pm$ 0.16	94.86 $\pm$ 0.10	94.51 $\pm$ 0.15	OOM
GraphGPS	88.94 $\pm$ 0.16	55.76 $\pm$ 0.23	OOM	95.06 $\pm$ 0.13	93.93 $\pm$ 0.12	OOM
NAGphormer	<b>89.70 <math>\pm</math> 0.19</b>	<b>71.51 <math>\pm</math> 0.13</b>	<b>91.22 <math>\pm</math> 0.14</b>	<b>95.49 <math>\pm</math> 0.11</b>	<b>95.75 <math>\pm</math> 0.09</b>	<b>97.34 <math>\pm</math> 0.03</b>

Table 1: Comparison of all models in terms of mean accuracy  $\pm$  stdev (%) on small-scale datasets. The best results appear in **bold**. OOM indicates the out-of-memory error.

# Experiments

Table 2: Comparison of all models in terms of mean accuracy  $\pm$  stdev (%) on large-scale datasets. The best results appear in **bold**.

Method	AMiner-CS	Reddit	Amazon2M
PPRGo	49.07 $\pm$ 0.19	90.38 $\pm$ 0.11	66.12 $\pm$ 0.59
GraphSAINT	51.86 $\pm$ 0.21	92.35 $\pm$ 0.08	75.21 $\pm$ 0.15
GRAND+	54.67 $\pm$ 0.25	92.81 $\pm$ 0.03	75.49 $\pm$ 0.11
NAGphormer	<b>56.21 <math>\pm</math> 0.42</b>	<b>93.58 <math>\pm</math> 0.05</b>	<b>77.43 <math>\pm</math> 0.24</b>

Table 3: The accuracy (%) with or without structural encoding.

	Pubmed	CoraFull	CS	Computer	Photo	Physics	AMiner-CS	Reddit	Amazon2M
W/O-SE	89.06	70.42	95.52	90.44	95.02	97.10	55.64	93.47	76.98
With-SE	89.70	71.51	95.75	91.22	95.49	97.34	56.21	93.58	77.43
Gain	+0.64	+1.09	+0.23	+0.78	+0.47	+0.24	+0.57	+0.11	+0.45



# Experiments

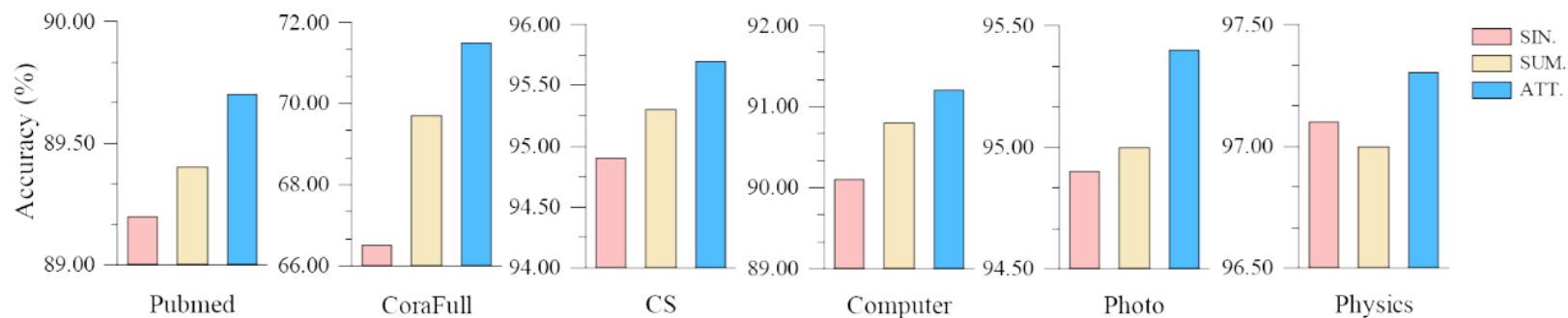
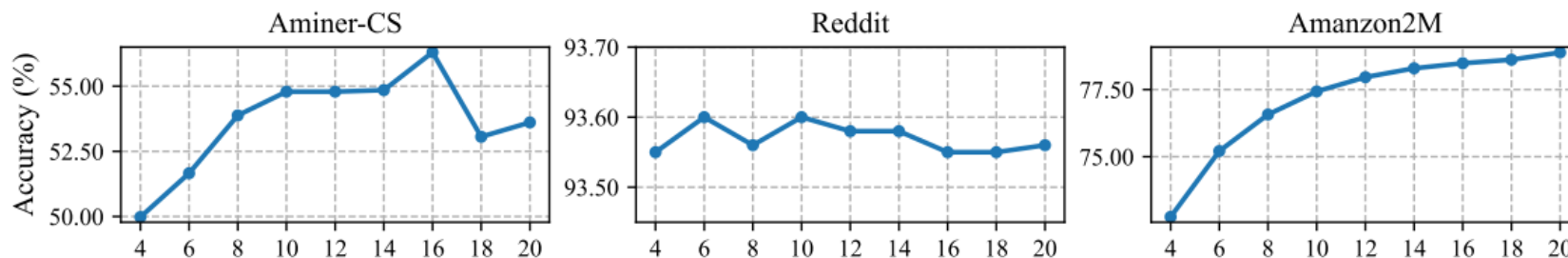
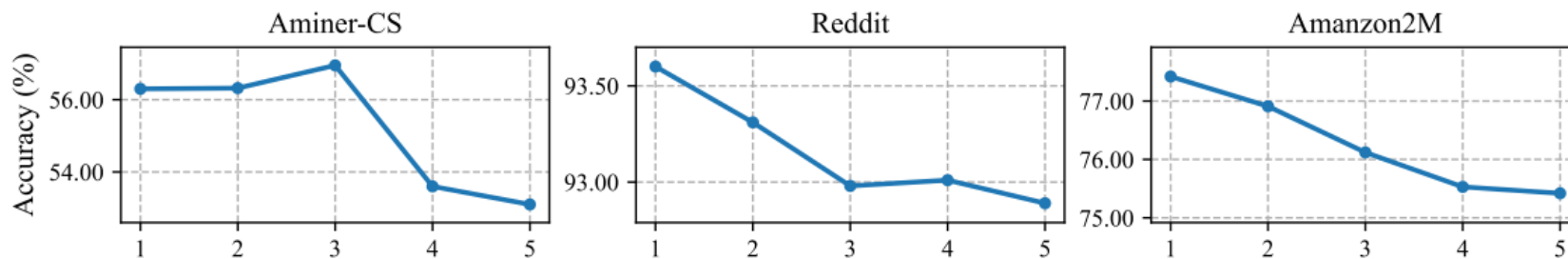


Figure 2: The performance of NAGphormer via different readout functions.

# Experiments



(a) On the number of propagation steps  $K$



(b) On the number of Transformer layers  $L$

Figure 3: Performance of NAGphormer on different parameters.



**THANKS**