

TOWARDS SCALABLE DISTANCE-ENHANCED GRAPH NEURAL NETWORK

Anonymous authors

Paper under double-blind review

ABSTRACT

Graph neural networks (GNNs) have demonstrated significant advantages in graph mining tasks, but often suffer from limited expressive power. Among existing expressive GNNs, distance-enhanced GNNs (DE-GNNs) arise as promising ones due to their conceptual simplicity and alignment with the expressive needs of real-world applications. However, scalability remains a key challenge for DE-GNNs, as constructing pairwise distance features requires quadratic complexity. Additionally, while existing work has shown that specialized distance features enable strong expressiveness, the expressive power of simpler distance metrics remains less understood. In this paper, we propose a new **Scalable Distance-Enhanced Graph Neural Network** (termed SDE-GNN) to tackle the above issues. SDE-GNN introduces a distance-aware message-passing framework, where message weights are computed by a learnable distance feature mapping. It first linearly projects the adjacency-power-based distance vector to a scalar, then applies a polynomial expansion. To efficiently scale to large graphs, we reformulate the distance features as the product of two asymmetric node encodings and apply Randomized SVD for dimensionality reduction, lowering the computational complexity from quadratic in the number of nodes to linear in the number of edges. Additionally, we leverage the sparsity of the adjacency matrix to directly compute the first-order term of the distance feature mapping, further mitigating distortion from dimensionality reduction. Theoretically, we show that the adopted adjacency-power-based distance outperforms other commonly used distance features. Empirically, we conduct experiments on 17 datasets and verify the effectiveness, efficiency, and scalability of SDE-GNN.

1 INTRODUCTION

Graph-structured data, which represents entities and their relationships as nodes and edges, is ubiquitous across domains ranging from social networks to natural sciences. Although Graph Neural Networks (GNNs) (Kipf & Welling, 2017; Hamilton et al., 2017; Velickovic et al., 2018) have achieved remarkable success in graph mining tasks, it is widely acknowledged that standard GNNs have suffered from limited expressive power. For instance, they are constrained by the 1-Wefeiler-Lehman test in their ability to distinguish graph isomorphisms (Xu et al., 2019), and are incapable of counting certain simple substructures (Chen et al., 2020), and fail to distinguish specific node pairs (Zhang & Chen, 2018). These limitations hinder their effectiveness in tasks that demand higher expressive capacity.

Improving the expressiveness of GNNs has been extensively explored in recent years (Morris et al., 2020b; 2019; Cotta et al., 2021; Bevilacqua et al., 2022; Abboud et al., 2021; Zhang et al., 2024; 2023). Among these approaches, distance-enhanced GNNs (DE-GNNs) have attracted remarkable attention due to their conceptual simplicity, which enhances the message-passing procedure by incorporating distance features. Specifically, when node u aggregates features from its neighbor node v , DE-GNNs concatenate an additional distance feature $d(u, v)$ with the original features transmitted from v to u . Promising results have been shown on the expressiveness of DE-GNNs, including the capability to distinguish graphs beyond 1-WL (Li et al., 2020), the ability to count specific graph substructures (Ma et al., 2023), the identification of critical nodes and edges (Zhang et al., 2023), and awareness of the affinity between nodes (Vellingker et al., 2023), which closely align with the demands of practical applications.

054 Despite recent progress, research on DE-GNNs remains unsatisfactory. One major limitation lies
 055 in scalability. As shown in Table 1, existing methods often require at least quadratic computational
 056 complexity to construct pairwise distance features among nodes, making them suitable only for graphs
 057 with fewer than a few thousand nodes. Besides, the theoretical understanding of DE-GNNs is also
 058 insufficient. Existing quantitative expressiveness analyses only demonstrate strong expressiveness for
 059 DE-GNN with specialized distances like eigen projection distance and resistance distance (Zhang
 060 et al., 2024; 2023). This naturally leaves a question regarding the expressiveness of DE-GNNs when
 061 equipped with more basic distance features, such as the power of the adjacency matrix.

062 To preserve the expressive power of DE-GNNs while ensuring scalability, we propose a **Scalable**
 063 **Distance-Enhanced Graph Neural Network**, termed SDE-GNN. We begin by introducing a distance-
 064 aware message-passing framework, where message weights between nodes are modulated by their
 065 pairwise distance features, defined as entries from powers of the adjacency matrix. Specifically, we
 066 map these features into scalar weights via a learnable distance feature mapping, which performs linear
 067 projection followed by a nonlinear polynomial expansion, effectively capturing complex patterns in
 068 the distance features. Instead of explicitly computing the distance feature mapping, we reformulate
 069 the distance features as the product of two asymmetric node encodings and then perform Randomized
 070 SVD (Halko et al., 2011) to compress the asymmetric node encodings, thereby reducing the overall
 071 complexity from quadratic in the number of nodes to linear in the number of nodes and edges.
 072 Additionally, we leverage the sparsity of the adjacency matrix to directly compute the first-order
 073 term of the distance feature mapping, further mitigating distortion from dimensionality reduction.
 074 To validate the effectiveness of the proposed method, we first show that the expressive power of the
 075 adopted adjacency-power-based distance—a widely used metric—theoretically upper-bounds that
 076 of the eigenspace projection distance, which has been previously shown to be more expressive than
 077 other commonly used distance features. Then, empirically, we evaluate SDE-GNN against 15 popular
 078 baselines on 17 widely used datasets. The experimental results verify the effectiveness, efficiency,
 and scalability of SDE-GNN. Our contributions can be summarized as follows:

- 079 • **Scalable DE-GNN.** By introducing decouplable polynomial distance encoding and adaptive
 080 dimensionality reduction mechanisms, we optimize the computational complexity of DE-
 081 GNNs from quadratic w.r.t number of nodes into linear w.r.t number of nodes and edges,
 082 making it scalable to larger graphs while preserving strong expressive power.
- 083 • **Theoretical Analysis.** We theoretically analyze the role of the adopted adjacency-power-
 084 based distance features in message passing, and prove that the expressive power of such
 085 distance features upper-bounds that of other commonly used ones.
- 086 • **Empirical Validation.** Extensive experiments on 17 benchmarked datasets demonstrate
 087 that SDE-GNN achieves superior performance compared to 15 popular baselines in terms of
 088 effectiveness, efficiency and scalability.

091 2 RELATED WORKS

092 Expressive graph neural networks have been a key focus in the graph learning community in recent
 093 years. Existing methods can broadly be divided into high-order GNNs, subgraph GNNs, and positional
 094 encoding-enhanced GNNs (Zhang et al., 2024). High-order GNNs (Morris et al., 2020b; 2019; Maron
 095 et al., 2019b;a) achieve expressiveness comparable to high-dimensional WL tests by performing
 096 message passing between node tuples but are hindered by prohibitive computational and storage costs.
 097 Subgraph GNNs (Cotta et al., 2021; Bevilacqua et al., 2022), on the other hand, offer a more practical
 098 approach by dividing the original graph into a collection of subgraphs and performing message
 099 passing on each. However, these methods often require a large number of subgraphs—typically
 100 equal to the number of nodes — which restricts their scalability to small graphs (Bevilacqua et al.,
 101 2024). Positional encoding-enhanced GNNs represent a conceptually simpler way to improve
 102 expressiveness, as they augment node and edge features with positional encodings without altering
 103 the model architecture, which can be roughly classified into absolute positional encodings (APEs)
 104 and relative positional encodings (RPEs). APEs, inspired by positional encodings in 1-D sequence
 105 models like Transformers, are identifier vectors of nodes to represent their positions within the graph.
 106 Examples include random node features (Abboud et al., 2021; Sato et al., 2021) and eigenvectors
 107 of the Laplacian matrix (Kreuzer et al., 2021; Dwivedi & Bresson, 2020). However, the absence of
 a canonical node ordering in graphs often results in ambiguity (Wang et al., 2022) or necessitates

complex preprocessing (Lim et al., 2023; Huang et al., 2024). More naturally for graphs, RPEs capture pairwise relationships between nodes, serving as additional edge features, with examples including random walk probabilities (Ma et al., 2023), shortest path distance (Li et al., 2020), resistance distance (Velingker et al., 2023), and eigenspace projection distance (Zhang et al., 2024). We refer to GNNs equipped with RPEs as distance-enhanced GNNs in this paper. One limitation of DE-GNNs is that constructing pairwise distance features requires quadratic computational complexity, which limits their application to small graphs. To address this limitation, we propose a scalable DE-GNN that exhibits linear complexity and is suitable for large-scale graphs.

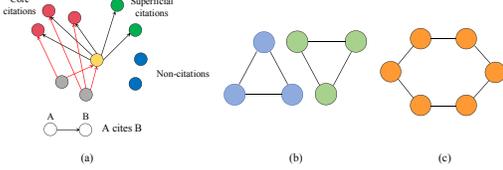


Figure 1: Examples to show the expressiveness of DE-GNNs. (a) depicts a real-world citation network with core citations, superficial citations, and non-citations. (b) and (c) illustrate two regular graphs, where the objective is to detect nodes within triangles.

Table 1: Time complexity comparison of distance-enhanced GNNs for computing all node representations, where n , e , and p denote the number of nodes, edges, and reduced dimension of our method (with $p \ll n$), respectively.

Model	Time Complexity
DE-GNN (Li et al., 2020)	$O(ne + n^2)$
Graphormer (Ying et al., 2021b)	$O(ne + n^2)$
RD-WL (Zhang et al., 2023)	$O(n^3)$
GRIT (Ma et al., 2023)	$O(ne + n^2)$
EPWL (Zhang et al., 2024)	$O(n^3)$
SDE-GNN (ours)	$O(ep + np^2)$

3 METHODOLOGY

In this section, we first introduce the overall framework of SDE-GNN in Section 3.1 and then show how to efficiently compute such a framework in Section 3.2. The final computation procedure of SDE-GNN is illustrated in Algorithm 1.

Algorithm 1 The Computation Procedure of SDE-GNN in Matrix Form.

Require: A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ with its normalized adjacency matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{n \times n}$ and node features $\mathbf{H} \in \mathbb{R}^{n \times d}$, reduced embedding dimension p , learned parameters $\mathbf{w} \in \mathbb{R}^{k+1}$, $\{\alpha_i\}_{i=1}^m$, θ^{f_m} .

▽ Dimensionality reduction

- 1: $\mathbf{U} \in \mathbb{R}^{n \times p}$, $\Sigma \in \mathbb{R}^{p \times p}$, $\mathbf{V} \in \mathbb{R}^{n \times p} \leftarrow \text{RandomSVD}(\tilde{\mathbf{A}})$
- 2: $\mathbf{E} \leftarrow [\mathbf{V} \parallel \mathbf{U}\Sigma \parallel \tilde{\mathbf{A}}\mathbf{U}\Sigma \parallel \dots \parallel \tilde{\mathbf{A}}^{k-1}\mathbf{U}\Sigma] \in \mathbb{R}^{n \times (k+1)p}$

▽ Compute the first-order representations

- 3: $\mathbf{H}'_1 \leftarrow \alpha_1 \sum_{i=0}^k \mathbf{w}_i \tilde{\mathbf{A}}^i f_m(\mathbf{H})$

▽ Compute higher-order representations

- 4: $\mathbf{P} \leftarrow \sum_{u \in \mathcal{V}} \mathbf{V}_u \otimes f_m(\mathbf{H}_u) \in \mathbb{R}^{p \times d}$
- 5: $\mathbf{S} \leftarrow \sum_{i=0}^k \mathbf{w}_i \mathbf{E}^{(i)}$; # Let $\mathbf{E}^{(0)}, \dots, \mathbf{E}^{(k)} \in \mathbb{R}^{n \times p}$ be the $k+1$ matrices forming \mathbf{E}
- 6: $\mathbf{H}' \leftarrow \mathbf{H}'_1 + \sum_{j=2}^m \alpha_j (\mathbf{S}^{\odot j}) \mathbf{P}$
- 7: **Return:** \mathbf{H}' .

3.1 DISTANCE-AWARE MESSAGE PASSING FRAMEWORK

Overall framework. Motivated by prior work on distance-enhanced GNNs, we introduce a distance-aware message passing mechanism to enhance the expressive power of standard GNNs. Unlike conventional fixed and local schemes (Kipf & Welling, 2017; Rampásek et al., 2022), our framework performs global and adaptive aggregation: each node gathers information from all other nodes rather than its immediate neighbors, with each contribution modulated by a distance-based feature, allowing the model to selectively incorporate relevant signals during message passing. Formally, given a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, our message passing framework can be formulated as

$$\mathbf{h}_u^{(l+1)} = \sum_{v \in \mathcal{V}} f_e^{(l)}(\mathbf{d}_{u,v}) \cdot f_m^{(l)}(\mathbf{h}_v^{(l)}), \quad (1)$$

where $\mathbf{h}_v^{(l)} \in \mathbb{R}^d$ is the representation of node v at layer l , $\mathbf{h}_u^{(l+1)} \in \mathbb{R}^d$ is the representation of node u at layer $l + 1$, $\mathbf{d}_{u,v}$ is the distance features between u and v , $f_e^l(\cdot)$ is a distance learning function that maps the distance feature into a scalar to measure the affinity between u and v , and $f_m^{(l)}(\cdot)$ is the feature transformation function. For notation simplicity, we omit the superscripts of $(\cdot)^{(l)}$ and replace $(\cdot)^{(l+1)}$ with $(\cdot)'$ in the following statements.

Detailed configurations. The three key components of Equation (1) are the distance features $\mathbf{d}_{u,v}$, the distance encoding function $f_e(\cdot)$, and the feature transformation function $f_m(\cdot)$. As the primary objective of this paper is to efficiently leverage distance features to enhance GNN expressiveness, we adopt a simple MLP for $f_m(\cdot)$ and focus on the choices of $\mathbf{d}_{u,v}$ and $f_e(\cdot)$. For $\mathbf{d}_{u,v}$, we adopt a basic choice—the powers of the adjacency matrix. Formally, let $\tilde{\mathbf{A}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{-\frac{1}{2}}$ denote the normalized adjacency matrix, the distance feature is defined as

$$\mathbf{d}_{u,v} = [\tilde{\mathbf{A}}_{u,v}^0, \tilde{\mathbf{A}}_{u,v}^1, \dots, \tilde{\mathbf{A}}_{u,v}^k]^\top \in \mathbb{R}^{k+1}, \quad (2)$$

where k is a hyperparameter to control the power of the adjacency matrix. Here we adopt the power of the normalized adjacency matrix $\tilde{\mathbf{A}}^k$ rather than that of the original adjacency matrix (i.e., \mathbf{A}^k), as the entries of \mathbf{A}^k will grow rapidly with k . The motivation behind our choice of the adjacency-power-based distance is that it serves as the most basic distance measure in the graph, from which other distances may be derived. For instance, the shortest path distance between u and v is given by the smallest i for which the entry $\tilde{\mathbf{A}}_{u,v}^i$ is non-zero. Aligning with such intuition, in Section 4, we will theoretically prove that the expressive power of the adopted adjacency-power-based distance upper-bounds that of the eigenspace projection distance, which (Zhang et al., 2024) has shown to possess the strongest expressive power among other commonly used distances. For the distance learning function $f_e(\cdot)$, we define it as $f_e(\mathbf{d}_{u,v}) = \sigma(\mathbf{w}^\top \mathbf{d}_{u,v})$, where $\mathbf{w} \in \mathbb{R}^{k+1}$ is a learnable weight vector that maps the $(k+1)$ -dimensional adjacency-power-based distance vector to a scalar, bias term is omitted for brevity, and σ is a nonlinear function. Such a design of $f_e(\mathbf{d}_{u,v})$ enables the adaptive utilization of distinct components in distance features, thereby allowing the model to selectively integrate relevant signals during the message passing process. To characterize the nonlinear function in $f_e(\mathbf{d}_{u,v})$, we introduce a polynomial expansion scheme grounded in the Weierstrass approximation theorem, which guarantees that any continuous function defined on a closed interval can be uniformly approximated by polynomials to arbitrary precision (Pérez & Quintana, 2008). Specifically, we redefine $f_e(\mathbf{d}_{u,v})$ with a learnable polynomial expansion as

$$f_e(\mathbf{d}_{u,v}) = \sum_{i=1}^m \alpha_i (\mathbf{w}^\top \mathbf{d}_{u,v})^i, \quad (3)$$

where m is hyperparameter that denotes the maximum order of polynomial expansion, and $\{\alpha_i\}_{i=1}^m$ are learnable parameters to flexibly control the coefficients for each polynomial term.

Illustration of expressive advantages. In Figure 1, we present illustrative examples to demonstrate the expressive advantage of the distance features in capturing nuanced structural relationships, illustrated by the example of triangle detection. Specifically, Figure 1 (b) and (c) illustrate two graphs where vanilla GNNs (Kipf & Welling, 2017) fail to determine whether a node is part of a triangle due to identical local structures. In contrast, the proposed distance-enhanced method can distinguish such cases by identifying whether there exists another node that is both one-hop and two-hop away—an indicator of triangle membership and can be determined by the distance features between nodes. Figure 1 (a) further demonstrates how this triangle-detection capability can be used to differentiate important and less important neighbors in a real-world citation network, where each node represents a paper. In this example, the yellow node makes two types of citations: core citations (e.g., the red node, which forms a triangle with the yellow node) and superficial citations (which do not form such triangles). Due to space limitations, we leave more details of Figure 1 in Appendix A.

3.2 EFFICIENT IMPLEMENTATION OF DISTANCE-AWARE MESSAGE PASSING FRAMEWORK

One limitation of the distance feature mapping in Equation (3) is the requirement of explicitly computing the pairwise distances between all node pairs, which incurs a quadratic time complexity. This limits its scalability to large-scale graphs. To address this, we propose an efficient implementation of distance-aware message passing framework, including decouplable message passing framework and efficient dimensionality reduction mechanism.

Decouplable message passing framework. We first represent each node by two asymmetric encodings \mathbf{E}_u and \mathbf{e}_v derived from the adjacency-power-based distance features as

$$\mathbf{E}_u = \left[\tilde{\mathbf{A}}_{u,:}^0 \parallel \tilde{\mathbf{A}}_{u,:}^1 \parallel \cdots \parallel \tilde{\mathbf{A}}_{u,:}^k \right], \mathbf{e}_u = \tilde{\mathbf{A}}_{u,:}^0, \quad (4)$$

where $\mathbf{E}_u \in \mathbb{R}^{n \times (k+1)}$ comprises the distance features from node u to all other nodes, $\mathbf{e}_u \in \mathbb{R}^n$ is a unit vector extracted from $\tilde{\mathbf{A}}^0$, and \parallel denotes the concatenation operation. Since $\mathbf{d}_{u,v} = \mathbf{E}_u^\top \mathbf{e}_v$, we reformulate the distance feature mapping in Equation (3) and decouple it as

$$f_e(\mathbf{d}_{u,v}) = \sum_{i=1}^m \alpha_i (\mathbf{w}^\top \mathbf{E}_u^\top \mathbf{e}_v)^i = \sum_{i=1}^m \alpha_i (\mathbf{w}^\top \mathbf{E}_u^\top)^{\odot i} \mathbf{e}_v, \quad (5)$$

where $(\cdot)^{\odot i}$ denotes the Hadamard power of order i . The polynomial is decouplable as \mathbf{e}_v is a one-hot vector. Then the overall message passing framework can be formalized as

$$\mathbf{h}'_u = \sum_{v \in \mathcal{V}} f_e(\mathbf{d}_{u,v}) \cdot f_m(\mathbf{h}_v) = \sum_{i=1}^m \alpha_i (\mathbf{w}^\top \mathbf{E}_u^\top)^{\odot i} \cdot \sum_{v \in \mathcal{V}} \mathbf{e}_v \otimes f_m(\mathbf{h}_v), \quad (6)$$

where \otimes denotes the outer product of two vectors.

Efficient dimensionality reduction mechanism. To reduce the computational complexity of the message passing framework from quadratic to linear, we propose an efficient dimensionality reduction mechanism without compromising performance. Building upon the theoretical foundation that Singular Value Decomposition (SVD) provides the optimal low-rank approximation of a matrix in both Frobenius and spectral norms (Eckart & Young, 1936), we adopt its faster variant—Randomized SVD (Halko et al., 2011), to decompose the normalized adjacency matrix $\tilde{\mathbf{A}}$ for dimensionality reduction. Specifically, $\tilde{\mathbf{A}}$ is decomposed as $\tilde{\mathbf{A}} \approx \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$, where $\mathbf{U} \in \mathbb{R}^{n \times p}$, $\mathbf{\Sigma} \in \mathbb{R}^{p \times p}$, and $\mathbf{V} \in \mathbb{R}^{n \times p}$ are the factor matrices. The reduced dimensionality p is empirically set to 1,000 to strike a balance between efficiency and efficacy. Then the two asymmetric representations \mathbf{E}_u and \mathbf{e}_u can be redefined, and we take their matrix form as

$$\mathbf{E} = \left[\mathbf{V} \parallel \mathbf{U}\mathbf{\Sigma} \parallel \tilde{\mathbf{A}}\mathbf{U}\mathbf{\Sigma} \parallel \cdots \parallel \tilde{\mathbf{A}}^{k-1}\mathbf{U}\mathbf{\Sigma} \right], \mathbf{e} = \mathbf{V}, \quad (7)$$

where $\mathbf{E} \in \mathbb{R}^{n \times k \times p}$ and $\mathbf{e} \in \mathbb{R}^{n \times p}$. A key efficient computational technique is avoiding explicit $n \times n$ matrix construction by iteratively computing $\tilde{\mathbf{A}}^{k-1}\mathbf{U}\mathbf{\Sigma}$ from right to left—starting with $\mathbf{U}\mathbf{\Sigma}$, then $\tilde{\mathbf{A}}\mathbf{U}\mathbf{\Sigma}$, $\tilde{\mathbf{A}}^2\mathbf{U}\mathbf{\Sigma}$, till $\tilde{\mathbf{A}}^{k-1}\mathbf{U}\mathbf{\Sigma}$ —ensuring all intermediate results remain $n \times p$. Leveraging the sparsity of $\tilde{\mathbf{A}}$, the computational complexity of $\tilde{\mathbf{A}}^{k-1}\mathbf{U}\mathbf{\Sigma}$ is linear in $\text{nnz}(\tilde{\mathbf{A}})$, which corresponds to the number of edges e . Consequently, the complexity of constructing \mathbf{E} and \mathbf{e} is linear in e , and the total computational complexity of $\{\mathbf{h}'_u\}_{u \in \mathcal{V}}$ in Equation (6) is linear in e .

Separation of first- and higher-order polynomial terms. The dimensionality reduction of $\tilde{\mathbf{A}}$ inevitably incurs information loss, which may degrade model performance. To mitigate the information loss, we propose a mechanism that separates the first- and higher-order terms in the polynomial expansion, preserving the first-order term without dimensionality reduction. Specifically, the first-order term without dimensionality reduction in Equation (6) can be reformulated as:

$$\mathbf{h}'_{u,1} = \sum_{v \in \mathcal{V}} \alpha_1 (\mathbf{w}^\top \mathbf{E}_u^\top \mathbf{e}_v) \cdot f_m(\mathbf{h}_v) = \alpha_1 \sum_{i=0}^k w_i \sum_{v=1}^n \tilde{\mathbf{A}}_{u,v}^i f_m(\mathbf{h}_v), \quad (8)$$

where w_i denotes the i -th entry of \mathbf{w} and n denotes the number of nodes. By stacking the initial node features $\{\mathbf{h}_u\}_{u \in \mathcal{V}}$ into a matrix $\mathbf{H} \in \mathbb{R}^{n \times d}$, Equation (8) can be rewritten in matrix form as $\mathbf{H}'_1 = \alpha_1 \sum_{i=0}^k w_i \tilde{\mathbf{A}}^i f_m(\mathbf{H})$, where $f_m(\mathbf{H})$ denotes row-wise application of $f_m(\cdot)$, and $\mathbf{H}'_1 \in \mathbb{R}^{n \times d}$ represents the first-order node representations. Notably, this computation can be efficiently performed by exploiting the sparsity of $\tilde{\mathbf{A}}$, similar to Equation (7). Through this separation mechanism, we reformulate the message passing procedure in Equation (6) into a matrix form as

$$\mathbf{h}'_u = \mathbf{h}'_{u,1} + \sum_{j=2}^m \alpha_j (\mathbf{w}^\top \mathbf{E}_u^\top)^{\odot j} \cdot \sum_{v \in \mathcal{V}} \mathbf{e}_v \otimes f_m(\mathbf{h}_v) \quad (9)$$

Time complexity analysis. We analyse the time complexity of SDE-GNN and compare it with other distance-enhanced GNNs. The time complexity of Randomized SVD on $\tilde{\mathbf{A}}$ is $O(ep + np^2)$, that of constructing \mathbf{E} and e is $O(ep)$, and that of the message passing framework in Equation (9) is $O(np^2)$. Thus, the overall time complexity of SDE-GNN is $O(ep + np^2)$, where n and e denote the number of nodes and edges respectively, and p denotes the reduced dimension. Given that the number of layers L , the power of adjacency matrix k , the maximum order of polynomial expansion m , and the node feature dimension d are small constants, their contributions are omitted from the complexity analysis. We compare SDE-GNN’s time complexity with other distance-enhanced GNNs in Table 1.

4 THEORETICAL ANALYSIS

In this section, we compare the expressive power of the proposed method with existing DE-GNNs. Specifically, we adopt the *Generalized Distance Weisfeiler-Lehman* (GD-WL) test (Zhang et al., 2023) as a theoretical abstraction of various DE-GNNs and compare the corresponding GD-WL variants in terms of their ability to distinguish non-isomorphic graphs. For a given graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, GD-WL iteratively refines node colors according to the update rule

$$\chi_G^{t+1}(u) = \text{hash}(\{(\chi_G^t(v), d_G(u, v)) : v \in V\}), \quad (10)$$

where $\chi_G^t(u)$ denotes the color of node u at round t , $\text{hash}(\cdot)$ is an injective function that assigns a distinct color to each unique input, and $d_G(\cdot, \cdot)$ is a distance feature mapping that maps a pair of nodes to a distance feature, which will be described in detail later. After running GD-WL for a sufficiently large round T , it determines whether two graphs $\mathcal{G}_1(\mathcal{V}_1, \mathcal{E}_1)$ and $\mathcal{G}_2(\mathcal{V}_2, \mathcal{E}_2)$ are isomorphic by comparing the multisets of node colors $\{\{\chi_{G_1}^T(u) : u \in \mathcal{V}_1\}\}$ and $\{\{\chi_{G_2}^T(u) : u \in \mathcal{V}_2\}\}$. Intuitively, Equation (10) can be viewed as a global GNN augmented with distance features, where each node u updates its representation $\chi_G^{t+1}(u)$ by aggregating information from all nodes $v \in \mathcal{V}$ in the graph, with each message comprising both the representation $\chi_G^t(v)$ and the corresponding distance feature $d_G(u, v)$. By specifying different distance feature mapping d , Equation (10) can instantiate various DE-GNNs. For instance, for DE-GNNs that construct additional edge features $\mathbf{f}_{u,v} \in \mathbb{R}^D$ on existing edges, we can define the distance as $d_G(u, v) = [1, \mathbf{f}_{u,v}] \in \mathbb{R}^{D+1}$ for $\{u, v\} \in E$, and set $d_G(u, v)$ to a zero vector in \mathbb{R}^{D+1} otherwise. More details about the GD-WL can be found in the Appendix B.

Distance Feature Mapping. Given a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, the distance feature mapping d outputs a function $d_G : V \times V \rightarrow \mathbb{R}^{D_G}$ that maps a pair of nodes within V into a distance vector. The dimension of the distance vector D_G can vary across graphs. (Zhang et al., 2024) shows that GD-WL equipped with the eigenspace projection distance achieves the highest expressiveness among commonly used distances such as resistance distance and shortest path distance. To illustrate the advantages of our method over existing DE-GNNs, we adopt GD-WL with the adjacency-power-based distance d^{AP} as a proxy, and compare its expressiveness against that of GD-WL with the eigenspace projection distance d^{EP} . Specifically, given a graph $\mathcal{G}(\mathcal{V}, \mathcal{E})$, let $\tilde{\mathbf{L}} = \mathbf{D}^{-1/2} \mathbf{L} \mathbf{D}^{-1/2}$ denote the normalized Laplacian matrix and $\tilde{\mathbf{L}} = \sum_{i=1}^m \lambda_i \mathbf{P}_i$ denotes its spectral decomposition, where $\lambda_1, \dots, \lambda_m$ are the eigenvalues of $\tilde{\mathbf{L}}$ and $\mathbf{P}_1, \dots, \mathbf{P}_m \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ are the corresponding eigenspace projections. The eigenspace projection distance between two nodes $u, v \in V$ is

$$d_G^{EP}(u, v) = [\lambda_1, (P_1)_{u,v}, \lambda_2, (P_2)_{u,v}, \dots, \lambda_m, (P_m)_{u,v}] \in \mathbb{R}^{2m}, \quad (11)$$

where $(P_i)_{u,v}$ denotes the u -th row and v -th column of \mathbf{P}_i . Let $\tilde{\mathbf{A}} = \tilde{\mathbf{D}}^{-1/2} \mathbf{A} \tilde{\mathbf{D}}^{-1/2}$ be the normalized adjacency matrix of the given graph \mathcal{G} and $n = |\mathcal{V}|$ be the number of nodes. The adjacency-powered-based distance feature between two nodes u, v is defined as

$$d_G^{AP}(u, v) = [\tilde{\mathbf{A}}_{u,v}^0, \tilde{\mathbf{A}}_{u,v}^1, \dots, \tilde{\mathbf{A}}_{u,v}^n] \in \mathbb{R}^{n+1}, \quad (12)$$

where $\tilde{\mathbf{A}}^i$ denotes the i -th power of $\tilde{\mathbf{A}}$.

In the following parts, we first establish the connection between GD-WL with d^{AP} and d^{EP} using Theorem 1, and then demonstrate that GD-WL with d^{AP} bounds the expressiveness of GD-WL with d^{EP} using Theorem 2. The Proofs of Theorem 1 and Theorem 2 can be found in Appendix C.1 and Appendix C.2, respectively.

Theorem 1 Let $\sigma(\mathbf{M}) = \{\{\lambda_1, \dots, \lambda_n\}\}$ denote the multiset of eigenvalues of a matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$, and let $\tilde{\mathbf{A}}_G = \mathbf{D}_G^{-1/2} \mathbf{A}_G \mathbf{D}_G^{-1/2}$ be the normalized adjacency matrix of a graph \mathcal{G} . For any two

graphs $\mathcal{G}_1(\mathcal{V}_1, \mathcal{E}_1)$ and $\mathcal{G}_2(\mathcal{V}_2, \mathcal{E}_2)$, if they cannot be distinguished by GD-WL equipped with the adjacency-power-based distance d^{AP} , then the following must hold:

- $|\mathcal{V}_1| = |\mathcal{V}_2|$ and $\sigma(\tilde{\mathbf{A}}_{\mathcal{G}_1}) = \sigma(\tilde{\mathbf{A}}_{\mathcal{G}_2})$
- if $d_{\mathcal{G}_1}^{AP}(u, v) = d_{\mathcal{G}_2}^{AP}(x, y)$ then $d_{\mathcal{G}_1}^{EP}(u, v) = d_{\mathcal{G}_2}^{EP}(x, y)$

Remark 1. While the power of the adjacency matrix has been used as a distance feature in prior works, its role in message passing remains insufficiently understood. Unlike previous case-specific analyses (Ma et al., 2023; Li et al., 2020), Theorem 1 offers a principled understanding. Specifically, the first statement shows that if two graphs \mathcal{G}_1 and \mathcal{G}_2 cannot be distinguished by GD-WL with d^{AP} , then their adjacency spectra must be identical, i.e., $\sigma(\tilde{\mathbf{A}}_{\mathcal{G}_1}) = \sigma(\tilde{\mathbf{A}}_{\mathcal{G}_2})$. Therefore, the features (i.e., the final multiset of node colors) extracted by GD-WL with d^{AP} encode the graph spectra, which partition graphs into different equivalence classes that share the same adjacency spectrum. The second statement serves as a bridge connecting d^{AP} to the eigenspace projection distance d^{EP} , thereby laying the foundation for comparing the two GD-WL variants in Theorem 2.

Theorem 2 For any two graphs \mathcal{G}_1 and \mathcal{G}_2 that cannot be distinguished by GD-WL equipped with d^{AP} , they also cannot be distinguished by GD-WL equipped with d^{EP} .

Remark 2. From Theorem 2, we know that if two graphs cannot be distinguished by GD-WL with d^{AP} , then they also cannot be distinguished by GD-WL with d^{EP} . Conversely, those distinguishable by d^{EP} are also distinguishable by d^{AP} . This theorem plays a key role in characterizing the expressiveness of GD-WL with d^{AP} . It establishes that GD-WL with d^{AP} upper-bounds the expressiveness of GD-WL with other distance metrics, as d^{EP} is known to achieve the highest expressiveness among commonly used metrics. Consequently, GD-WL with d^{AP} inherits known expressiveness results of GD-WL, such as being strictly more powerful than 1-WL (Zhang et al., 2024) and can distinguish block cut-vertex/edge trees (Zhang et al., 2023). This positions d^{AP} within the broader landscape of expressiveness studies.

5 EXPERIMENTAL RESULTS

5.1 EXPERIMENT SETTING

We mainly perform experiments on 17 widely used graph datasets, including 14 node classification datasets and 3 graph classification datasets. We compare our model with 15 state-of-the-art baselines, including 4 classic graph neural networks, 4 efficient graph transformer models, and 7 expressive graph neural networks. Our code is released anonymously at <https://anonymous.4open.science/r/SDE-GNN-7853>. More information about datasets, baselines, and experimental settings like hyperparameter configurations can be found in Appendix D.

5.2 PERFORMANCE COMPARISON

Effectiveness. As shown in Table 2, 3, and 4, SDE-GNN model yields consistently strong performance across all 17 datasets with different properties, which achieve the best performance on 11 out of them, demonstrating the effectiveness of the proposed distance-enhanced message passing framework. Besides the above observations, specific strengths of SDE-GNN can also be concluded from the performance of different methods: 1) In Table 2, the rank of SDE-GNN is significantly lower than that of the second one, Exphormer, on the five datasets that exhibit complex structural patterns, empirically demonstrating its strong expressive power. 2) In Table 3, SDE-GNN achieves the best performance on two heterophilous datasets, indicating its ability to capture high-order structural patterns—a key factor in effectively modeling heterophilous graphs. 3) In Table 4, SDE-GNN achieves the highest average rank on three graph classification datasets, demonstrating its effectiveness in capturing the overall graph structure.

Efficiency and Scalability. For scalability, the results from Table 2 and Table 3 show that SDE-GNN successfully scales to large-scale graphs, with node counts ranging from over 100K to approximately 3 million (e.g., ogbn-arxiv, twitch-gamers, pokec, and snap-patents). While existing expressive GNNs can only run on small graphs with a few thousand nodes and will encounter the out-of-memory error

Table 2: Test performance on 5 complex structural datasets is presented as the mean \pm standard deviation over 5 runs with different random seeds. “OOM” indicates out of memory. Highlighted are the top **first** and **second** results. “Acc” denotes the test accuracy metric, and “Rank” denotes the average rank across different datasets.

Model	CLUSTER	PATTERN	USA-Airports	Europe-Airports	Brazil-Airports	Rank \downarrow
	Acc \uparrow					
GCN	67.46 \pm 0.84	84.37 \pm 0.03	64.03 \pm 0.96	60.34 \pm 2.53	75.23 \pm 2.14	7.80
GraphSAGE	64.28 \pm 0.75	82.46 \pm 0.12	63.46 \pm 1.13	58.42 \pm 3.16	77.77 \pm 3.32	9.20
GAT	69.53 \pm 0.53	78.68 \pm 0.10	63.98 \pm 2.06	59.67 \pm 3.44	77.89 \pm 3.45	7.60
GIN	65.31 \pm 0.96	84.79 \pm 0.03	64.04 \pm 0.88	59.12 \pm 1.84	75.56 \pm 2.56	7.80
GraphGPS	78.03 \pm 0.21	86.19 \pm 0.06	43.10 \pm 2.44	46.35 \pm 0.52	52.22 \pm 1.57	10.80
Expformer	78.07 \pm 0.04	86.74 \pm 0.02	60.96 \pm 1.22	56.20 \pm 1.03	76.67 \pm 1.57	6.60
SGFormer	59.76 \pm 0.86	83.31 \pm 0.15	61.53 \pm 1.11	50.95 \pm 1.40	69.78 \pm 3.37	12.20
GOAT	58.61 \pm 0.67	82.76 \pm 0.11	38.65 \pm 2.41	36.77 \pm 3.88	37.78 \pm 4.24	14.60
DE-GNN	76.56 \pm 0.92	85.21 \pm 0.08	64.16 \pm 0.95	60.69 \pm 2.73	76.47 \pm 2.03	5.80
Graphormer	76.79 \pm 1.13	85.58 \pm 0.16	63.74 \pm 1.03	58.72 \pm 2.87	74.31 \pm 2.31	7.60
ESAN	OOM	OOM	OOM	OOM	OOM	16.00
GRIT	79.43 \pm 0.33	86.62 \pm 0.08	63.36 \pm 1.53	53.29 \pm 2.76	71.11 \pm 2.38	7.20
RD-WL	78.22 \pm 0.98	86.35 \pm 0.18	66.53 \pm 1.24	61.48 \pm 3.03	77.28 \pm 3.75	3.20
SPE	70.37 \pm 0.45	85.73 \pm 0.05	56.28 \pm 2.47	51.37 \pm 3.53	68.79 \pm 5.03	10.60
NeuralWalker	78.19 \pm 0.20	86.98 \pm 0.02	54.78 \pm 2.62	56.20 \pm 2.47	73.33 \pm 1.26	7.40
SDE-GNN (ours)	78.68 \pm 0.24	86.82 \pm 0.03	67.57 \pm 1.03	63.02 \pm 0.84	82.22 \pm 1.42	1.40

Table 3: Test performance on 5 real-world medium-scale datasets and 4 large-scale datasets is presented as the mean \pm standard deviation over 5 runs with different random seeds. “OOM” indicates out of memory. Highlighted are the top **first** and **second** results. “Acc” denotes the test accuracy metric, and “Rank” denotes the average rank across datasets.

Model	Cora	CiteSeer	PubMed	Squirrel	Chameleon	ogbn-arxiv	twitch-gamers	pokec	snap-patents	Rank \downarrow
	Acc \uparrow									
GCN	85.10 \pm 0.48	72.94 \pm 0.41	80.96 \pm 0.49	45.04 \pm 1.59	45.55 \pm 3.35	73.11 \pm 0.23	63.39 \pm 0.34	85.20 \pm 0.23	46.17 \pm 0.18	2.22
GraphSAGE	83.68 \pm 0.50	72.04 \pm 0.56	78.66 \pm 0.50	40.56 \pm 1.46	44.08 \pm 4.62	72.98 \pm 0.27	62.18 \pm 0.42	84.78 \pm 0.27	45.62 \pm 0.34	6.33
GAT	82.70 \pm 0.53	70.96 \pm 0.68	80.42 \pm 0.66	40.78 \pm 1.95	41.27 \pm 4.06	73.26 \pm 0.25	60.23 \pm 0.53	79.35 \pm 0.65	43.76 \pm 1.14	7.00
GIN	83.82 \pm 0.62	71.43 \pm 0.70	79.70 \pm 0.59	43.24 \pm 1.74	44.26 \pm 4.28	71.68 \pm 0.31	61.85 \pm 0.47	78.56 \pm 0.42	45.88 \pm 0.71	6.11
GraphGPS	83.81 \pm 0.93	72.64 \pm 1.35	79.88 \pm 0.37	39.82 \pm 2.14	41.55 \pm 4.06	71.22 \pm 0.62	61.59 \pm 0.58	OOM	OOM	8.00
Expformer	83.24 \pm 1.23	71.75 \pm 1.32	79.62 \pm 0.83	39.10 \pm 1.02	44.88 \pm 2.86	72.14 \pm 0.48	OOM	OOM	OOM	8.33
SGFormer	84.50 \pm 0.82	72.54 \pm 0.23	80.31 \pm 0.58	41.85 \pm 2.41	44.93 \pm 3.72	72.68 \pm 0.16	65.92 \pm 0.19	76.17 \pm 1.86	40.13 \pm 2.13	5.11
GOAT	76.83 \pm 1.07	54.30 \pm 1.25	78.64 \pm 0.69	41.15 \pm 1.86	42.27 \pm 3.87	71.88 \pm 0.42	63.38 \pm 0.26	71.43 \pm 3.47	42.53 \pm 3.65	8.89
DE-GNN	84.67 \pm 0.53	71.69 \pm 0.43	OOM	44.42 \pm 1.93	46.51 \pm 3.64	OOM	OOM	OOM	OOM	7.44
Graphormer	82.43 \pm 0.86	70.16 \pm 0.77	OOM	41.87 \pm 1.86	43.77 \pm 4.13	OOM	OOM	OOM	OOM	9.44
ESAN	OOM	11.89								
GRIT	79.45 \pm 0.78	68.77 \pm 0.89	OOM	39.57 \pm 1.58	44.23 \pm 3.69	OOM	OOM	OOM	OOM	10.56
RD-WL	84.76 \pm 0.94	72.86 \pm 0.92	OOM	44.69 \pm 2.05	47.03 \pm 3.82	OOM	OOM	OOM	OOM	6.56
SPE	OOM	11.78								
NeuralWalker	82.10 \pm 1.23	69.52 \pm 1.06	76.70 \pm 0.83	40.71 \pm 2.33	43.30 \pm 4.59	64.59 \pm 2.64	53.33 \pm 0.94	OOM	OOM	9.89
SDE-GNN (ours)	85.92 \pm 0.42	73.83 \pm 0.55	80.63 \pm 0.51	46.32 \pm 1.53	47.24 \pm 3.56	72.90 \pm 0.26	66.04 \pm 0.35	85.23 \pm 0.32	46.80 \pm 0.25	1.44

on these larger graphs, highlighting the superior scalability of SDE-GNN over existing expressive GNNs. To further validate the efficiency advantage of SDE-GNN over existing expressive and efficient GNNs, we report the running time, memory usage, and accuracy of different methods on twitch-gamers dataset in Figure 2. We exclude NeuralWalker as its performance is not competitive on this dataset. As shown in Figure 2, SDE-GNN achieves the highest accuracy and the lowest running time, while maintaining a moderate level of memory consumption, demonstrating its overall efficiency.

5.3 ABLATION STUDIES

To verify the effectiveness of the proposed components, we compare SDE-GNN with the following two variants of 1) w/o polynomial expansion by keeping only the first order in Equation (3) and 2) w/o separation that also computes the first order term in Equation (9) using dimension reduction. As shown in Table 5, there is a significant performance drop of w/o separation, highlighting the advantage of performing computations explicitly rather than relying on dimensionality reduction. The performance of w/o separation is N/A on CLUSTER and PATTERN because the maximum number of nodes in these datasets does not exceed 200, and dimension reduction is not applied. For

Table 4: Test performance on 3 real-world graph-classification datasets.

Model	IMDB-BINARY	IMDB-MULTI	REDDIT-BINARY	Rank ↓
	Accuracy ↑	Accuracy ↑	Accuracy ↑	
GCN	70.66 ± 0.95	46.77 ± 1.41	75.49 ± 1.11	9.67
GraphSAGE	69.92 ± 1.24	46.14 ± 1.44	75.03 ± 1.23	11.33
GAT	70.87 ± 1.38	47.23 ± 1.59	74.87 ± 1.58	10.00
GIN	70.95 ± 1.15	50.32 ± 1.53	75.16 ± 1.52	6.67
GraphGPS	66.36 ± 1.43	41.77 ± 1.67	86.38 ± 1.64	11.67
Expformer	65.48 ± 1.51	41.96 ± 1.88	87.04 ± 1.75	11.33
SGFormer	64.25 ± 1.08	42.23 ± 1.45	85.73 ± 1.42	12.33
GOAT	65.36 ± 2.13	41.75 ± 2.38	83.53 ± 2.18	13.33
DE-GNN	70.93 ± 1.02	48.49 ± 1.85	91.97 ± 1.17	4.67
Graphormer	67.67 ± 1.38	45.69 ± 2.34	87.79 ± 2.32	9.67
ESAN	73.05 ± 1.54	50.78 ± 1.95	OOM	6.00
GRIT	72.13 ± 1.37	49.56 ± 1.74	89.64 ± 2.31	3.67
RD-WL	69.48 ± 2.16	49.86 ± 1.64	89.45 ± 1.63	6.33
SPE	70.44 ± 1.67	48.12 ± 1.83	OOM	10.67
NeuralWalker	72.00 ± 1.16	48.64 ± 1.77	79.66 ± 1.45	6.67
SDE-GNN (ours)	73.33 ± 1.34	50.66 ± 1.62	91.78 ± 1.48	1.67

Table 5: Ablation study results.

Model	CLUSTER	PATTERN	ogbn-arxiv	twitch-gamers
	Accuracy ↑	Accuracy ↑	Accuracy ↑	Accuracy ↑
SDE-GNN	78.68	86.82	72.90	66.04
w/o polynomial	78.12	86.74	72.81	65.91
w/o separation	N/A	N/A	69.91	64.83

Table 6: The trade-offs between test accuracy (%), runtime (s/epoch), and GPU memory usage (MB) by varying the projection dimension p on twitch-gamers dataset (with 168,114 nodes). “rank” indicates SDE-GNN’s test accuracy ranking compared to other baseline methods.

Metric \ Param p	100	200	300	400	500	600	700	800	900	1000
Runtime (s/epoch)	0.38	0.40	0.42	0.44	0.46	0.48	0.50	0.52	0.54	0.54
GPU Memory (MB)	6,204	7,094	7,874	8,764	9,788	10,320	11,344	11,708	13,732	15,460
Test Acc	65.03	65.42	65.74	65.82	65.88	65.96	66.02	65.94	66.02	66.04
Rank	2	2	2	2	2	1	1	1	1	1

w/o polynomial, its performance consistently drops on these four datasets, verifying the necessity of adopting the polynomial term in the learning process of the distance feature mapping. Besides, as people focus on the trade-off between effectiveness and cost, we resort to empirical analysis to gain insights into how dimensionality (i.e., p) affects performance. In Table 6, we vary from $p = 100$ to 1,000 on twitch-gamers dataset and report the accuracy, runtime, GPU memory usage, and the rank of our method among all compared baselines. From the results, we conclude that both performance and the required computational resources generally increase as dimensionality grows, but at different rates: p has a much smaller impact on performance compared to memory usage and runtime. For example, as p decreases from 1,000 to 100, memory and runtime decrease by 60% and 30%, respectively. While performance drops by only 1.01, and the rank drops by just one. For practical usage, a dimensionality between 100 and 300 may provide a favorable balance between performance and efficiency. More experimental results of the time and memory consumption at different node scales and the analysis of dimensionality reduction methods are deferred to Appendix E.

6 CONCLUSION

In this work, we have addressed the critical challenge of enhancing the expressiveness of DE-GNNs while ensuring scalability to large graphs. By proposing SDE-GNN, a scalable distance-enhanced GNN framework leveraging adjacency-power-based distance features, learnable polynomial distance encoding, and randomized dimensionality reduction, we significantly reduce computational complexity from quadratic to near-linear with respect to graph size. Our theoretical analysis establishes that the chosen distance features possess superior expressive power, exceeding that of commonly used alternatives. Extensive empirical evaluations across diverse benchmark datasets validate that SDE-GNN consistently outperforms strong baselines in accuracy, efficiency, and scalability. This study not only advances the understanding of distance features in GNNs but also offers a practical solution for applying expressive GNN models to large-scale graph mining tasks.

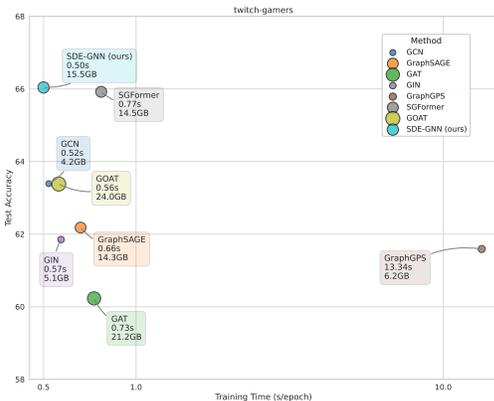


Figure 2: Comparison of training time (s/epoch) and test accuracy (%) of methods on twitch-gamers dataset (with 168,114 nodes). The size of the node reflects the GPU memory cost.

486 7 ETHICS STATEMENT

487

488 The datasets used for the experiments were publicly available and fully anonymized. We rigorously
 489 evaluated our model for potential biases, societal implications, and the safety of its generated content.
 490 The authors declare no conflicts of interest. For transparency and reproducibility, the code and data
 491 are open-sourced.

492

493 8 REPRODUCIBILITY STATEMENT

494

495 All experimental code and data are open-sourced available in an anonymous repository, allowing for
 496 the complete reproduction of our experimental results.

497

498 REFERENCES

499

500 Ralph Abboud, İsmail İlkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. The surprising
 501 power of graph neural networks with random node initialization. In *Proceedings of the Thirtieth
 502 International Joint Conference on Artificial Intelligence, IJCAI 2021, Virtual Event / Montreal,
 503 Canada, 19-27 August 2021*, pp. 2112–2118, 2021.

504 Robert Ackland et al. Mapping the us political blogosphere: Are conservative bloggers more
 505 prominent? In *BlogTalk Downunder 2005 Conference, Sydney*. BlogTalk Downunder 2005
 506 Conference, Sydney, 2005.

507

508 Beatrice Bevilacqua, Fabrizio Frasca, Derek Lim, Balasubramaniam Srinivasan, Chen Cai, Gopinath
 509 Balamurugan, Michael M. Bronstein, and Haggai Maron. Equivariant subgraph aggregation
 510 networks. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual
 511 Event, April 25-29, 2022*, 2022.

512 Beatrice Bevilacqua, Moshe Eliasof, Eli A. Meirum, Bruno Ribeiro, and Haggai Maron. Efficient
 513 subgraph gnns by learning effective selection policies. In *The Twelfth International Conference on
 514 Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024.

515 Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In
 516 *International Colloquium on Automata, Languages, and Programming*, pp. 693–703. Springer,
 517 2002.

518

519 Dexiong Chen, Till Hendrik Schulz, and Karsten M. Borgwardt. Learning long range dependencies
 520 on graphs via random walks. *CoRR*, abs/2406.03386, 2024. doi: 10.48550/ARXIV.2406.03386.
 521 URL <https://doi.org/10.48550/arXiv.2406.03386>.

522 Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count
 523 substructures? In *Advances in Neural Information Processing Systems 33: Annual Conference on
 524 Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.

525 Leonardo Cotta, Christopher Morris, and Bruno Ribeiro. Reconstruction for powerful graph rep-
 526 resentations. In *Advances in Neural Information Processing Systems 34: Annual Conference on
 527 Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
 528 1713–1726, 2021.

529

530 Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs.
 531 *CoRR*, abs/2012.09699, 2020.

532 Vijay Prakash Dwivedi, Chaitanya K. Joshi, Anh Tuan Luu, Thomas Laurent, Yoshua Bengio, and
 533 Xavier Bresson. Benchmarking graph neural networks. *J. Mach. Learn. Res.*, 24:43:1–43:48, 2023.
 534 URL <https://jmlr.org/papers/v24/22-0567.html>.

535

536 Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychome-
 537 trika*, 1(3):211–218, 1936.

538 Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding structure with randomness:
 539 Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):
 217–288, 2011. doi: 10.1137/090771806. URL <https://doi.org/10.1137/090771806>.

- 540 William L. Hamilton, Zitao Ying, and Jure Leskovec. Inductive representation learning on large
541 graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural*
542 *Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 1024–1034,
543 2017.
- 544 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,
545 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *Advances in*
546 *neural information processing systems*, 33:22118–22133, 2020.
- 547 Yinan Huang, William Lu, Joshua Robinson, Yu Yang, Muhan Zhang, Stefanie Jegelka, and Pan Li.
548 On the stability of expressive positional encodings for graphs. In *The Twelfth International Confer-*
549 *ence on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net,
550 2024. URL <https://openreview.net/forum?id=xAqcJ9XoTf>.
- 551 William B Johnson, Joram Lindenstrauss, et al. Extensions of lipschitz mappings into a hilbert space.
552 *Contemporary mathematics*, 26(189-206):1, 1984.
- 553 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.
554 In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April*
555 *24-26, 2017, Conference Track Proceedings, 2017*.
- 556 Kezhi Kong, Jiuhai Chen, John Kirchenbauer, Renkun Ni, C. Bayan Bruss, and Tom Goldstein. Goat:
557 A global transformer on large-scale graphs. In *International Conference on Machine Learning*,
558 2023. URL <https://api.semanticscholar.org/CorpusID:260927574>.
- 559 Devin Kreuzer, Dominique Beaini, William L. Hamilton, Vincent L etourneau, and Pruden-
560 cio Tossou. Rethinking graph transformers with spectral attention. In Marc’Aurelio Ran-
561 zato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan
562 (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neu-*
563 *ral Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp.
564 21618–21629, 2021. URL [https://proceedings.neurips.cc/paper/2021/hash/](https://proceedings.neurips.cc/paper/2021/hash/b4fd1d2cb085390fbbadae65e07876a7-Abstract.html)
565 [b4fd1d2cb085390fbbadae65e07876a7-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/b4fd1d2cb085390fbbadae65e07876a7-Abstract.html).
- 566 Pan Li, Yanbang Wang, Hongwei Wang, and Jure Leskovec. Distance encoding: Design provably
567 more powerful neural networks for graph representation learning. In *Advances in Neural Informa-*
568 *tion Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020,*
569 *NeurIPS 2020, December 6-12, 2020, virtual, 2020*.
- 570 Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and
571 Ser-Nam Lim. Large scale learning on non-homophilous graphs: New benchmarks and strong
572 simple methods. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang,
573 and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34:*
574 *Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December*
575 *6-14, 2021, virtual*, pp. 20887–20902, 2021. URL [https://proceedings.neurips.cc/](https://proceedings.neurips.cc/paper/2021/hash/ae816a80e4c1c56caa2eb4e1819cbb2f-Abstract.html)
576 [paper/2021/hash/ae816a80e4c1c56caa2eb4e1819cbb2f-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/ae816a80e4c1c56caa2eb4e1819cbb2f-Abstract.html).
- 577 Derek Lim, Joshua David Robinson, Lingxiao Zhao, Tess E. Smidt, Suvrit Sra, Haggai Maron, and
578 Stefanie Jegelka. Sign and basis invariant networks for spectral graph representation learning. In
579 *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda,*
580 *May 1-5, 2023, 2023*.
- 581 Yuankai Luo, Lei Shi, and Xiao-Ming Wu. Classic gnns are strong baselines: Reassess-
582 ing gnns for node classification. In Amir Globersons, Lester Mackey, Danielle Belgrave,
583 Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances*
584 *in Neural Information Processing Systems 38: Annual Conference on Neural Informa-*
585 *tion Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 -*
586 *15, 2024, 2024*. URL [http://papers.nips.cc/paper_files/paper/2024/](http://papers.nips.cc/paper_files/paper/2024/hash/b10ed15ff1aa864f1be3a75f1ffc021b-Abstract-Datasets_and_Benchmarks_Track.html)
587 [hash/b10ed15ff1aa864f1be3a75f1ffc021b-Abstract-Datasets_and_](http://papers.nips.cc/paper_files/paper/2024/hash/b10ed15ff1aa864f1be3a75f1ffc021b-Abstract-Datasets_and_Benchmarks_Track.html)
588 [Benchmarks_Track.html](http://papers.nips.cc/paper_files/paper/2024/hash/b10ed15ff1aa864f1be3a75f1ffc021b-Abstract-Datasets_and_Benchmarks_Track.html).
- 589 Liheng Ma, Chen Lin, Derek Lim, Adriana Romero-Soriano, Puneet K Dokania, Mark Coates, Philip
590 Torr, and Ser-Nam Lim. Graph inductive biases in transformers without message passing. In
591 *International Conference on Machine Learning*, pp. 23321–23337. PMLR, 2023.

- 594 Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph
595 networks. In *Advances in Neural Information Processing Systems 32: Annual Conference on*
596 *Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver,*
597 *BC, Canada, 2019a.*
- 598 Haggai Maron, Heli Ben-Hamu, Nadav Shamir, and Yaron Lipman. Invariant and equivariant graph
599 networks. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans,*
600 *LA, USA, May 6-9, 2019, 2019b.*
- 602 Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav
603 Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks.
604 In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Inno-*
605 *vative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium*
606 *on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January*
607 *27 - February 1, 2019*, pp. 4602–4609, 2019.
- 608 Christopher Morris, Nils M. Kriege, Franka Bause, Kristian Kersting, Petra Mutzel, and Marion
609 Neumann. TUDataset: A collection of benchmark datasets for learning with graphs. *CoRR*,
610 abs/2007.08663, 2020a. URL <https://arxiv.org/abs/2007.08663>.
- 611 Christopher Morris, Gaurav Rattan, and Petra Mutzel. Weisfeiler and leman go sparse: Towards
612 scalable higher-order graph embeddings. In *Advances in Neural Information Processing Systems*
613 *33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December*
614 *6-12, 2020, virtual, 2020b.*
- 615 Dilcia Pérez and Yamilet Quintana. A survey on the weierstrass approximation theorem. *Divulga-*
616 *ciones Matemáticas*, 16(1):231–247, 2008.
- 618 Ladislav Rampásek, Michael Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and
619 Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In Sanmi
620 Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), *Advances*
621 *in Neural Information Processing Systems 35: Annual Conference on Neural Information*
622 *Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December*
623 *9, 2022, 2022.* URL http://papers.nips.cc/paper_files/paper/2022/hash/5d4834a159f1547b267a05a4e2b7cf5e-Abstract-Conference.html.
- 625 Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-scale attributed node embedding. *J.*
626 *Complex Networks*, 9(2), 2021. doi: 10.1093/COMNET/CNAB014. URL [https://doi.org/](https://doi.org/10.1093/comnet/cnab014)
627 [10.1093/comnet/cnab014](https://doi.org/10.1093/comnet/cnab014).
- 628 Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Random features strengthen graph neural
629 networks. In *Proceedings of the 2021 SIAM International Conference on Data Mining, SDM 2021,*
630 *Virtual Event, April 29 - May 1, 2021*, pp. 333–341, 2021.
- 632 Hamed Shirzad, Ameya Velingker, B. Venkatachalam, Danica J. Sutherland, and Ali Kemal Sinop.
633 Exphormer: Sparse transformers for graphs. In *International Conference on Machine Learning,*
634 *2023.* URL <https://api.semanticscholar.org/CorpusID:257482539>.
- 635 Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
636 Bengio. Graph attention networks. In *6th International Conference on Learning Representations,*
637 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings, 2018.*
- 638 Ameya Velingker, Ali Kemal Sinop, Ira Ktena, Petar Velickovic, and Sreenivas Golla-
639 pudi. Affinity-aware graph networks. In Alice Oh, Tristan Naumann, Amir Globerson,
640 Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in Neural In-*
641 *formation Processing Systems 36: Annual Conference on Neural Information Pro-*
642 *cessing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*
643 *2023, 2023.* URL [http://papers.nips.cc/paper_files/paper/2023/hash/](http://papers.nips.cc/paper_files/paper/2023/hash/d642b0633afad94f660554e05b40608e-Abstract-Conference.html)
644 [d642b0633afad94f660554e05b40608e-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d642b0633afad94f660554e05b40608e-Abstract-Conference.html).
- 645 Haorui Wang, Haoteng Yin, Muhan Zhang, and Pan Li. Equivariant and stable positional encoding
646 for more powerful graph neural networks. In *The Tenth International Conference on Learning*
647 *Representations, ICLR 2022, Virtual Event, April 25-29, 2022, 2022.*

648 Qitian Wu, Wentao Zhao, Chenxiao Yang, Hengrui Zhang, Fan Nie, Haitian Jiang, Yatao Bian, and
649 Junchi Yan. Sgformer: Simplifying and empowering transformers for large-graph representations.
650 *Advances in Neural Information Processing Systems*, 36:64753–64773, 2023.
651

652 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
653 networks? In *7th International Conference on Learning Representations, ICLR 2019, New Orleans,
654 LA, USA, May 6-9, 2019*, 2019.

655 Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning with
656 graph embeddings. In *International conference on machine learning*, pp. 40–48. PMLR, 2016.
657

658 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen,
659 and Tie-Yan Liu. Do transformers really perform badly for graph representation? In
660 Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman
661 Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference
662 on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,
663 pp. 28877–28888, 2021a. URL [https://proceedings.neurips.cc/paper/2021/
664 hash/f1c1592588411002af340cbaedd6fc33-Abstract.html](https://proceedings.neurips.cc/paper/2021/hash/f1c1592588411002af340cbaedd6fc33-Abstract.html).

665 Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-
666 Yan Liu. Do transformers really perform badly for graph representation? In *Neural Information
667 Processing Systems*, 2021b. URL [https://api.semanticscholar.org/CorpusID:
668 265104899](https://api.semanticscholar.org/CorpusID:265104899).

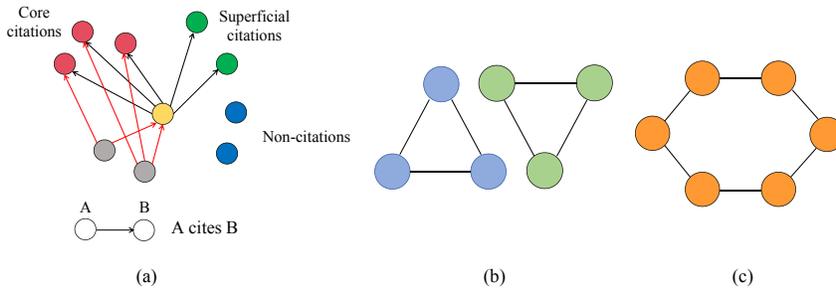
669 Bohang Zhang, Shengjie Luo, Liwei Wang, and Di He. Rethinking the expressive power of gnns
670 via graph biconnectivity. In *The Eleventh International Conference on Learning Representations,
671 ICLR 2023, Kigali, Rwanda, May 1-5, 2023*, 2023.

672 Bohang Zhang, Lingxiao Zhao, and Haggai Maron. On the expressive power of spectral invariant
673 graph neural networks. In *Forty-first International Conference on Machine Learning, ICML 2024,
674 Vienna, Austria, July 21-27, 2024*, 2024.
675

676 Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In *Advances in
677 Neural Information Processing Systems 31: Annual Conference on Neural Information Processing
678 Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 5171–5181, 2018.
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

702 A EXAMPLES TO SHOW THE EXPRESSIVE ADVANTAGE OF DE-GNNs

703
704 The figure below demonstrates the expressive advantage of DE-GNNs through the task of triangle
705 detection. In Figure 3 (b) and (c), vanilla GNNs fail to determine whether a node is part of a triangle
706 because all nodes share the same local structure—each connected to two neighbors. As a result,
707 standard message passing assigns identical representations to all 12 nodes, making it impossible to
708 distinguish the blue and green nodes (which are part of triangles) from the orange ones (which are
709 not). In contrast, the proposed distance-enhanced method determines whether a node is part of a
710 triangle by checking if there exists another node that is both one-hop and two-hop away—an indicator
711 of the triangle and is a relation captured by the distance features. Figure 3 (a) further illustrates the
712 practical utility of triangle detection in a real-world citation network, where nodes represent papers.
713 In this scenario, not all citations carry equal importance—some are core citations, while others are
714 superficial. For example, the red nodes are core citations of the yellow node, as all papers citing the
715 yellow one also cite the red ones. In contrast, the green nodes are superficial citations lacking this
716 pattern. Triangle-aware expressiveness provides an effective means to capture such relationships, as
717 it allows one to infer whether two cited nodes form a triangle and thus identify core versus superficial
718 citations more effectively.



719
720
721
722
723
724
725
726
727
728
729 Figure 3: Examples to show the expressiveness of DE-GNNs. (a) depicts a real-world citation
730 networks with core citations, superficial citations, and non-citations. (b) and (c) illustrate two regular
731 graphs, where the objective is to detect nodes within triangles.

732
733
734 B INTRODUCTION OF GD-WL

735
736 **Notations.** We use $G(V, E)$ to denote a graph. We use \mathcal{C} to denote a set of colors and $\text{hash}(\cdot)$ to note
737 a bijective hash function that maps an arbitrary object x into a color $\text{hash}(x) \in \mathcal{C}$. We use the notation
738 $\{\{\}\}$ to denote a multiset. We use d to denote a distance metric, where given a graph $G(V, E)$, it
739 outputs a function $d_G : V \times V \rightarrow \mathbb{R}^{D_G}$ that maps a pair of nodes within V into a distance vector.
740 Here, the dimension of the distance vector D_G may vary across graphs, and we will introduce it
741 later in this section. A partition of a set S is a collection Q of nonempty, pairwise disjoint subsets
742 of S such that $\bigcup_{S' \in Q} S' = S$. We use $Q_1 \preceq Q_2$ to denote that partition Q_1 is finer than Q_2 , where
743 for every $S_1 \in Q_1$, there exists $S_2 \in Q_2$ such that $S_1 \subseteq S_2$. We denote $Q_1 \prec Q_2$ if $Q_1 \preceq Q_2$ and
744 $Q_1 \neq Q_2$.

745 **GD-WL.** GD-WL is an extension of the WL-test for the graph isomorphism test, which incorporates
746 a distance metric in the color refinement procedure. Specifically, given a graph $G(V, E)$, GD-WL
747 begins by assigning a predefined $c \in \mathcal{C}$ to all nodes and iteratively applies the following formula to
748 update the color of nodes

$$749 \chi_G^{t+1}(u) = \text{hash}(\{\{\chi_G^t(v), d_G(u, v)\} : v \in V\}), \tag{13}$$

750
751 where $\chi_G^t(u)$ denotes the color of node u at round t . The set of node colors obtained at iteration t ,
752 denoted as $\{\{\chi_G^t(u) : u \in V\}\}$, can induce a partition of the node set as $Q^t(V) = \{(\chi_G^t)^{-1}(c) : c \in \mathcal{C}\}$,
753 where $(\chi_G^t)^{-1}(c) = \{u \in V : \chi^t(u) = c\}$ represents the set of nodes assigned the same color
754 c . The GD-WL will terminate at round $t + 1$, once the partition is stable, i.e., $Q^{t+1}(V) = Q^t(V)$.
755 One important property of the partition is that the obtained partition $\{Q^t(V)\}_{t \in \mathbb{N}}$ is a sequence
of refinements, where there exists a $T \in \mathbb{N}$ such that $Q^t(V) \prec Q^{t+1}(V)$ for all $t \leq T$ and

756 $Q^t(V) = Q^{t+1}(V)$ for all $t \geq T^1$. Since the finest partition is $\{\{u\} : u \in V\}$, the T will be less than
 757 $|V|$ and the GD-WL will terminate at finite step. To determine whether two graphs $G_1(V_1, E_1)$ and
 758 $G_2(V_2, E_2)$ are isomorphic, we need to run GD-WL on two graphs in parallel and compare their final
 759 node colors. In this situation, we use notation $Q^t(V_1 \cup V_2)$ to denote the partition of $V_1 \cup V_2$ at round
 760 t , where each element of $Q^t(V_1 \cup V_2)$ is a subset of $V_1 \cup V_2$ consisting of nodes that share the same
 761 color. Algorithm 2 outlines the detailed procedure by which GD-WL determines whether two graphs
 762 are isomorphic.

763 **Distance Feature Mapping.** Various choices for d are possible. For instance, it can be the shortest-
 764 path distance, where $D_G \equiv 1$ and $d_G(u, v) \in \mathbb{R}^1$ is the shortest path from u to v ; or it can be a
 765 fixed-length random walk distance, where $D_G \equiv K$ and $d_G(u, v) = [p_{u,v}^1, \dots, p_{u,v}^k] \in \mathbb{R}^k$, with
 766 $p_{u,v}^k$ denoting the probability of reaching v from u via a k -step random walk. (Zhang et al., 2024)
 767 demonstrates that GD-WL equipped with the eigenspace projection distance exhibits the highest
 768 expressiveness among commonly used distance feature mappings. Therefore, in our analysis, we
 769 focus on the eigenspace projection distance and our adjacency-power-based distance, denoting them
 770 by d^{EP} and d^{AP} , respectively. Specifically, given a graph $G(V, E)$, let $\tilde{L} = D^{-1/2}LD^{-1/2}$ denote
 771 the normalized Laplacian matrix and $\tilde{L} = \sum_{i=1}^m \lambda_i P_i$ denotes its spectral decomposition, where
 772 $\lambda_1, \dots, \lambda_m$ are the eigenvalues of \tilde{L} and $P_1, \dots, P_m \in \mathbb{R}^{|V| \times |V|}$ are the corresponding eigenspace
 773 projections. The eigenspace projection distance between two nodes $u, v \in V$ is

$$774 d_G^{EP}(u, v) = [\lambda_1, (P_1)_{u,v}, \lambda_2, (P_2)_{u,v}, \dots, \lambda_m, (P_m)_{u,v}] \in \mathbb{R}^{2m}, \quad (14)$$

775 where $(P_i)_{u,v}$ denote the u -th row and v -th column of P_i . Let $\tilde{A} = \tilde{D}^{-1/2}A\tilde{D}^{-1/2}$ be the
 776 normalized adjacency matrix of the given graph G and $n = |V|$ be the number of nodes. The
 777 adjacency-power-based distance feature between two nodes u, v is defined as

$$778 d_G^{AP}(u, v) = [\tilde{A}_{u,v}^0, \tilde{A}_{u,v}^1, \dots, \tilde{A}_{u,v}^n] \in \mathbb{R}^{n+1}, \quad (15)$$

779 where \tilde{A}^i denotes the i -th power of \tilde{A} .

783 Algorithm 2 The Generalized Distance Weisfeiler-Lehman Algorithm

784 **Require:** Graphs $G_1 = (V_1, E_1)$, graph $G_2 = (V_2, E_2)$, distance metric d

785 **Ensure:** Whether G_1 and G_2 are isomorphism

786 1: Initialize $t = 0$, $\chi_{G_1}^0(u) = c$ and $\chi_{G_2}^0(x) = c$ for all $u \in V_1$ and $x \in V_2$ correspondingly

787 2: **while** *True* **do**

788 3: $t \leftarrow t + 1$

789 4: **for** $u \in V_1$ **do**

790 5: $\chi_{G_1}^t(u) = \text{hash}(\{(d_{G_1}(u, v), \chi_{G_1}^{t-1}(v)) : v \in V_1\})$

791 6: **end for**

792 7: **for** $x \in V_2$ **do**

793 8: $\chi_{G_2}^t(x) = \text{hash}(\{(d_{G_2}(x, y), \chi_{G_2}^{t-1}(y)) : y \in V_2\})$

794 9: **end for**

795 10: **if** $Q^t(V_1 \cup V_2) = Q^{t-1}(V_1 \cup V_2)$ **then**

796 11: **break**

797 12: **end if**

798 13: **end while**

799 14: **Return:** $\{\{\chi_{G_1}^t(u) : u \in V_1\}\} = \{\{\chi_{G_2}^t(x) : x \in V_2\}\}$

802 C PROOFS OF THEOREMS

803
 804 In this section, we first give lemmas that are required for the proofs of Theorem 1 and Theorem 2,
 805 and then give the detailed proofs.

806
 807 **Lemma 1** For any $\lambda_1, \dots, \lambda_n \in \mathbb{C}$, any $\sigma_1, \dots, \sigma_n \in \mathbb{C}$, if $\sum_{i=1}^n \lambda_i^k = \sum_{i=1}^n \sigma_i^k$ holds for $0 \leq k \leq n$,
 808 then we will have $\{\{\lambda_1, \dots, \lambda_n\}\} = \{\{\sigma_1, \dots, \sigma_n\}\}$.

809 ¹More details can be found in Appendix B.1 of (Zhang et al., 2023).

Proof. Let $e_i(a_1, \dots, a_n)$ denote $\sum_{1 \leq j_1 < j_2 < \dots < j_i \leq n} (-1)^i a_{j_1} a_{j_2} \dots a_{j_n}$, and $p_i(a_1, \dots, a_n)$ denote $\sum_{j=1}^n a_j^i$. According to Newton's Identities, for $1 \leq i \leq n$, we have

$$e_i(a_1, \dots, a_n) = \frac{-1}{i} \sum_{j=0}^{i-1} e_j(a_1, \dots, a_n) * p_{i-j}(a_1, \dots, a_n), \quad (16)$$

where $e_0(a_1, \dots, a_n) = 0$. Thus if $p_i(\lambda_1, \dots, \lambda_n) = p_i(\sigma_1, \sigma_2, \dots, \sigma_n)$ for $0 \leq i \leq n$, then by recursion, we will have $e_i(\lambda_1, \dots, \lambda_n) = e_i(\sigma_1, \dots, \sigma_n)$ for $0 \leq i \leq n$. Then we can construct two polynomials $f(x) = \prod_{i=1}^n (x - \lambda_i)$, $g(x) = \prod_{i=1}^n (x - \sigma_i)$, where the coefficients of x^{n-i} in $f(x)$ and $g(x)$ will be $e_i(\lambda_1, \dots, \lambda_n)$ and $e_i(\sigma_1, \dots, \sigma_n)$ respectively. Since $e_i(\lambda_1, \dots, \lambda_n) = e_i(\sigma_1, \dots, \sigma_n)$ for all $0 \leq i \leq n$, we will have $f(x) = g(x)$, so $\{\{\lambda_1, \dots, \lambda_n\}\} = \{\{\sigma_1, \dots, \sigma_n\}\}$.

Lemma 2 Let $\sigma(M)$ denote the multiset of eigenvalues for a matrix $M \in \mathbb{C}^{n \times n}$, and $\text{tr}(M)$ denote the trace of M . For any two matrices $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}$ with $\text{tr}(\mathbf{X}^i) = \text{tr}(\mathbf{Y}^i)$ for $0 \leq i \leq n$, we will have $\sigma(\mathbf{X}) = \sigma(\mathbf{Y})$.

Proof. Let x_1, \dots, x_n denote the eigenvalues of \mathbf{X} and y_1, \dots, y_n denote the eigenvalues of \mathbf{Y} . Since $\text{tr}(\mathbf{X}^i) = \text{tr}(\mathbf{Y}^i)$ for $0 \leq i \leq n$, we will have $\sum_{j=1}^n x_j^i = \sum_{j=1}^n y_j^i$ for $1 \leq i \leq n$. According to Lemma 1, we will have $\{\{x_1, \dots, x_n\}\} = \{\{y_1, \dots, y_n\}\}$.

C.1 PROOF OF THEOREM 1

We refer Algorithm 2 as the specific execution procedure of GD-WL. Assume that Algorithm 2 runs for T rounds. Since it cannot distinguish between G_1 and G_2 , we will have $\{\{\chi_{G_1}^T(u) : u \in V_1\}\} = \{\{\chi_{G_2}^T(x) : x \in V_2\}\}$. Thus, the number of nodes must be the same, i.e., $|V_1| = |V_2|$. Furthermore, there exists a bijective mapping: $f : V_1 \rightarrow V_2$ that satisfies $\chi_{G_1}^T(u) = \chi_{G_2}^T(f(u))$. Considering any two nodes that $x = f(u)$, since $\chi_{G_1}^T(u) = \chi_{G_2}^T(x)$, we will have

$$\begin{aligned} \text{hash}(\{\{\chi_{G_1}^{T-1}(v), d_{G_1}^{AP}(u, v) : v \in V_1\}\}) &= \text{hash}(\{\{\chi_{G_2}^{T-1}(y), d_{G_2}^{AP}(x, y) : y \in V_2\}\}) \\ \xrightarrow{I_1} \{\{\chi_{G_1}^{T-1}(v), d_{G_1}^{AP}(u, v) : v \in V_1\}\} &= \{\{\chi_{G_2}^{T-1}(y), d_{G_2}^{AP}(x, y) : y \in V_2\}\}. \end{aligned} \quad (17)$$

I_1 comes from the injectivity of the hash function. Therefore, there exists a bijective mapping $f_1 : V_1 \rightarrow V_2$ satisfy that $d_{G_1}^{AP}(u, v) = d_{G_2}^{AP}(x, f_1(v))$. Among all nodes $v \in V_1$, the adjacency-power-based distance $d_{G_1}^{AP}(u, u)$ is the only one whose first entry equals one, i.e., $d_{G_1}^{AP}(u, u)[1] = 1$, since only the diagonal elements of $\tilde{\mathbf{A}}_{G_1}^0 = \mathbf{I}$ are nonzero. The same reasoning applies to $d_{G_2}^{AP}(x, y)$ for $y \in V_2$. Consequently, we have $d_{G_1}^{AP}(u, u) = d_{G_2}^{AP}(x, x)$. Based on the above reasoning, we have

$$\chi_{G_1}^T(u) = \chi_{G_2}^T(x) \implies d_{G_1}^{AP}(u, u) = d_{G_2}^{AP}(x, x) \quad (18)$$

Then we will have

$$\begin{aligned} \{\{\chi_{G_1}^T(u) : u \in V_1\}\} &= \{\{\chi_{G_2}^T(v) : v \in V_2\}\} \\ \xrightarrow{I_1} \sum_{u \in V_1} d_{G_1}^{AP}(u, u) &= \sum_{x \in V_2} d_{G_2}^{AP}(x, x) \\ \xrightarrow{I_2} \text{tr}(\tilde{\mathbf{A}}_{G_1}^i) &= \text{tr}(\tilde{\mathbf{A}}_{G_2}^i), 0 \leq i \leq n \\ \xrightarrow{I_3} \sigma(\tilde{\mathbf{A}}_{G_1}) &= \sigma(\tilde{\mathbf{A}}_{G_2}) \end{aligned} \quad (19)$$

where $n = |V_1| = |V_2|$ is the number of nodes, I_1 comes by applying equation 18, I_2 comes from the observation that for a graph $G(V, E)$ with $|V| = m$, we have

$$\sum_{u \in V} d_G^{AP}(u, u)[i] = \sum_{j=1}^m \tilde{\mathbf{A}}_G^i[j, j] = \text{tr}(\tilde{\mathbf{A}}_G^i), 1 \leq i \leq m, \quad (20)$$

where $d_G^{AP}(u, u)[i]$ denotes the i -th element of $d_G^{AP}(u, u)$ and $\tilde{\mathbf{A}}_G^i[j, j]$ denotes element in the j -th row and j -th column of $\tilde{\mathbf{A}}_G$. I_3 comes by applying Lemma 2. This completes the proof of the first statement of Theorem 1.

864 Consider the second statement. According to the first statement, we know that $\tilde{\mathbf{A}}_{G_1}$ and $\tilde{\mathbf{A}}_{G_2}$ have
 865 the same eigenvalues. Let m be the number of distinct eigenvalues, and denote them by $\lambda_1, \dots, \lambda_m$.
 866 Since for any graph G it holds that $\tilde{\mathbf{L}}_G = \mathbf{I} - \tilde{\mathbf{A}}_G$, the values $1 - \lambda_1, \dots, 1 - \lambda_m$ are also the distinct
 867 eigenvalues of both $\tilde{\mathbf{L}}_{G_1}$ and $\tilde{\mathbf{L}}_{G_2}$. We denote the corresponding eigenspace projections of $\tilde{\mathbf{L}}_{G_1}$ by
 868 $\mathbf{P}_{G_1,1}, \dots, \mathbf{P}_{G_1,m}$, and those of $\tilde{\mathbf{L}}_{G_2}$ by $\mathbf{P}_{G_2,1}, \dots, \mathbf{P}_{G_2,m}$. Then, since $d_{G_1}^{AP}(u, v) = d_{G_2}^{AP}(x, y)$,
 869 we have
 870

$$\begin{aligned}
 871 \quad & d_{G_1}^{AP}(u, v) = d_{G_2}^{AP}(x, y) \implies \tilde{\mathbf{A}}_{G_1}^i[u, v] = \tilde{\mathbf{A}}_{G_2}^i[u, v], \quad 0 \leq i \leq n \\
 872 \quad & \implies \tilde{\mathbf{L}}_{G_1}^i[u, v] = \tilde{\mathbf{L}}_{G_2}^i[u, v], \quad 0 \leq i \leq n \\
 873 \quad & \implies \sum_{j=1}^m \lambda_j^i \mathbf{P}_{G_1,j}[u, v] = \sum_{j=1}^m \lambda_j^i \mathbf{P}_{G_2,j}[u, v], \quad 0 \leq i \leq n \\
 874 \quad & \implies \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \\ \lambda_1 & \lambda_2 & \dots & \lambda_m \\ \lambda_1^2 & \lambda_2^2 & \dots & \lambda_m^2 \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_1^n & \lambda_2^n & \dots & \lambda_m^n \end{bmatrix}}_{\mathbf{V}} \cdot \begin{bmatrix} \mathbf{P}_{G_1,1}[u, v] - \mathbf{P}_{G_2,1}[u, v] \\ \mathbf{P}_{G_1,2}[u, v] - \mathbf{P}_{G_2,2}[u, v] \\ \vdots \\ \mathbf{P}_{G_1,m}[u, v] - \mathbf{P}_{G_2,m}[u, v] \end{bmatrix} = \mathbf{0} \quad (21) \\
 877 \quad & \implies \begin{bmatrix} \mathbf{P}_{G_1,1}[u, v] \\ \mathbf{P}_{G_1,2}[u, v] \\ \vdots \\ \mathbf{P}_{G_1,m}[u, v] \end{bmatrix} = \begin{bmatrix} \mathbf{P}_{G_2,1}[u, v] \\ \mathbf{P}_{G_2,2}[u, v] \\ \vdots \\ \mathbf{P}_{G_2,m}[u, v] \end{bmatrix} \\
 881 \quad & \implies d_{G_1}^{EP}(u, v) = d_{G_2}^{EP}(x, y),
 \end{aligned}$$

882 where I_1 follows from the properties of spectral decomposition, and I_2 holds because the columns
 883 of \mathbf{V} are linearly independent (the submatrix formed by the first m rows of \mathbf{V} is the transpose of a
 884 Vandermonde matrix) and I_3 comes from the definition of d^{EP} .
 885

886 C.2 PROOF OF THEOREM 2

887 For notational simplicity, let $\chi_G^t(u)$ and $\bar{\chi}_G^t(u)$ denote the color of node u in graph G at round t of
 888 Algorithm 2 using d^{AP} and d^{EP} , respectively. Given two graphs $G_1(V_1, E_1)$ and $G_2(V_2, E_2)$, let t_1
 889 denote the final round of Algorithm 2 with d^{AP} . We first show that for any two nodes $u \in V_1$ and
 890 $x \in V_2$, if $\chi_{G_1}^{t_1}(u) = \chi_{G_2}^{t_1}(x)$, we will have $\bar{\chi}_{G_1}^t(u) = \bar{\chi}_{G_2}^t(x)$ for all $t \in \mathbb{N}$, which is

$$891 \quad \chi_{G_1}^{t_1}(u) = \chi_{G_2}^{t_1}(x) \implies \bar{\chi}_{G_1}^t(u) = \bar{\chi}_{G_2}^t(x), t \in \mathbb{N} \quad (22)$$

892 We prove this by mathematical induction. At $t = 0$, all nodes have the same color, and Equation (22)
 893 holds. Assume Equation (22) holds at round t . Since t_1 is the final round, the node partition is stable
 894 after round t_1 , and we will have

$$\begin{aligned}
 895 \quad & \chi_{G_1}^{t_1}(u) = \chi_{G_2}^{t_1}(x) \implies \chi_{G_1}^{t_1+1}(u) = \chi_{G_2}^{t_1+1}(x) \\
 896 \quad & \xrightarrow{I_1} \{(\chi_{G_1}^{t_1}(v), d_{G_1}^{AP}(u, v)) : v \in V_1\} = \{(\chi_{G_2}^{t_1}(y), d_{G_2}^{AP}(x, y)) : y \in V_2\} \\
 897 \quad & \xrightarrow{I_2} \{(\bar{\chi}_{G_1}^t(v), d_{G_1}^{EP}(u, v)) : v \in V_1\} = \{(\bar{\chi}_{G_2}^t(y), d_{G_2}^{EP}(x, y)) : y \in V_2\} \\
 898 \quad & \implies \bar{\chi}_{G_1}^{t+1}(u) = \bar{\chi}_{G_2}^{t+1}(x),
 \end{aligned} \quad (23)$$

899 where I_1 comes from the definition of $\chi_{G_1}^{t_1+1}(u)$ and $\chi_{G_2}^{t_1+1}(x)$, I_2 comes from the condition that
 900 Equation (22) holds at round t and the second statement of Theorem 1. This finishes the proof of
 901 Equation (22). Let t_2 denote the final round of Algorithm 2 with d^{EP} . If Algorithm 2 with d^{AP} can
 902 not distinguish G_1 and G_2 , by applying Equation (22), we will have

$$\begin{aligned}
 903 \quad & \{\{\chi_{G_1}^{t_1}(u) : u \in V_1\}\} = \{\{\chi_{G_2}^{t_1}(x) : x \in V_2\}\} \\
 904 \quad & \implies \{\{\bar{\chi}_{G_1}^{t_2}(u) : u \in V_1\}\} = \{\{\bar{\chi}_{G_2}^{t_2}(x) : x \in V_2\}\},
 \end{aligned} \quad (24)$$

905 which finishes the proof of Theorem 2.
 906

D EXPERIMENTAL DETAILS

D.1 DATASETS

Table 7: Overview of the graph datasets used in this study.

Dataset	# Graphs	Avg. # nodes	Avg. # edges	# Classes	Prediction task	Metric
PATTERN	14,000	118.9	3,039.3	2	Node Classification	Accuracy
CLUSTER	12,000	117.2	2,150.9	6	Node Classification	Accuracy
USA-Airports	1	1,190	13,599	4	Node Classification	Accuracy
Europe-Airports	1	399	5,995	4	Node Classification	Accuracy
Brazil-Airports	1	131	1,074	4	Node Classification	Accuracy
Cora	1	2,708	10,556	7	Node Classification	Accuracy
CiteSeer	1	3,327	9,104	6	Node Classification	Accuracy
PubMed	1	19,717	88,648	3	Node Classification	Accuracy
Squirrel	1	2,223	46,998	5	Node Classification	Accuracy
Chameleon	1	2,277	36,101	5	Node Classification	Accuracy
ogbn-arxiv	1	169,343	1,166,243	40	Node Classification	Accuracy
twitch-gamers	1	168,114	6,797,557	2	Node Classification	Accuracy
pokec	1	1,632,803	30,622,564	2	Node Classification	Accuracy
snap-patents	1	2,923,922	13,975,788	5	Node Classification	Accuracy
IMDB-BINARY	1,000	19.8	193.1	2	Graph Classification	Accuracy
IMDB-MULTI	1,500	13.0	131.8	3	Graph Classification	Accuracy
REDDIT-BINARY	2,000	429.6	995.5	2	Graph Classification	Accuracy

CLUSTER and PATTERN (Dwivedi et al., 2023) are synthetic datasets sampled from Stochastic Block Model. Unlike other datasets, the prediction task here is an inductive node-level classification. In PATTERN the task is to recognize which nodes in a graph belong to one of 100 possible sub-graph patterns that were randomly generated with different SBM parameters than the rest of the graph. In CLUSTER, every graph is composed of 6 SBM-generated clusters, each drawn from the same distribution, with only a single node per cluster containing a unique cluster ID. The task is to infer which cluster ID each node belongs to.

USA-Airports, Europe-Airports, and Brazil-Airports (Ackland et al., 2005) are three air traffic networks, which were collected from the government websites throughout the year 2016 and were used to evaluate algorithms to learn structural representations of nodes. Networks are built such that nodes represent airports and there exists an edge between two nodes if there are commercial flights between them. In each dataset, the airports are divided into 4 different levels according to the annual passengers flow distribution by 3 quantiles: 25%, 50%, 75%. The goal is to infer the level of an airport using solely the connectivity pattern of them.

Cora, CiteSeer, and PubMed (Yang et al., 2016) are three text classification datasets. Each dataset contains bag-of-words representation of documents and citation links between the documents. The bag-of-words are embedded as feature vectors \mathbf{x} , and the graph is constructed based on the citation links. The goal is to classify each document into one class.

Squirrel and Chameleon (Rozemberczki et al., 2021) are two heterophilous graph datasets. They are page-page networks on specific topics in Wikipedia. In those datasets, nodes represent web pages and edges are mutual links between pages. And node features correspond to several informative nouns in the Wikipedia pages. We classify the nodes into five categories in term of the number of the average monthly traffic of the web page. In this study, we adopt their filtered version described in (Luo et al., 2024).

ogbn-arxiv (Hu et al., 2020) is a directed graph, representing the citation network between all Computer Science (CS) arXiv papers indexed by MAG. Each node is an arXiv paper and each directed edge indicates that one paper cites another one. Each paper comes with a 128-dimensional feature vector obtained by averaging the embeddings of words in its title and abstract. The embeddings of individual words are computed by running the skip-gram model over the MAG corpus.

twitch-gamers (Lim et al., 2021) is a connected undirected graph of relationships between accounts on the streaming platform Twitch. Each node is a Twitch account, and edges exist between accounts that are mutual followers. The binary lassification task is to predict whether the channel has explicit content.

pokec (Lim et al., 2021) is an online social network, its task is to predict reported gender, certain account labels, or use of explicit content on user accounts, which is the same to twitch-gamers.

snap-patents (Lim et al., 2021) is a large-scale patent citation network dataset. Its primary goal is to predict the year a patent was granted, treating it as a node-level prediction task where each patent is a node and each citation is a directed edge.

IMDB-BINARY, **IMDB-MULTI**, and **REDDIT-BINARY** (Morris et al., 2020a) are benchmark datasets for graph-level classification. The IMDB datasets provide actor collaboration graphs for binary and multi-class movie genre prediction, respectively. **REDDIT-BINARY** consists of graphs from online discussion threads, with the task being to classify the community type.

D.2 BASELINES

In our study, we mainly compare our SDE-GNN model with 15 baselines. The descriptions of these baselines are as following:

GCN (Kipf & Welling, 2017), **GAT** (Velickovic et al., 2018), **GraphSAGE** (Hamilton et al., 2017), and **GIN** (Xu et al., 2019) are classic graph neural network models that rely on a basic message-passing framework without explicit distance-aware mechanisms, which limits their expressive capacity.

GraphGPS (Rampásek et al., 2022) proposes a modular, scalable graph Transformer architecture with linear complexity by decoupling local edge-based message passing from global Transformer attention mechanisms, while maintaining universal function approximation capabilities.

Exphormer (Shirzad et al., 2023) constructs a graph Transformer architecture with linear complexity and superior theoretical properties by introducing a sparse attention mechanism based on virtual global nodes and expander graphs. With core components including expander graphs and virtual nodes, the computational complexity is reduced from the quadratic level of traditional graph Transformers to linear scaling with graph size.

SGFormer (Wu et al., 2023) addresses the scalability challenge of graph Transformers by proposing a simplified architecture that discards complex multi-layer multi-head attention structures. Instead, it leverages a single-layer attention mechanism capable of propagating information across arbitrary nodes with linear complexity, eliminating the need for positional encodings, feature/graph pre-processing, or augmented loss functions.

GOAT (Kong et al., 2023) introduces a scalable global Transformer framework that enables adaptive learning of homophily/heterophily relationships by allowing each node to attend to all nodes through an approximate, theoretically-justified global self-attention mechanism. This design eliminates domain-specific inductive biases while remaining compatible with both homophilic and heterophilic settings.

ESAN (Bevilacqua et al., 2022) enhances MPNN expressiveness by aggregating information from subgraphs (instead of direct node interactions), using predefined policies and equivariant architectures to distinguish MPNN-indistinguishable graphs. Theoretically grounded in novel 1D-WL test variants, it introduces stochastic subgraph sampling to reduce computational costs while maintaining flexibility in subgraph selection and design.

Graphormer (Ying et al., 2021a) addresses the challenge of adapting the powerful Transformer architecture for graph representation learning. It augments the standard Transformer with several novel structural encoding methods—such as Centrality Encoding, Spatial Encoding, and Edge Encoding—which are incorporated directly into the self-attention mechanism.

GRIT (Ma et al., 2023) introduces a graph Transformer that integrates inductive biases without message-passing, addressing prior limitations. Key innovations include: relative positional encodings initialized via random walk probabilities, an attention mechanism updating both node and node-pair representations, and degree injection in each layer. Theoretically, GRIT captures shortest-path distances and propagation matrices.

SPE (Huang et al., 2024) addresses Laplacian-based positional encoding limitations in graph Transformers by introducing a stable framework using a "soft partition" mechanism guided by eigenvalues. It resolves non-uniqueness and instability from hard eigenspace partitions, ensuring robustness to perturbations while maintaining universal expressiveness for symmetry-preserving functions.

NeuralWalker (Chen et al., 2024) bridges the gap between message-passing GNNs (which excel at local relations but fail at long-range dependencies) and graph Transformers (which oversimplify graphs via fixed-length vectors). It combines random walks (for long-range context) with local message passing by treating random walks as sequences and leveraging sequence models to capture dependencies.

DE-GNN (Li et al., 2020) introduces Distance Encoding (DE), a general class of structure-related features that measures the distance from a target set of nodes to every other node in the graph, using metrics like shortest path or PageRank. DE enhances GNNs by being used either as extra node features or as controllers of the message aggregation step, allowing the model to gain expressive power beyond the 1-WL test.

RD-WL. **GD-WL** (Zhang et al., 2023) overcomes the expressive limitation beyond the standard Weisfeiler-Lehman test by introducing a principled framework based on generalized graph distances, such as resistance distance (for RD-WL), to systematically enhance expressive power.

D.3 HYPERPARAMETERS

For our SDE-GNN model, we employ grid search to identify optimal hyperparameters. The reduced dimension of dimensionality reduction is set to 1000 for two large-scale datasets—ogbn-arxiv and twitch-gamers. While for other medium- or small-scale datasets, the dimensionality reduction is not applicable. The power of adjacency matrix k is searched within the range of 1 to 5. The maximum order of polynomial m is searched within the range of 1 to 5. The number of layers L is searched within the range of 1 to 10 (For CLUSTER and PATTERN datasets, the range is 1 to 40). The node embedding dimension is searched in $\{64, 128, 256, 512, 1024\}$. And the learning rate is searched in $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$.

For the baseline models to compare with, we also employ grid search to optimize key hyperparameters and report the performance of the best configuration.

D.4 EVALUATION SETTINGS

We uniformly employ classification accuracy on the test set to measure models' performance. Note that for CLUSTER and PATTERN dataset, since their node classes are imbalanced, we use the average node-level accuracy weighted with respect to the class sizes following (Dwivedi et al., 2023). The IMDB-BINARY, IMDB-MULTI, and REDDIT-BINARY datasets are randomly partitioned into training, validation, and test sets in a 5:2:3 ratio. While the data splits for the remaining datasets are kept consistent with those of the baseline methods (Dwivedi et al., 2023; Luo et al., 2024). Each model is run five times under identical conditions, with the mean and standard deviation of the results reported for statistical reliability. Evaluations are conducted on a machine with 192GB RAM, two 28-core Intel Xeon CPUs (2.2GHz), and an NVIDIA GeForce RTX 3090 GPU (24GB memory).

E MORE EXPERIMENTAL RESULTS

E.1 TIME AND MEMORY CONSUMPTION AT DIFFERENT NODE SCALES

To assess the scalability of different methods, we generate a series of synthetic graphs with node counts ranging from 200 to 200,000. The edges are randomly generated, maintaining a fixed 10:1 edge-to-node ratio. As illustrated in Figure 4, both the inference time and GPU memory footprint of our SDE-GNN model scale linearly with the number of nodes, empirically verifying the linear computation complexity of SDE-GNN.

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

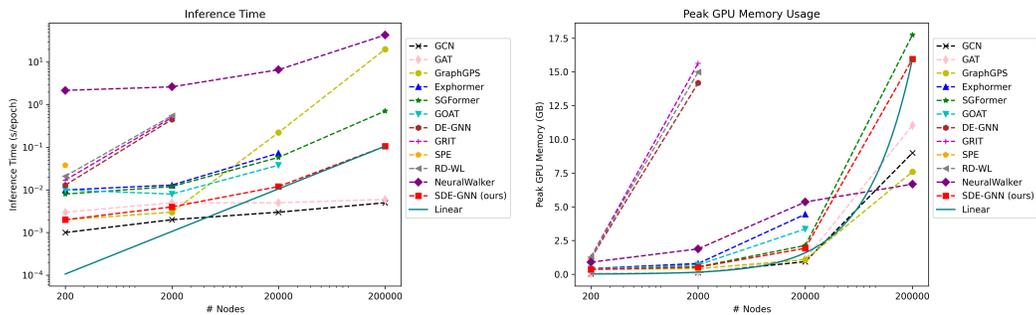


Figure 4: The comparison of models’ inference time (s/epoch) and GPU memory usage (GB) at different node scales.

E.2 ANALYSIS OF DIMENSIONALITY REDUCTION METHODS

We also investigate the impact of the selected dimensionality reduction method and the target dimension size on ogbn-arxiv dataset. As illustrated in Figure 5, the results show that 1) Randomized SVD (Halko et al., 2011) consistently outperforms other techniques, namely Count-Sketch (Charikar et al., 2002) and Gaussian Random Projection (Johnson et al., 1984), and 2) Model performance is positively correlated with the target dimensionality, while reducing the dimension from 1,000 to 200 does not lead to a significant drop in performance.

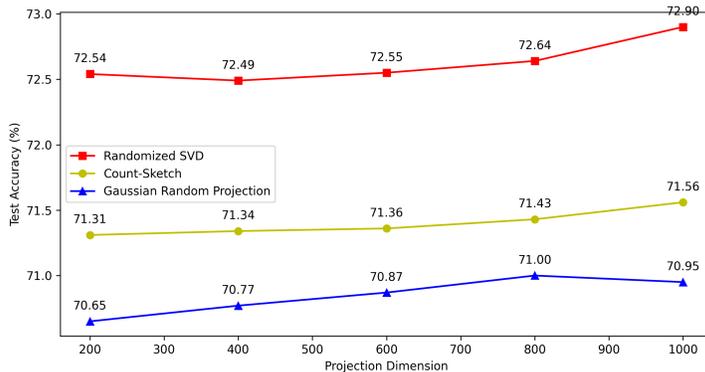


Figure 5: Comparison of different dimensionality reduction methods and reduced dimension.

F LIMITATION

In this paper, we show that the adopted adjacency-power-based distance feature theoretically upper bounds the expressiveness of the eigenspace projection distance. However, whether this bound is strict or the two are equally expressive remains an open question, which we leave for future investigation.

G USE OF LLMs

This article used LLMs as a tool to enhance writing clarity and correct grammatical errors. It is important to note that the fundamental ideas and the overall framework of this research are the original contributions of the authors.