Graph Triple Attention Network: A Decoupled Perspective

Xiaotang Wang*
Huazhong University of Science
and Technology
Wuhan, China
wangxiaotang0906@foxmail.com

Yun Zhu*[†] Zhejiang University Hangzhou, China zhuyun_dcd@zju.edu.cn Haizhou Shi Rutgers University New Brunswick, New Jersey, US haizhou.shi@rutgers.edu

Yongchao Liu Ant Group Hangzhou, China yongchao.ly@antgroup.com

Chuntao Hong
Ant Group
Hangzhou, China
chuntao.hct@antgroup.com

Abstract

Graph Transformers (GTs) have recently achieved significant success in the graph domain by effectively capturing both long-range dependencies and graph inductive biases. However, these methods face two primary challenges: (1) multi-view chaos, which results from coupling multi-view information (positional, structural, attribute), thereby impeding flexible usage and the interpretability of the propagation process. (2) local-global chaos, which arises from coupling local message passing with global attention, leading to issues of overfitting and over-globalizing. To address these challenges, we propose a high-level decoupled perspective of GTs, breaking them down into three components and two interaction levels: positional attention, structural attention, and attribute attention, alongside local and global interaction. Based on this decoupled perspective, we design a decoupled graph triple attention network named DeGTA, which separately computes multiview attentions and adaptively integrates multi-view local and global information. This approach offers three key advantages: enhanced interpretability, flexible design, and adaptive integration of local and global information. Through extensive experiments, DeGTA achieves state-of-the-art performance across various datasets and tasks, including node classification and graph classification. Comprehensive ablation studies demonstrate that decoupling is essential for improving performance and enhancing interpretability. Our code is available at: https://github.com/wangxiaotang0906/DeGTA.

1 Introduction

Recently, Graph Transformers (GT) [10, 22, 33] have shown great potential in handling graph-structured data such as social networks [32], drug discovery [34], and traffic networks [9]. The success of GTs is built upon their ability to leverage inductive bias and handle long-range dependencies simultaneously. For example, Graphormer [43] utilizes global attention to capture long-range dependencies and uses structural biases to influence attention scores, thereby introducing inductive bias. GraphGPS [33] employs GNNs before global attention to enhance graph topology information. Besides, these methods

will compensate positional encodings and structural encodings into node features to enhance the performance. However, existing GTs mainly face two challenges:

- 1) multi-view chaos: Multi-view encodings, such as positional encodings, structural encodings, and attribute encodings, are crucial for building a GT [10]. Existing methods couple these types of information to obtain node features, effectively harnessing graph topology and inductive biases. However, multi-view chaos constrains the separability of positional, structural, and attribute information during propagation, thus impeding flexible usage and the interpretability of the propagation process.
- 2) local-global chaos: Existing GTs couple local message passing and global attention, which complicates the adjustment of importance weights between local and global information. This uniformity potentially precipitates issues of overfitting and over-globalizing [40]. In some graphs, local information may be crucial, while in others, global information may dominate. Coupling these sources of information makes it challenging to adapt the weights appropriately and lacks interpretability regarding which information contributes to success.

To address the aforementioned challenges, we propose a high-level decoupled perspective of Graph Transformers, breaking down GTs into three components and two message interaction levels, namely: (i) **Positional Attention**: Self-attention on positional information. (ii) **Structural Attention**: Self-attention on structural information. (iii) **Attribute Attention**: Self-attention on attribute information. Additionally, the interaction levels are: (a) **Local-level**: Aggregating messages between neighbors. (b) **Global-level**: Aggregating messages between all nodes. Based on this decoupled perspective, we design a decoupled graph triple attention network, coined as DeGTA, that separately computes multi-view attentions, including positional attention, structural attention, and attribute attention, and adaptively integrates local and global information. This approach offers three key advantages:

1. *Enhanced Interpretability*: Aggregation becomes more interpretable as we can visualize the attention scores separately, enabling analysis of which information contributes the most.

^{*}Contributed equally.

[†]Corresponding author.

- 2. Flexible and Adaptive Design: The attention mechanism for multi-view information can be flexibly designed and combined in an adaptive manner.
- 3. Adaptive Integration of Local and Global Information: By employing a global sampling strategy, DeGTA can capture long-range dependencies with global positional and structural attention, and then adaptively integrate global and local information.

Our contributions can be summarized as follows:

- We identify the main challenges of GTs and propose a decoupled perspective that provides a framework for designing novel decoupled models.
- Based on this decoupled perspective, we introduce DeGTA, a decoupled graph triple attention network that enhances interpretability, enables flexible design of multi-view attention, and adaptively integrates local and global information.
- Through extensive experiments, DeGTA achieves stateof-the-art performance across various datasets and tasks, such as node-level classification and graph classification. Comprehensive ablation studies demonstrate that decoupling is essential for improving performance and enhancing interpretability.

2 Related Work

2.1 Graph Attention Networks & GTs

Graph attention-based networks endeavor to ascertain the relational significance between node pairs. These models can be classified into two main categories: (1) edge-attention, exemplified by the Graph Attention Network (GAT) [36] and its variants [3, 19], in which each source node aggregates features from its neighbors based on the deduced importance of the edges; and (2) hop-attention [8, 27], which discerns the relative importance of each hop. Hop-attention models apply attention weights to each hop's information and then compute node features through a weighted summation of different hop information.

The above models conduct message passing based on graph topology, e.g., MPNN-based [17, 23, 46, 47], which can be considered a local interaction between nodes. It's well known that MPNN-based GNNs face challenges of over-smoothing [1], oversquashing [35], and limited expressive power limitations [30]. To extend the receptive field, Graph Transformers (GTs) [10, 22] have emerged, utilizing global attention combined with positional and structural information to simultaneously capture long-range dependencies and graph inductive bias. For instance, Graphormer [43] uses centrality encoding to enhance structural information and spatial matrix as positional biases, SAN [24] uses a learnable Laplacian PE as input to a branch Transformer layer, and GRIT [28] uses a random walk encoding to integrate positional information. GraphGPS [33] explicitly integrate various other types of MPNN modules into their architectures. Other works such as NAGphormer [5] and SAT [4] use subgraph to introduce neighborhood inductive biases to node features. However, GTs face two significant

problems: (1) local-global chaos, and (2) multi-view chaos. These issues make it challenging for GTs to adapt weights appropriately, lack interpretability regarding which information contributes to success, and impede flexible usage and the interpretability of the propagation process.

In this work, we propose a high-level decoupled perspective of GTs. Based on this perspective, we design a novel decoupled graph triple attention network to address the above problems.

2.2 PE and SE

In the graph domain, PE and SE have been studied for enhancing the expressive power and performance of both MPNNs and GTs, particularly for GTs to introduce inductive bias of graph. For instance, some works incorporate Laplacian encoding [24] random walk encoding [28], shortest-path-distance [26] or centrality encoding [43] as PE/SE to capture important positional or structural relationships between nodes, with certain studies emphasizing the learnability of PE/SE. However, there is a paucity of work that clearly distinguishes between the definitions of PE and SE, many so-called positional or structural encodings comprising both positional and structural information, resulting in the coupling of information. Additionally, these encodings are often concatenated or integrated with attribute encodings to compute overall attention, leading to multi-view chaos. Section 3.2 contains a detailed revisit including the usage of PE and SE in previous works.

3 Revisiting Graph Attention Mechanism: A Decoupled Perspective

In this section, we introduce our decoupled perspective for the graph attention mechanism from two core aspects: multi-view attentions and message interactions. Then, we revisit previous work from the decoupled perspective and draw conclusions.

3.1 Decoupled Perspective

- 3.1.1 Decoupled Perspective of Multi-View Attentions. In this work, we refer to positional, structural, and attribute information as multi-view information. First, we provide definitions of these types of information, as few previous works [33] offer clear concepts of them. Subsequently, we introduce several optional encodings from our decoupled perspective.
- (i) Structural Encoding (SE) focuses on a node's capacity to perceive its surrounding structure. A node does not concern itself with the specific identities of its neighboring nodes; rather, it focuses exclusively on topological information, such as degree information, the shape of its subgraph, and other topological characteristics like triangle counting and cycle counting. In our method, we provide several strategies to achieve this objective: (1) Random-Walk Structural Encoding (RWSE) [12], (2) Node Degree Encoding (DSE) [43], and (3) Topology Counting Encoding (TCSE) [2]. The details of these encodings can be found in Appendix B.
- (ii) Positional Encoding (PE) focuses on the capability of nodes to perceive their relative positions with respect to other specific nodes and their positions within the entire graph. In

Attention Type	Position	Positional Attention		Structural Attention		Attribute Attention		Decoupling	
interaction	local	global	local	global	local	global	multi-view	local-global	
GCN [23]	Х	Х	Х	Х	√	Х	Х	Х	
GAT [36]	X	X	X	X	✓	X	X	X	
ADSF [45]	X	Х	/	X	✓	X	✓	X	
GT [10]	X	✓	X	X	X	✓	X	X	
Graphormer [43]	X	✓	X	✓	X	✓	X	X	
SAN [24]	X	✓	X	X	X	✓	✓	X	
LSPE [12]	/	Х	X	Х	1	Х	✓	X	
GraphGPS [33]	X	✓	X	✓	1	✓	X	1	
SAT [4]	X	✓	X	✓	Х	✓	X	X	
NAGphormer [5]	✓	X	✓	X	✓	X	X	X	
DeGTA	√	√	√	√	√	√	√	√	

Table 1: Revisiting previous work through our decoupled perspective. If the last two columns of "Decoupling" are not checked, the ticked attentions are considered coupled to one another.

contrast to structural information, positional information is more specific, as it can discern the identity information of other nodes. For instance, it can sense the distance to other nodes and intersections with other nodes within n-hop subgraphs. We provide several off-the-shelf strategies that can reflect a node's distance or intersection relationships with other specific nodes: (1) Laplacian Positional Encoding (LapPE) [24, 49], (2) Relative Random Walk Probabilities (RWPE) [28], and (3) Jaccard Encoding (JaccardPE) [45]. The introduction of these strategies can be found in Appendix B.

(iii) Attribute Encoding (AE) aim to transform raw node attributes, such as text or images, into numeric vectors. Recently, with the significant success of LLMs in language understanding, text-attributed graphs (TAGs) have become a prominent topic in the graph domain. For TAGs, LLMs can be employed to enhance attribute encodings [20, 48]. Thus we can select AE from original features [16] or LLM-enhanced features [37, 48].

Previous methods, like GraphGPS [33] and Graphormer [43], couple these information types to calculate attention, leading to *multi-view chaos*, which impedes flexible usage and the interpretability of the propagation process. The clear definitions can guide us in checking if the information is adequate and help address the *multi-view chaos* of attention in graphs by separately computing structural, positional, and attribute attention through different encodings of node pairs.

- 3.1.2 Decoupled Perspective of Message Interaction. Message interaction means how does messages of nodes interact with each other. There are definitions of different interactions:
- (a) Local Interaction (LI) is the information aggregation based on the original graph topology, where message passing occurs only along the edges present in the original graph.
- (b) Global Interaction (GI) refers to the message passing between all pairs of nodes whether there have edges in the original graph or not.

According to these definitions, we can identify *local-global chaos* in current GTs, which makes it challenging to adapt the weights appropriately and lacks interpretability regarding which information contributes to success. Decoupling these

interactions is essential for building a robust and interpretable graph model.

3.2 Revisiting Previous Work

In this subsection, we will revisit previous classical studies from our decoupled perspective. Table 1 shows the comparisons of different methods from this viewpoint. GCN uses AE for message passing along the graph structure and averages the aggregated neighbors. While GAT uses an adaptive attention aggregation. ADSF uses a Structural Fingerprint, indeed a Jaccard encoding as PE to compute a fixed positional attention, and adaptively integrate it with attribute attention. GT replaces GAT's local attention of AE with global attention. Graphormer uses centrality encoding as SE, and it uses spatial matrix directly as positional attentions to serve as a bias term in the self-attention module, which coupled AE and PE. SAN generates a learnable PE from a branch Transformer layer, then it concats this PE with AE, and computes global attention. This attention score is equal to the sum of positional attention and attribute attention. LSPE proposes learnable positional and structural encoding, however, its so-called structural encoding is node features and then concat with RWSE, which is equivalent to the coupling of AE and SE from our perspective, while its positional encoding is RWSE, indeed a SE as we defined. GPS explicitly splits positional encodings and structural encodings and it uses MPNN before transformer to incorporate local information into global. SAT uses subgraph extractor to couple SE and PE with AE through the subgraph information, and then input the coupled encoding to Transformer layers to get global attention. NAGphormer uses the multi-hop representation based on subgraph as the token sequence of the Transformer layer, which couples the PE, SE, and AE, and computes attention between these local representations.

As shown in Table 1 and the above discussions, previous works struggle to decouple multi-view information and different-level message interactions, leading to *multi-view chaos* and *local-global chaos*. Additionally, some methods do not utilize sufficient information and interactions, resulting in sub-optimal performance. In order to solve these issues, we need to answer the following research questions:

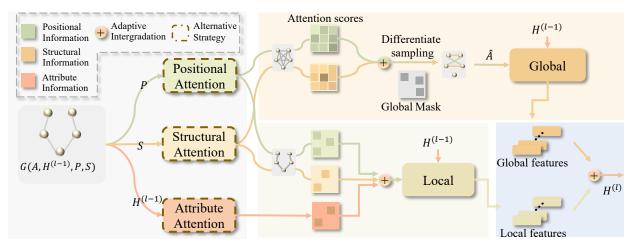


Figure 1: Framework of DeGTA. The framework comprises four main components: Decoupled Multi-View Encoder, Local Channel, Global Channel, and Local-Global Integration. The strategies of encoders for multi-view and attention mechanism for local and global are all optional.

RQ1: How to rationally utilize all the information in the graph, while avoiding the multi-view chaos?

RQ2: How to avoid the local-global chaos while capturing the long-range dependencies and maintaining the inductive biases?

We will design our model, named DeGTA, by addressing these problems to resolve existing issues. The details are illustrated in the next section.

4 Decoupled Graph Triple Attention Network

In this section, we will introduce our designed model, named DeGTA, by addressing RQ1 in Section 4.2 and RQ2 in Section 4.3. Then, we discuss DeGTA from a theoretical perspective in Section 4.4. The overall framework of our method is shown in Figure 1.

4.1 Notations

Let $G=(\mathcal{V},A,X)$ denote a graph, where \mathcal{V} denote the sets of nodes, respectively. $X\in\mathbb{R}^{N\times d}$ represents the node attribute feature matrix with N nodes and d attributes. $A\in\mathbb{R}^{N\times N}$ is the adjacency matrix. For positional and structural information, we denote the initial PE as $P\in\mathbb{R}^{N\times K}$, and the initial SE as $S\in\mathbb{R}^{N\times K}$. Additionally, we denote the diagonal degree matrix as $D\in\mathbb{R}^{N\times N}$. The matrices with self-loops are denoted as $\tilde{A}\in\mathbb{R}^{N\times N}$ and $\tilde{D}\in\mathbb{R}^{N\times N}$, respectively. Thus, the graph Laplacian with self-loops is $\tilde{L}=\tilde{D}-\tilde{A}$, the normalized adjacency matrix is $\hat{A}=\tilde{D}^{-1}\tilde{A}$, and the symmetrically normalized graph Laplacian is $\hat{L}=\tilde{D}^{-\frac{1}{2}}\tilde{L}\tilde{D}^{-\frac{1}{2}}$. For clarity in notations, we will use A to denote the original graph (local level). The sampling graph in global level will be represented by \hat{A} .

4.2 Decoupling Multi-view Attention (RQ1)

To rationally utilize all the information in the graph and avoid multi-view chaos, we meticulously design multi-view encodings and decouple multi-view attentions. First, we introduce the initialization strategies for structural, positional, and attribute information. Then, we illustrate how to achieve decoupled multi-view attention.

4.2.1 Multi-view Information Encodings. In this subsection, we employ the available encoding strategies for positional, structural, and attribute information that introduced in Section 3.1.1. Because our method decouples multi-view attention, it allows for flexible selection of encoding strategies.

Strategy Selection. Different graphs may require different positional and structural encoding strategies, making it challenging to choose a specific encoding strategy that performs best across all datasets. In this work, considering simplicity and effectiveness, we choose Jaccard Positional Encoding (JaccardPE) and Random-Walk Structural Encoding (RWSE) as our positional and structural encodings, respectively. For node attribute encodings, we use the provided node encodings by PyG [16]. Exploration of more complex strategy selection will be left for future work.

4.2.2 Decoupled Multi-view Attention. To avoid multi-view chaos, unlike previous works that operate overall attention with mixed encodings, we handle them separately with independent encoders. This decoupled approach allows us to flexibly choose the type of encoders for each type of information as follows:

$$S^{(l)} = \mathcal{E}_s^{(l)}(S), \quad P^{(l)} = \mathcal{E}_p^{(l)}(P), \quad H^{(l)} = \mathcal{E}_a(H^{(l-1)})$$
 (1)

where $\mathcal{E}_s(\cdot): \mathbb{R}^{N \times K} \mapsto \mathbb{R}^{N \times d_s}$, $\mathcal{E}_p(\cdot): \mathbb{R}^{N \times K} \mapsto \mathbb{R}^{N \times d_p}$, and $\mathcal{E}_a(\cdot): \mathbb{R}^{N \times d} \mapsto \mathbb{R}^{N \times d}$ are alternative encoders for processing structural, positional and attribute encodings respectively. Note that $d_s, d_p \ll d$ because K is always small to

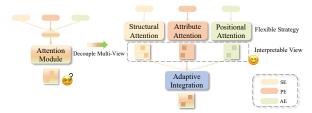


Figure 2: Comparison of traditional attention and decoupled multi-view attention. Our method enables the flexible design of distinct attention mechanisms for various encodings, and enhances interpretability by the capacity to visualize the attention scores independently.

control the receptive field, generally no more than 8. To retain the pure structural and positional information, we feed the initial SE and PE into each layer. For AE, to obtain high-level features, we use the output of the previous layer as input.

In this work, we design a graph triple attention network that utilizes independent attention modules to process positional, structural, and attribute encodings. As shown in Figure 2, this approach allows us to visualize the attention scores separately, facilitating an analysis of which information is similar in a node pair, thereby improving interpretability. We then adaptively learn the weights of the three types of attention for information aggregation, enabling an analysis of which information contributes the most through the learned weights.

4.3 Decoupling Message Interaction (RQ2)

To avoid the local-global chaos while capturing long-range dependencies and maintaining graph inductive bias (RQ2), we decouple the local and global interaction in this work by designing separate attention mechanisms for each level of interaction. This decoupling allows us to effectively manage the unique characteristics and importance of local and global information, enhancing both performance and interpretability.

4.3.1 Decoupling Local-level Interaction. The local-level interaction focuses on message passing between neighboring nodes, leveraging the immediate graph topology. This interaction is crucial for capturing fine-grained local structures and details within the graph. In our work, we use MPNN-based attention methods to compute local attention as follows:

$$s_{i,j} = q_s^T \sigma([W_{\text{str}} S_i \parallel W_{\text{str}} S_j])$$

$$p_{i,j} = q_p^T \sigma([W_{\text{pos}} P_i \parallel W_{\text{pos}} P_j])$$

$$a_{i,j} = q_a^T \sigma([W_{\text{atr}} H_i \parallel W_{\text{atr}} H_i])$$
(2)

where \cdot^T represents transposition and \parallel is the concatenation operation, σ is a non-linear activation, e.g., LeakyReLU. $W_{\text{str}} \in \mathbb{R}^{d' \times d_s}$, $W_{\text{pos}} \in \mathbb{R}^{d' \times d_p}$, $W_{\text{atr}} \in \mathbb{R}^{d'' \times d}$ denote feature transformation matrices, and $q_s, q_p \in \mathbb{R}^{2d'}$ and $q_a \in \mathbb{R}^{2d''}$ are learnable weight vectors. The decoupled multi-view attention scores, i.e., $s_{i,j}$, $p_{i,j}$, and $a_{i,j}$, can assist in the interpretability through visualizing the attention scores separately.

Local Integration. Then, we adaptively combine the scores of local structural, positional, and attribute attention:

$$z_{i,j} = \alpha p_{i,j} + \beta s_{i,j} + \gamma a_{i,j} \tag{3}$$

where α , β , and γ are learnable weights for the multi-view attention, enabling adaptive use of decoupled multi-view information to compute the total weight of local aggregation. This adaptive mechanism allows for an analysis of which type of information contributes the most to the graph through the learned weights, enhancing both the flexibility and interpretability of the model.

Last, we normalize the attention scores and aggregate the local information through AE along the original graph:

$$\hat{H}_i^{\text{local}} = \sum_{i \in \mathcal{N}_i} \hat{z}_{i,j} H_j, \text{ where } \hat{z}_{i,j} = \frac{e^{z_{i,j}}}{\sum_{k \in \mathcal{N}_i} e^{z_{i,k}}}$$
(4)

where H_i , H_j is attribute encoding of node i and j. and \mathcal{N}_i is the set of neighbours of node i.

4.3.2 Decoupling Global-level Interaction. The global-level interaction captures long-range dependencies and overall graph structures by employing global attention mechanisms that can link distant nodes and aggregate information across the entire graph. To achieve this, we will sample a new graph topology that can effectively capture long-range dependencies from a global viewpoint. This new topology will enable the model to consider connections beyond immediate neighbors, ensuring comprehensive integration of global information. First, we compute global attention by Transformer attention layer:

$$U_s = \text{Softmax}(\frac{Q_s K_s^T}{\sqrt{d_s}}), \quad U_p = \text{Softmax}(\frac{Q_p K_p^T}{\sqrt{d_p}})$$
 (5)

where $Q_s = SW_s^q, Q_p = PW_p^q$ are the global query features of structural and positional encodings, and $K_s = SW_s^k, K_p = PW_p^k$ are the global key features of structural and positional encodings, $W_s^q, W_s^k \in \mathbb{R}^{d_s \times d'}$ and $W_p^q, W_p^k \in \mathbb{R}^{d_p \times d'}$ are learnable weight matrices.

Then, we adaptively activate long-range dependency edges for each node according to the results of global structural attention scores and positional attention scores. We sample the long-range nodes using two alternative strategies: (1) Top-K and (2) Threshold control. Below, we illustrate the threshold control strategy:

$$M = \text{Softmax}((\alpha U_s + \beta U_p) \odot (1 - A)), \tag{6}$$

Here \odot is element-wise multiplication with (1 - A) to ensure the sampling only focus on non-neighboring nodes. Then we use a differentiable sampling trick as follows:

$$\hat{A} = \mathbb{1}_{M > \tau} - \tilde{M} + M \tag{7}$$

where $\mathbbm{1}_{M>\tau}$ is an indicator function that returns 1 if $M>\tau$. \tilde{M} results from truncating the gradient of M. It ensures the output is one-hot and maintains the original gradients. Last, \hat{A} represents an active mask used to determine which long-range dependency nodes are sampled.

Different from global attention aggregation in previous works [10, 22, 33], we employ differentiable hard sampling to restrict global message passing to key pairs of long-range nodes. This approach captures critical long-range dependencies while avoiding the overfitting and over-globalizing problems [40] associated with the high degree of freedom in information propagation. Specifically, we adaptively combine the scores of global structural attention and positional attention with the attribute attention among sampled long-range nodes:

$$z_{i,j} = \alpha U_s[i,j] + \beta U_p[i,j] + \gamma \hat{U}_a[i,j], \ \forall j \in \mathcal{K}_i$$
 (8)

where \mathcal{K}_i is the sampled long-range nodes set of node i, $\hat{U}_a = \operatorname{Softmax}\left(\frac{Q_a K_a^T}{\sqrt{d}}\right)$ represents the global attention score of attribute encodings, where Q_a and K_a are obtained using the sampled global topology \hat{A} . By focusing on the sampled edges rather than all node pairs, we significantly reduce the computational complexity of computing attribute attention. After obtaining these scores, we normalize the attention scores and aggregate the global information, efficiently capturing long-range dependencies:

$$\hat{H}_{i}^{global} = \sum_{j \in \mathcal{K}_{i}} \hat{z}_{ij} H_{j}, \text{ where } \hat{z}_{i,j} = \frac{e^{z_{i,j}}}{\sum_{k \in \mathcal{K}_{i}} e^{z_{i,k}}}$$
(9)

4.3.3 Adaptive Local-Global Integration. Finally, an additional adaptive module will be employed to integrate the local information and global information to obtain the output of AE $H^{(l)}$ in l-th layer:

$$H^{(l)} = W_{l}^{(l)} \hat{H}^{\text{local},(l)} + W_{g}^{(l)} \hat{H}^{\text{global},(l)}$$
 (10)

where $\hat{H}^{\mathrm{local},(l)}$, $\hat{H}^{\mathrm{global},(l)}$ are the local and global attribute encodings in l-th layer, $W_{\mathrm{l}}^{(l)}$ and $W_{\mathrm{g}}^{(l)} \in \mathbb{R}^{d \times d}$ are learnable weight matrices in the l-th layer. These matrices enable us to adaptively adjust the weights for local and global information, catering to graphs with different characteristics. This adaptive weighting helps avoid local-global chaos by dynamically balancing the contributions of local and global interactions.

4.4 Discussion

4.4.1 Complexity Analysis. In this subsection, we will present a theoretical complexity analysis of DeGTA.

Time complexity. For simplicity, we assume that the feature dimension remains unchanged and that the number of model layers is set to 1. The time complexity of DeGTA mainly depends on four modules. The complexity of the encoder module for multi-view encodings is $O(N(d^2+2K^2))$. The complexity of the local attention module is $O(E(2K+d)+N(d^2+2K^2))$. The complexity of our global attention module is $O(2N^2K+NKd+N(d^2+2K^2))$. Finally, the complexity of multi-view attention integration and local-global integration is O(N) and $O(Nd^2)$ respectively. Thus the total time complexity of our method is $O(2N^2K+E(d+2K)+N(4d^2+Kd+6K^2))$. Given that $K\ll d$ and $E\ll N^2$, and for the sake of clarity, we omit the smaller variables. The time complexity is simplified as $O(N^2K+Ed)$, which is more efficient compared to GT's $O(N^2d)$.

Space complexity. The space complexity of the encoder module for multi-view encodings is $O(N(d+2K)+d^2+2K^2)$. The complexity of the local attention module is $O(3E+d^2+2K^2+N(d+2K))$. The complexity of the global attention module is $O(2N^2+E+2K^2+d^2+N(d+2K))$. Finally, the complexity of local-global integration is $O(d^2)$. Thus the total space complexity is $O(2N^2+4E+3d^2+6K^2+3N(d+2K))$. Given that $K\ll d$ and $E\ll N^2$, and for the sake of clarity, we omit the smaller variables. The space complexity can be simplified to $O(N^2+E+d^2+Nd)$, which is comparable to GT's $O(N^2+d^2+Nd)$.

4.4.2 Analysis for Expressivity and Over-smoothing. We analyze the expressive power of DeGTA and present a case study in Appendix C. DeGTA exhibits greater expressive power than the 1-WL test through its positional and structural encodings and global attention mechanism, and the case study shows that DeGTA can produce correct results in scenarios where the coupled encoding approach yields incorrect outcomes.

Additionally, we provide both empirical and theoretical evidence demonstrating that DeGTA offers stronger resistance to over-smoothing compared to GT. Details are in Appendix D.

5 Experiments

In this section, we present a comprehensive empirical investigation. Specifically, we aim to address the following research questions: $\mathbf{RQ1}$: How does the proposed DeGTA perform in node classification tasks, including both homophilic and heterophilic datasets? $\mathbf{RQ2}$: How does DeGTA perform in graphlevel tasks? $\mathbf{RQ3}$: How effectively does DeGTA capture longrange dependencies? $\mathbf{RQ4}$: How does the decoupling of multiview and local-global influence performance? Additionally, we provide a parameter study on K in Appendix A.2, and the case study that demonstrates our enhanced expressive power and interpretability in Figure 6.

5.1 Experimental Setup

5.1.1 Datasets. We employ 10 benchmark datasets for node classification, consisting of 4 homophilic, 4 heterophilic, and 2 large-scale datasets. Additionally, we utilize 5 benchmark datasets for graph-level tasks, of which 2 are long-range graph benchmarks. Detailed information is provided in Appendix A.

5.1.2 Baselines. For node classification tasks, the baselines primarily consist of three categories: (1) MPNN methods without attention such as GCN [23], GCNII [6], GraphSAGE [18], and GPRGNN [8], (2) attention-based MPNN methods including GAT [36], GATv2 [3], ADSF [45], and AERO-GNN [25], (3) Graph Transformer methods like GT [10], NodeFormer [38], GraphGPS [33], NAGphormer [5], and SGFormer [39].

For graph classification and regression tasks, we compare our method with (1) MPNN-based methods including GCN [23], GIN [41], and GAT [36], (2) GT-based methods including GT [10], EGT [22], SAN [24], GraphGPS [33], GRIT [28].

5.2 Node Classification (RQ1)

Table 2: Node classification performance on homophilic and heterophilic graphs. The boldface and underscore show
the best and the runner-up, respectively.

Туре	Homophilic graphs					Heterophilic graphs			
Dataset	Pubmed	Citeseer	Cora	Arxiv	Texas	Cornell	Wisconsin	Actor	Avg
GCN [23]	79.54 ± 0.4	72.10 ± 0.5	82.15 ± 0.5	71.74 ± 0.3	65.65 ± 4.8	58.41 ± 3.3	62.02 ± 5.9	30.57 ± 0.7	65.27
GCNII [6]	80.14 ± 0.7	72.80 ± 0.5	84.33 ± 0.5	72.74 ± 0.2	78.59 ± 6.6	78.84 ± 6.6	81.41 ± 4.7	35.76 ± 1.0	73.08
GraphSAGE [18]	78.67 ± 0.4	71.85 ± 0.6	83.76 ± 0.5	71.49 ± 0.3	82.43 ± 6.1	75.95 ± 5.3	81.18 ± 5.5	34.23 ± 1.0	72.45
GPRGNN [8]	75.68 ± 0.4	71.60 ± 0.8	84.20 ± 0.5	71.86 ± 0.3	81.51 ± 6.1	80.27 ± 8.1	84.06 ± 5.2	35.58 ± 0.9	73.10
GAT [36]	78.91 ± 0.4	71.89 ± 0.8	83.18 ± 0.5	71.95 ± 0.4	60.46 ± 6.2	58.22 ± 3.7	63.59 ± 6.1	30.36 ± 0.9	64.82
GATv2 [3]	79.12 ± 0.3	71.15 ± 1.1	83.88 ± 0.6	72.14 ± 0.5	60.32 ± 7.0	58.35 ± 3.8	61.94 ± 4.7	30.27 ± 0.8	64.65
ADSF [45]	80.21 ± 0.4	73.00 ± 0.4	84.29 ± 0.5	72.64 ± 0.5	78.15 ± 6.1	77.52 ± 5.9	69.24 ± 4.1	34.68 ± 0.9	71.22
AERO-GNN [25]	80.59 ± 0.5	73.20 ± 0.6	83.90 ± 0.5	72.41 ± 0.4	84.35 ± 5.2	81.24 ± 6.8	84.80 ± 3.3	36.57 ± 1.1	74.63
GT [10]	79.08 ± 0.4	70.16 ± 0.8	82.22 ± 0.6	70.63 ± 0.4	84.18 ± 5.4	80.16 ± 5.2	82.74 ± 6.0	34.28 ± 0.7	72.93
NodeFormer [38]	79.90 ± 1.0	72.50 ± 1.1	82.20 ± 0.9	71.24 ± 0.6	81.61 ± 5.4	82.15 ± 6.7	83.17 ± 5.8	36.28 ± 1.2	73.63
GraphGPS [33]	79.94 ± 0.3	72.75 ± 0.6	82.44 ± 0.6	70.97 ± 0.4	82.21 ± 6.9	82.06 ± 5.1	85.36 ± 4.2	36.01 ± 0.9	73.97
NAGphormer [5]	80.57 ± 0.3	72.43 ± 0.8	84.20 ± 0.5	70.13 ± 0.6	80.12 ± 5.5	79.89 ± 7.1	82.97 ± 3.2	34.24 ± 0.9	73.07
SGFormer [39]	80.30 ± 0.6	72.60 ± 0.2	$\underline{84.50\pm0.8}$	72.63 ± 0.1	84.29 ± 5.2	81.64 ± 5.0	85.29 ± 5.7	37.90 ± 1.1	74.89
DeGTA (ours)	81.19 ± 0.7	73.70 ± 0.4	84.79 ± 0.7	73.26 ± 0.2	85.44 ± 4.8	83.19 ± 4.8	86.95 ± 5.9	37.87 ± 1.0	75.80

5.2.1 Experiments on Small and Medium-scale Datasets. For homophilic graphs, MPNNs benefit from their inductive bias based on the homophily assumption, resulting in higher performance compared to GTs. In contrast to MPNNs, our approach enhances the utilization of positional, structural, and attribute information in graphs by decoupling multi-view information and introducing long-range dependencies through a global sampling strategy that MPNNs cannot capture. Unlike GTs, which rely on global attention aggregation, our sampling strategy focuses on capturing only the most important and plausible long-range dependencies. Additionally, by decoupling local and global interactions and adaptively integrating them, DeGTA places greater emphasis on local information, which is more crucial for homophilic graphs. Therefore, DeGTA outperforms both advanced MPNNs and GTs on homophilic graphs.

For heterophilic datasets, GTs always outperform MPNNs due to their ability to utilize global attention, which effectively filters out inter-class edges from neighboring nodes and captures disconnected yet informative nodes in the graph. Our hard sampling strategy enhances the accuracy of dependency capture, thereby mitigating the overfitting associated with global attention on small graphs. Thus DeGTA achieves significant improvement on Texas, Cornell, and Wisconsin. For the Actor dataset, our results are on par with the SOTA model, likely because the dataset is less reliant on local information, making global attention sufficient.

In summary, DeGTA outperforms the SOTA baselines across 7 out of 8 datasets, achieving an average absolute improvement of 1% over the runner-up. This underscores the importance of addressing multi-view and local-global chaos.

5.2.2 Experiments on Large-scale Datasets. To evaluate the performance of DeGTA on large datasets, we conducted experiments on Aminer-CS and Amazon2M. As shown in Table 3, comparing with 2 scalable GNNs and NAGphormer,

Table 3: Node classification performance on large graphs.

Dataset	Aminer-CS	Amazon2M
GraphSAINT[44]	51.91 ± 0.20	75.20 ± 0.18
GRAND+[14]	54.76 ± 0.23	75.83 ± 0.21
NAGphormer[5]	56.21 ± 0.42	77.43 ± 0.24
DeGTA (ours)	56.38 ± 0.51	78.49 ± 0.29

DeGTA achieves the best performance on all datasets, surpasses the runner-up over 1% absolute improvement on Amazon2M, demonstrating its effectiveness for large-scale graphs.

This further validates the effectiveness of decoupling multiview attention and message interaction in capturing information and learning implicit long-range dependencies at a large scale among numerous nodes. Moreover, as detailed in Appendix A.2, a larger value of K is essential for achieving optimal performance. This is particularly crucial for large graphs, where effectively capturing long-range dependencies and maintaining a broader receptive field for both positional and structural information is necessary.

5.3 Graph Classification (RQ2&RQ3)

5.3.1 Experiments on Classical Datasets (RQ2). To evaluate the effectiveness of DeGTA across different tasks, we generalize it to both classification and regression tasks at the graph level. As shown in Table 4, DeGTA consistently outperforms the baselines across all datasets, achieving significant improvements. For instance, DeGTA enhances accuracy by an absolute improvement of 4.5% compared to the runner-up on the CIFAR10 dataset. By leveraging multi-view encodings and facilitating the adaptive interaction of local-global information, DeGTA learns suitable representations for all nodes throughout the graph, enabling mean pooling to effectively yield strong results.

Table 4: Results for graph classification tasks.

		•		
Dataset	ZINC	MNIST	CIFAR10	
Metric	MAE↓	ACC↑		
GCN[23]	0.367 ± 0.011	90.705 ± 0.218	55.710 ± 0.381	
GIN[41]	0.526 ± 0.051	96.485 ± 0.252	55.255 ± 1.527	
GAT[36]	0.474 ± 0.007	95.535 ± 0.205	64.223 ± 0.455	
GT[10]	0.226 ± 0.014	-	-	
EGT[22]	0.108 ± 0.009	98.173 ± 0.087	68.702 ± 0.409	
SAN[24]	0.139 ± 0.006	-	-	
Graphormer[43]	0.122 ± 0.006	-	-	
GraphGPS[33]	0.070 ± 0.004	98.051 ± 0.126	72.298 ± 0.356	
DeGTA (ours)	0.059 ± 0.004	98.230 ± 0.112	76.756 ± 0.927	

Table 5: Results for graph classification tasks of longrange graph benchmarks (LRGB).

Dataset Metric	Peptides-func AP↑	Peptides-struct MAE↓
GCN[23] GIN[41] GAT[36]	0.5930 ± 0.0023 0.5498 ± 0.0079 0.5842 ± 0.0046	0.3496 ± 0.0013 0.3547 ± 0.0045 0.3504 ± 0.0028
GT[10] SAN[24] GraphGPS[33] GRIT[28]	0.6326 ± 0.0126 0.6439 ± 0.0075 0.6535 ± 0.0041 0.6988 ± 0.0082	0.2529 ± 0.0016 0.2545 ± 0.0012 0.2500 ± 0.0012 0.2460 ± 0.0012
DeGTA (ours)	0.7023 ± 0.0101	0.2437 ± 0.0014

5.3.2 Experiments on Long-range Dependencies Datasets (RQ3). To further evaluate DeGTA's ability to capture long-range dependencies, we tested our method on the Long-Range Graph Benchmark (LRGB) [13]. Specifically, we conducted experiments on two peptide graph benchmarks from LRGB, namely Peptides-func and Peptides-struct. As shown in Table 5, GTs significantly outperform MPNNs in these datasets due to their global attention mechanism, which effectively captures longrange dependencies. However, our method more accurately captures these dependencies in the graph by decoupling local-global interactions and using an adaptive integration. This strategy shows the best performance on both datasets.

5.4 Ablation Study (RQ4)

We perform comprehensive ablation studies on the importance of designs in DeGTA, including: (1) The decoupling of multi-view encodings. (2) The hard sampling strategy. (3) The decoupling of local-global integration.

Decoupling Multi-view Encodings. We conduct in-depth experiments to assess the performance impact of coupling various encodings, underscoring the importance of decoupling multi-view information. We use DeGTA with only AE as the baseline. Parentheses (A+B) indicate coupled encoding, where the two encodings are concatenated and processed by a single encoder. In contrast, A+B without parentheses refers to decoupled encodings, which utilize different encoders, with attention computed separately and integrated adaptively. For example, AE+(SE+PE) indicates that we use AE to compute attribute attention, while the fused encoding of SE and PE computes another coupled attention. As shown in Figure 3,

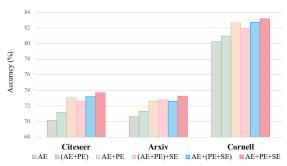


Figure 3: The results of ablation experiments for multiview decoupling.

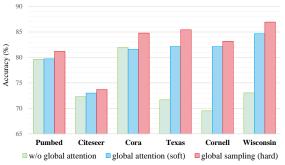


Figure 4: The results on different strategies of global information.

the decoupling of multi-view encoding facilitates a more comprehensive and adaptive usage of the information embedded in graphs, thus significantly improving the model's performance. Hard Sampling Strategy of DeGTA. Unlike the global attention aggregation in Graph Transformers (GTs), we confine global message passing to key pairs of long-range nodes using a differentiable hard sampling strategy. In this section, we compare the impact of different global attention aggregation methods on DeGTA's performance. We remove the global branch of DeGTA to serve as a baseline and evaluate the performance difference between global attention aggregation and hard-sampling aggregation. As shown in Figure 4, our experiments reveal that hard sampling more accurately captures long-distance dependencies and graph structures, leading to enhanced performance. Additionally, our sampling strategy reduces the time complexity associated with computing global attribute attention, as we only compute the attention for the sampled node pairs.

Decoupling Local-global Integration. To evaluate the effectiveness of local-global decoupling, we conducted ablation experiments in two steps. First, we use the results of global attention instead of local attention in the local channel of DeGTA, resulting in a coupled local-global attention strategy. Second, we employed different strategies for local and global attentions but simply summed the local features and global features, thereby coupling local message passing and global attentional aggregation. As shown in Table 6, both changes lead to a significant decrease in performance, showing the importance of local-global decoupling and adaptive integration.

Table 6: The ablation results for local-global decoupling. c_attn and de_attn denote the coupled and decoupled attention. c_inte denotes the local-global coupling integration and ada_inte denotes the adaptive integration.

	Pubmed	Citeseer	Cora	Texas	Cornell	Wisconsin
c_attn	79.97	72.40	82.44	83.59	80.46	83.27
de_attn + c_inte	80.06	72.85	82.64	82.26	82.17	84.49
de_attn + ada_inte	81.19	73.70	84.79	85.44	83.19	86.95

6 Conclusion

In this work, we propose a decoupled perspective to analyze attentions in graph, breaking them down into three components and two message interaction levels. This perspective helps us identify the issues of multi-view and local-global chaos in GTs. To address these challenges, we design DeGTA, a decoupled graph triple attention network. Extensive experiments demonstrate that DeGTA achieves SOTA performance across various datasets and tasks, highlighting the effectiveness of decoupling multi-view attention and local-global interaction.

References

- Uri Alon and Eran Yahav. 2021. On the Bottleneck of Graph Neural Networks and its Practical Implications. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021.
 OpenReview.net. https://openreview.net/forum?id=i80OPhOCVH2
- [2] Giorgos Bouritsas, Fabrizio Frasca, Stefanos Zafeiriou, and Michael M. Bronstein. 2023. Improving Graph Neural Network Expressivity via Subgraph Isomorphism Counting. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 45, 1 (Jan. 2023), 657–668. https://doi.org/10.1109/tpami.2022.3154319
- [3] Shaked Brody, Uri Alon, and Eran Yahav. 2022. How Attentive are Graph Attention Networks?. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net. https://openreview.net/forum?id=F72ximsx7C1
- [4] Dexiong Chen, Leslie O'Bray, and Karsten M. Borgwardt. 2022. Structure-Aware Transformer for Graph Representation Learning. In International Conference on Machine Learning. ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA (Proceedings of Machine Learning Research, Vol. 162), Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (Eds.). PMLR, 3469–3489. https://proceedings.mlr.press/v162/chen22r.html
- [5] Jinsong Chen, Kaiyuan Gao, Gaichao Li, and Kun He. 2023. NAGphormer: A Tokenized Graph Transformer for Node Classification in Large Graphs. arXiv:2206.04910 [cs.LG]
- [6] Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. 2020. Simple and Deep Graph Convolutional Networks. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event (Proceedings of Machine Learning Research, Vol. 119). PMLR, 1725–1735. http://proceedings.mlr.press/v119/chen20v.html
- [7] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-GCN: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. In Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '19). ACM. https://doi.org/10.1145/3292500.3330925
- [8] Eli Chien, Jianhao Peng, Pan Li, and Olgica Milenkovic. 2021. Adaptive Universal Generalized PageRank Graph Neural Network. In 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net. https://openreview.net/forum? id=n6il7fLxrP
- [9] Austin Derrow-Pinion, Jennifer She, David Wong, Oliver Lange, Todd Hester, Luis Perez, Marc Nunkesser, Seongjae Lee, Xueying Guo, Brett Wiltshire, Peter W. Battaglia, Vishal Gupta, Ang Li, Zhongwen Xu, Alvaro Sanchez-Gonzalez, Yujia Li, and Petar Velickovic. 2021. ETA Prediction with Graph Neural Networks in Google Maps. In Proceedings of the 30th ACM International Conference on Information and Knowledge Management (CIKM '21). ACM. https://doi.org/10.1145/3459637.3481916
- [10] Vijay Prakash Dwivedi and Xavier Bresson. 2020. A Generalization of Transformer Networks to Graphs. https://arxiv.org/abs/2012.09699

- [11] Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. 2020. Benchmarking Graph Neural Networks. https://arxiv.org/abs/2003.00982
- [12] Vijay Prakash Dwivedi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. 2022. Graph Neural Networks with Learnable Structural and Positional Representations. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net. https://openreview.net/forum?id=wfTjnvGphYj
- [13] Vijay Prakash Dwivedi, Ladislav Rampášek, Mikhail Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu, and Dominique Beaini. 2022. Long Range Graph Benchmark. https://arxiv.org/abs/2206.08164
- [14] Wenzheng Feng, Yuxiao Dong, Tinglin Huang, Ziqi Yin, Xu Cheng, Evgeny Kharlamov, and Jie Tang. 2022. GRAND+: Scalable Graph Random Neural Networks. arXiv:2203.06389 [cs.LG]
- [15] Wenzheng Feng, Jie Zhang, Yuxiao Dong, Yu Han, Huanbo Luan, Qian Xu, Qiang Yang, Evgeny Kharlamov, and Jie Tang. 2020. Graph Random Neural Network for Semi-Supervised Learning on Graphs. https://arxiv.org/abs/2005.11079
- [16] Matthias Fey and Jan Eric Lenssen. 2019. Fast Graph Representation Learning with PyTorch Geometric. https://arxiv.org/abs/1903.02428
- [17] Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. 2017. Neural Message Passing for Quantum Chemistry. In Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017 (Proceedings of Machine Learning Research, Vol. 70), Doina Precup and Yee Whye Teh (Eds.). PMLR, 1263–1272. http://proceedings.mlr.press/v70/gilmer17a.html
- [18] William L. Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 1024–1034. https://proceedings.neurips.cc/paper/2017/hash/5dd9db5e033da9c6fb5ba83c7a7ebea9-Abstract.html
- [19] Tiantian He, Yew Soon Ong, and Lu Bai. 2021. Learning Conjoint Attentions for Graph Neural Nets. In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 2641–2653. https://proceedings.neurips.cc/paper/2021/hash/1587965fb4d4b5afe8428a4a024feb0d-Abstract.html
- [20] Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. 2023. Harnessing Explanations: LLM-to-LM Interpreter for Enhanced Text-Attributed Graph Representation Learning. https://arxiv.org/abs/2305.19523
- [21] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open Graph Benchmark: Datasets for Machine Learning on Graphs. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). https://proceedings.neurips.cc/paper/2020/hash/fbo0d411a5c5b72b2e7d3527cfc84fd0-Abstract.html
- [22] Md Shamim Hussain, Mohammed J. Zaki, and Dharmashankar Subramanian. 2022. Global Self-Attention as a Replacement for Graph Convolution. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '22). ACM. https://doi.org/10.1145/3534678. 3539296
- [23] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net. https://openreview.net/forum? id=SJU4ayYgl
- [24] Devin Kreuzer, Dominique Beaini, William L. Hamilton, Vincent Létourneau, and Prudencio Tossou. 2021. Rethinking Graph Transformers with Spectral Attention. In Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 21618–21629. https://proceedings.neurips.cc/paper/2021/hash/b4fd1d2cb085390fbbadae65e07876a7-Abstract.html
- [25] Soo Yong Lee, Fanchen Bu, Jaemin Yoo, and Kijung Shin. 2023. Towards Deep Attention in Graph Neural Networks: Problems and Remedies. https://arxiv.org/abs/2306.02376
- [26] Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. 2020. Distance Encoding: Design Provably More Powerful Neural Networks for

- Graph Representation Learning. In Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). https://proceedings.neurips.cc/paper/2020/hash/2f73168bf3656f697507752ec592c437-Abstract.html
- [27] Meng Liu, Hongyang Gao, and Shuiwang Ji. 2020. Towards Deeper Graph Neural Networks. In Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '20). ACM. https://doi.org/10.1145/3394486.3403076
- [28] Liheng Ma, Chen Lin, Derek Lim, Adriana Romero-Soriano, Puneet K. Dokania, Mark Coates, Philip Torr, and Ser-Nam Lim. 2023. Graph Inductive Biases in Transformers without Message Passing. arXiv:2305.17589 [cs.LG]
- [29] Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. 2019. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. In The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019. AAAI Press, 4602–4609. https://doi.org/10.1609/aaai.v33i01.33014602
- [30] Kenta Oono and Taiji Suzuki. 2020. Graph Neural Networks Exponentially Lose Expressive Power for Node Classification. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net. https://openreview.net/forum?id=S1ldO2EFPr
- [31] Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. 2020. Geom-GCN: Geometric Graph Convolutional Networks. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net. https://openreview.net/forum?id=51e2agrFvS
- [32] Jiezhong Qiu, Jian Tang, Hao Ma, Yuxiao Dong, Kuansan Wang, and Jie Tang. 2018. DeepInf: Social Influence Prediction with Deep Learning. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '18). ACM. https://doi.org/10.1145/3219819.3220077
- [33] Ladislav Rampášek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. 2023. Recipe for a General, Powerful, Scalable Graph Transformer. arXiv:2205.12454 [cs.LG]
- [34] Jonathan M. Stokes, Kevin Yang, Kyle Swanson, Wengong Jin, and James J. Collins. 2020. A Deep Learning Approach to Antibiotic Discovery. *Cell* 181, 2 (2020), 475–483.
- [35] Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M. Bronstein. 2022. Understanding oversquashing and bottlenecks on graphs via curvature. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022. OpenReview.net. https://openreview.net/forum?id= TUmiRGzp-A
- [36] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In 6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings. OpenReview.net. https://openreview.net/forum?id=rIXMpikCZ
- [37] Yaoke Wang, Yun Zhu, Wenqiao Zhang, Yueting Zhuang, Yunfei Li, and Siliang Tang. 2024. Bridging Local Details and Global Context in Text-Attributed Graphs. https://arxiv.org/abs/2406.12608
- [38] Qitian Wu, Wentao Zhao, Zenan Li, David Wipf, and Junchi Yan. 2023. NodeFormer: A Scalable Graph Structure Learning Transformer for Node Classification. https://arxiv.org/abs/2306.08385
- [39] Qitian Wu, Wentao Zhao, Chenxiao Yang, Hengrui Zhang, Fan Nie, Haitian Jiang, Yatao Bian, and Junchi Yan. 2023. SGFormer: Simplifying and Empowering Transformers for Large-Graph Representations. https://arxiv.org/abs/2306.10759
- [40] Yujie Xing, Xiao Wang, Yibo Li, Hai Huang, and Chuan Shi. 2024. Less is More: on the Over-Globalizing Problem in Graph Transformers. https://arxiv.org/abs/2405.01102
- [41] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. 2019. How Powerful are Graph Neural Networks?. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net. https://openreview.net/forum?id=ryGs6iA5Km
- [42] Zhilin Yang, William W. Cohen, and Ruslan Salakhutdinov. 2016. Revisiting Semi-Supervised Learning with Graph Embeddings. In Proceedings of the 33nd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016 (JMLR Workshop and Conference Proceedings, Vol. 48), Maria-Florina Balcan and Kilian Q. Weinberger (Eds.). JMLR.org, 40-48. http://proceedings.mlr.press/v48/yanga16.html

- [43] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. 2021. Do Transformers Really Perform Bad for Graph Representation? arXiv:2106.05234 [cs.LG]
- [44] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor K. Prasanna. 2020. GraphSAINT: Graph Sampling Based Inductive Learning Method. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net. https://openreview.net/forum?id=BJe8pkHFwS
- [45] Kai Zhang, Yaokang Zhu, Jun Wang, and Jie Zhang. 2020. Adaptive Structural Fingerprints for Graph Attention Networks. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net. https://openreview.net/forum?id=BJsWx0NYPr
- [46] Yun Zhu, Jianhao Guo, Fei Wu, and Siliang Tang. 2022. RoSA: A Robust Self-Aligned Framework for Node-Node Graph Contrastive Learning. In Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22, Lud De Raedt (Ed.). International Joint Conferences on Artificial Intelligence Organization, 3795–3801. https://doi.org/10. 24963/ijcai.2022/527 Main Track.
- [47] Yun Zhu, Haizhou Shi, Zhenshuo Zhang, and Siliang Tang. 2024. MARIO: Model Agnostic Recipe for Improving OOD Generalization of Graph Contrastive Learning. In Proceedings of the ACM on Web Conference 2024 (Singapore, Singapore) (WWW '24). Association for Computing Machinery, New York, NY, USA, 300–311. https://doi.org/10.1145/3589334.3645322
- [48] Yun Zhu, Yaoke Wang, Haizhou Shi, and Siliang Tang. 2024. Efficient Tuning and Inference for Large Language Models on Textual Graphs. In Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI-24, Kate Larson (Ed.). International Joint Conferences on Artificial Intelligence Organization, 5734–5742. https://doi.org/10. 24963/ijcai.2024/634 Main Track.
- [49] Yun Zhu, Yaoke Wang, Haizhou Shi, Zhenshuo Zhang, Dian Jiao, and Siliang Tang. 2024. GraphControl: Adding Conditional Control to Universal Graph Pre-trained Models for Graph Domain Transfer Learning. In Proceedings of the ACM on Web Conference 2024 (Singapore, Singapore) (WWW '24). Association for Computing Machinery, New York, NY, USA, 539–550. https://doi.org/10.1145/3589334.3645439

Dataset	#Nodes	#Edges	#Features	#Classes	Split(%)
Pubmed	19,717	44,324	500	3	0.3/2.5/5.0
Citeseer	3,327	4,552	3,703	6	5.2/18/37
Cora	2,708	5,728	1,433	7	3.6/15/30
Arxiv	169,343	1,166,243	128	40	53.7/17.6/28.7
Texas	183	279	1,703	5	48/32/20
Cornell	183	277	1,703	5	48/32/20
Wisconsin	251	450	1,703	5	48/32/20
Actor	7,600	26,659	932	5	48/32/20
AMiner-CS	593,486	6,217,004	100	18	20/30/50
Amazon2M	2,449,029	61,859,140	100	47	20/30/50

Table 7: Details of node classification datasets.

Table 8: Details of graph classification datasets.

Dataset	#Graphs	Avg.#Nodes	Avg.#Edges	Prediction task	#Classes	Metric	Split(%)
ZINC	12,000	23	25	regression	-	Mean Abs. Error	83.3/8.3/8.3
MINST	70,000	71	565	classif.	10	Accuracy	78.6/7.1/14.2
CIFAR10	60,000	118	941	classif.	10	Accuracy	75/8.3/16.7
Peptides-func	15,535	151	307	classif.	10	Avg. Precision	70/15/15
Peptides-struct	15,535	151	307	regression	11	Mean Abs. Error	70/15/15

A More Experiment details

A.1 Details of Datasets

We employs 10 node classification benchmark datasets, among which 4 are homophilic, 4 are heterophilic, and 2 are large-scale datasets. on the other hand, we employs 5 graph classification benchmark datasets, among which 2 are long range graph benchmark. The train-validation-test splits utilized are those which are publicly accessible. Details of these benchmark are shown in Table 7 and Table 8.

The Pubmed, Citeseer, and Cora datasets are citation networks [42]. Each node represents a research article and two nodes are adjacent if there is a citation between two articles. The node attributes are the bag-of-words features, and the node label is the category of the research domain of the article

The ogbn-arxiv [21] dataset constitutes a directed graph typifying the citation network amid all Computer Science (CS) arXiv manuscripts as cataloged by the Microsoft Academic Graph (MAG). Each vertex within this graph is an arXiv document, while each directed edge signifies a citation from one document to another. Accompanying each paper is a 128-dimensional feature vector, derived from the mean of the embeddings corresponding to words in the document's title and abstract sections.

The Texas, Cornell, and Wisconsin datasets are extracted from the WebKB dataset [31]. Each node represents a webpage and two nodes are adjacent if there is a hyperlink between the two webpages. The node attributes are the bag-of-words features, and the node label is the category of the webpage.

The actor dataset is the actor-only induced subgraph of a film-director-actor-writer network obtained from Wikipedia webpages [31]. Each node represents an actor and two nodes are adjacent if the two corresponding actors appear on the same Wikipedia webpage. The node features are derived from the keywords on the Wikipedia webpage of the corresponding actor, and the node label is determined by the words on the webpage.

AMiner-CS [15] are citation networks in which nodes represent papers and edges represent citations. Amazon2M [7] are co-purchase networks, where nodes indicate goods and edges indicate that the two connected goods are frequently bought together. The splits of large-scale datasets are followed the settings from [14].

ZINC [11] consists of 12K molecular graphs from the ZINC database of commercially available chemical compounds. These molecular graphs consist of 9 to 37 nodes that each represents a heavy atom (28 possible atom types) and each edge represents a bond (3 possible types). The task is to regress constrained solubility (logP) of the molecule.

MNIST and CIFAR10 [11] are derived from like-named image classification datasets by constructing an 8 nearest-neighbor graph of SLIC superpixels for each image.

Peptides-func and Peptides-struct [13] are composed of atomic graphs of peptides retrieved from SATPdb. In Peptides-func the prediction is multi-label graph classification into 10 nonexclusive peptide functional classes. While for Peptides-struct the task is graph regression of 11 3D structural properties of the peptides.

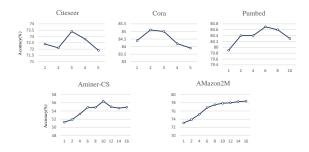


Figure 5: The results of experiments for the hyperparameter K

A.2 Details of hyperparameter

Experimental results are reported on the hyperparameter settings below, where we choose the settings that achieve the highest performance on the validation set. We choose hyperparameter grids that do not necessarily give optimal performance, but hopefully cover enough regimes so that each model is reasonably evaluated on each dataset.

- $lr \in \{5e-2, 1e-2, 5e-3, 1e-3, 5e-4\}$
- $K \in \{2, 3, 4, 6, 8, 12\}$
- pe_dim $(d_p) \in \{2, 4, 6, 8, 12, 16\}$
- se_dim $(d_s) \in \{2, 4, 6, 8, 12, 16\}$
- $ae_dim(d) \in \{32, 64, 128, 256, 512\}$
- pe&se_hidden_dim $(d') \in \{2, 4, 6, 8, 12, 16\}$
- ae_hidden_dim $(d'') \in \{32, 64, 128, 256, 512\}$
- dropout $\in \{0, 0.1, 0.2, 0.3, 0.5, 0.8\}$
- weight_decay $\in \{1e 2, 5e 3, 1e 3, 5e 4, 1e 4\}$
- activation ∈ {elu, relu, leakyrelu}
- layer_num $\in \{1, 2, 3, 4, 5, 6, 7, 8\}$

Parameter study on K. We conducted experiments on 3 small datasets and 2 large datasets for the hyperparameter K, e.g., the dim of initial SE, PE, and the numeber of sampled long-range nodes of each node. As shown in Figure 5, we can observe that the values of K are different for each dataset to achieve the optimal performance, since different graph exhibit different importance for positional information, structural information, and long-range dependency (global information).

For smaller datasets such as Citeseer, Cora, and Pubmed, optimal performance is achieved when K takes on a small value, as a small K is sufficient to capture the necessary field of view for positional and structural information. However, as K increases up to the graph diameter, nodes gain a perspective of nearly the entire graph, which results in a sharp drop in performance akin to over-smoothing. On the other hand, in homophilic graphs, the emphasis is on local information, thus a small K is enough to capture sufficient long-distance dependencies.

For large datasets e.g.Aminer-CS and AMazon2M, larger K is necessary to achieve optimal performance. This is because it is crucial for large graphs to capture long-range dependencies and maintain a broader receptive field for positional and structural information.

B Details of Initialization Strategies of PE/SE

In this section, we provide a detailed description of the encoding initialization strategies discussed in Section 1.

1) Strategies for Structural Encodings. As discussed in Section 3.1.1, structural encodings should capture the ability to perceive nodes' surrounding structure. In our method, we provide several strategies to achieve this objective: (1) Random-Walk Structural Encoding (RWSE), (2) Node Degree Encoding (DSE), and (3) Topology Counting Encoding (TCSE).

Specifically, RWSE is computed by the probability of arriving at the node itself with 0 step to K-1 step wandering. It reflect pure structure information of topology information of the K-hop receptive field.

$$s_i^{RWSE} = [I, \hat{A}, \hat{A}^2, \cdots \hat{A}^{K-1}]_{i,i} \in \mathbb{R}^K$$
 (11)

DSE is computed by the indegree and outdegree of the node, while TCSE is computed by counting the topological structure of its K-hop subgraph, such as the count of triangles, quads, and rings.

2) Strategies for Positional Encodings. To introduce the ability for a node to perceive its position relative to other specific nodes, we provide several off-the-shelf strategies that can reflect a node's distance or intersection relationships with other specific nodes: (1) Laplacian Positional Encoding (LapPE), (2) Relative Random Walk Probabilities (RWPE), and (3) Jaccard Encoding (JaccardPE).

Specifically, LapPE is a general method to encode node positions. For each node, its Laplacian PE is the K smallest non-trivial eigenvectors.

$$p_i^{LapPE} = [\phi_{0,i}, \phi_{1,i}, \phi_{2,i}, \cdots \phi_{K-1,i}] \in \mathbb{R}^K$$
 (12)

where $\phi_{m,i}$ is the *i*-th row of normalized eigenvector associated to the *m*-th lowest eigenvalue λ_m .

RWPE is computed by the probability of one node arriving other nodes within *K* steps. It reflect the positional interaction of a node with other nodes in the graph.

$$\check{p}_i^{RWPE} = [I, \hat{A}, \hat{A}^2, \cdots \hat{A}^{K-1}]_i \in \mathbb{R}^{K \times N}
p_i^{RWPE} = \mathcal{F}(\check{p}_i^{RWPE})$$
(13)

where $\mathcal{F}(\cdot):\mathbb{R}^{K\times N}\mapsto\mathbb{R}^K$ is a encoder to condense information.

Jaccard coefficient is used to compare similarities between different sets. It is the ratio of the size of the intersection of A and B to the size of the concatenation of A and B. In graph domain, we can use Jaccard coefficient as positional similarity between global node pairs. As follows:

$$J(p_i, p_j) = \frac{\sum_{k \in (V_i \cup V_j)} \min(p_{ik}, p_{jk})}{\sum_{k \in (V_i \cup V_j)} \max(p_{ik}, p_{jk})}$$
(14)

Similarly, we would like to represent our positional attention in terms of intersection, concatenation, and importance relationships between nodes, thus we use the shortest distance encoding as JaccardPE, and adopt the weight of node pairs with a Gaussian decay, i.e., $p_{i,j}^{Jaccard} = exp(-\frac{distance(i,j)^2}{2h^2})$. The algorithm is as follows:

Algorithm 1 Fast Jaccard encoding computing

Input: Normalized adjacency matrix \hat{A} ; Positional receptive field K

```
Output: Initialized positional encoding P^{Jaccard}

1: Initialization: P^{Jaccard} \leftarrow 0, B = I

2: for k = 0 to K - 1 do

3: for i, j in range(N) do

4: if B_{i,j} > 0 and P^{Jaccard}_{i,j} = 0 then

5: P^{Jaccard}_{i,j} = exp(-\frac{k^2}{2})

6: end if

7: end for

8: B = B\tilde{A}

9: end for

10: return P^{Jaccard}
```

C Expressive power of DeGTA

1-Weisfeiler-Leman (1-WL) Test and MPNNs. The 1-WL test is a node-coloring algorithm within the hierarchy of WL heuristics used for graph isomorphism. It operates by iteratively updating node colors based on their 1-hop local neighborhoods until no further changes occur in the node colors. The limitations of MPNNs in distinguishing non-isomorphic graphs were rigorously analyzed in the work of [41]. This analysis highlights the established equivalence between MPNNs and the 1-WL isomorphism test, as detailed by [29]. As a result, MPNNs may perform poorly on graphs that exhibit multiple symmetries within their original structure, including node and edge isomorphisms.

Expressive power of DeGTA. In the DeGTA architecture, we employ decoupled structural and positional encodings, both of which offer greater expressive power than the 1-WL test through the global positional and structural attention mechanism, as demonstrated in [24, 33]. Thus, DeGTA demonstrates excellent performance on graphs that exhibit multiple symmetries within their original structure, including node and edge isomorphisms. It is capable of distinguishing any pair of non-isomorphic graphs as well as CSL graphs, which cannot be learned by the 1-WL test or MPNNs. Further, compared to the GT with SE/PE, through the decoupling of PE, SE and AE, DeGTA can distinguish some graphs with more sensitive positional, structural and attribute information that coupled PS/SE cannot learn, as shown in case study.

Case study. As shown in Figure 6, we present additional visualization examples of the molecular graph propagation process to demonstrate the necessity of decoupling positional, structural, and attribute encodings. In these experimental cases, the DeGTA layer often produces correct results in scenarios where the coupled encoding approach yields incorrect outcomes.

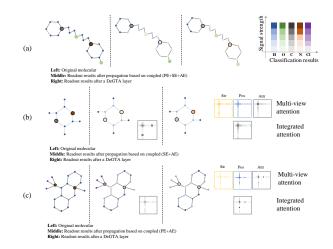


Figure 6: Our case study of molecular classification. For each case (left), we focus on pairs of nodes highlighted by a black border. These node pairs will produce incorrect classification results when using coupled encoding for global attention (middle), whereas DeGTA can produce correct results through decoupled multi-view attention and adaptive aggregation (right).

For instance, in case (b), we observe that the coupling of SE and AE results in high attention between the highlight pair, leading to erroneous classification outcomes (middle). However, with multi-view decoupling, we see that the pair exhibit high structural attention but low attribute and positional attention. This allows adaptive integration to learn a low overall attention between the two nodes, effectively suppressing message passing and yielding the correct classification (right). In case (c), the coupling of PE and AE results in low attention between the pair, again leading to incorrect classification (middle). In contrast, the decoupled attention shows that the node pairs possess high structural and attribute attention but low positional attention. The high structural attention activates long-range sampling between the pair, and the adaptive integration demonstrates a high importance to attribute information, resulting in a high overall attention between the two nodes and achieving the correct classification (right).

D Over-smoothing analysis of DeGTA

In this section, we provide both experimental and theoretical analyses of DeGTA's ability to mitigate the over-smoothing problem. Firstly, we introduce a quantitative metric for smoothness. We use the average Euclidean distance of the node attributes as a measure of smoothness and calculate this metric

at each layer as follows:

$$sm(G) = \frac{1}{n} \sum_{i \in V} \frac{1}{n-1} \sum_{j \in (V-i)} D(x_i, x_j)$$

$$= \frac{1}{n} \sum_{i \in V} \frac{1}{n-1} \sum_{j \in V} \frac{1}{2} \left\| \frac{x_i}{\|x_i\|} - \frac{x_j}{\|x_j\|} \right\|$$

$$= \frac{1}{2n(n-1)} \sum_{i} \sum_{j} \left\| \frac{x_i}{\|x_i\|} - \frac{x_j}{\|x_j\|} \right\|$$

Here, $\|\cdot\|$ denotes the Euclidean norm. sm(G) is negatively related to the overall smoothness of node attributes in graph G. A larger decline as the number of layers increases indicates that the model is less resistant to over-smoothing.

Experimental results. We conducted experiments on three datasets, evaluating accuracy and smoothness of DeGTA and GT as the number of layers increases. As shown in Figure 7, the experimental results demonstrate that DeGTA significantly enhances model performance and alleviates over-smoothing at each layer compared to GT, showing that the hard sampling strategy of DeGTA effectively captures the important longrange dependencies while avoiding the harm associated with the high degree of freedom in global information propagation of GT.

Theoretical analysis. In the following, we prove theoretically that DeGTA is further away from over-smoothing than GT when have equal attention to long-range dependencies.

First, we give the overall node update formula for DeGTA, where N(i) is the set of neighboring nodes and K(i) is the set of sampled nodes. For simplicity, we omit the details of the model here and frozen the adaptive integration of local-global. Ultimately, our node update can be written as $A_{Ni}X + A_{Si}X$, where A_{Ni} only pays attention to neighboring nodes and A_{Si} only pays attention to sampled nodes.

$$\hat{x}_i = \sum_{j \in N(i)} a_{ij} x_j + \sum_{j \in K(i)} a_{ij} x_j$$

There is a normalization method so that each row of A is equal to 1, which leads to that when GT and DeGTA have the same attention to important long-range dependencies, the nodes of DeGTA have relatively higher attention to themselves as well as to their neighboring nodes than GT, focusing on important long-range dependencies while better preserving

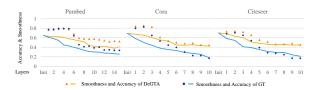


Figure 7: The results of the smoothness of DeGTA compared with GT of each layer. DeGTA demonstrates higher performance while maintaining further away from over-smoothing at each layer.

their own information and avoiding the phenomenon of over-smoothing.

$$sm(\hat{G})_{DeGTA} = \sum_{i,j} \left\| \frac{\hat{x}_i}{||\hat{x}_i||} - \frac{\hat{x}_j}{||\hat{x}_j||} \right\|_{DeGTA}$$

$$= \sum_{i,j} \left\| \hat{x}_i - \hat{x}_j \right\|_{DeGTA}$$

$$= \sum_{i,j} \sum_{x} \left\| \hat{x}_{ix} - \hat{x}_{jx} \right\|_{DeGTA}$$

$$= \sum_{i,j} \sum_{x} \left\| \sum_{y \in N(x)} x_{iy} A_{Nyx} + \sum_{y \in K(x)} x_{iy} A_{Syx} - \sum_{y \in N(x)} x_{jy} A_{Nyx} - \sum_{y \in K(x)} x_{jy} A_{Syx} \right\|$$

$$= \sum_{i,j} \sum_{x} \left\| \sum_{y \in N(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx} + \sum_{i,j} \sum_{x} \sum_{y \in K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$\geq \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$\geq \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{jy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{iy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{iy}) A_{Nyx}$$

$$= \sum_{i,j} \sum_{x} \sum_{x} \sum_{y \in N(x) \cup K(x)} (x_{iy} - x_{iy}) A_{Nyx}$$

$$= \sum_{x} \sum_{x$$

The inequality holds when the nodes of DeGTA and GT have the same attention to long-range dependencies in K(i), the sm metric of DeGTA is larger than that of GT, proving that our DeGTA captures the same long-range dependencies while being further away from over-smoothing than GT.