# Inductive Structural Role Embedding on Large-scale Graphs

### Anonymous Author(s) Affiliation Address email

### Abstract

1	Graph embedding methods have been proposed to extract structural identities in a
2	graph, but most of the existing structural role embedding methods are transductive
3	and the embedding cannot be generalized to unseen nodes. Here we introduce In-
4	SuRE, an inductive method to embed nodes' structural roles. Instead of leveraging
5	a diffusion process on the entire graph, we characterize a local diffusion kernel
6	with two learnable parameters, the local neighborhood radius and corresponding
7	diffusion scale. With the two parameters, the embedding of unseen nodes can be
8	efficiently generated based on their neighborhood topology. InSuRE is computa-
9	tionally efficient, provides discriminative structural features to improve GNN's
10	expressive power, and outperforms baseline methods in empirical experiments.

# 11 **1 Introduction**

Nodes in a graph may have different structural roles, reflected by their neighborhood topology. For example, managers in a company are usually the hubs of the communication network, whereas other staff are usually non-hub nodes. Identifying the structural roles of the nodes helps understand their identities and behaviors in many real-world applications.

Graph structured data analysis has traditionally fo-16 cused on predefined metrics [1, 2, 3], such as struc-17 tural hole value [4]. However, these approaches only 18 19 extract structural information based on the predefined metrics. Recent graph representation learning ap-20 proaches extract nodes' structural information in a 21 data-driven fashion. Especially, they encode the high-22 dimensional, non-Euclidean structural information 23 with a low-dimensional Euclidean vector, which can 24 be used in downstream analysis, such as node classi-25 fication and link prediction. 26



- Although most of the node embedding methods [5, 6,
  7, 8, 9, 10, 11, 12, 13, 14, 15, 16] focus on proximity
  preserving (i.e., nodes nearby have similar represen-
- <sup>30</sup> tations), several structural embedding methods [17,

Figure 1: Similar structural roles of nodes uand v are identified with their diffusion patterns  $h_u$  and  $h_v$  captured as distributions.

- 18, 19, 20, 21, 22, 23, 24, 25] have been proposed to embed nodes in terms of their structural roles (see Fig. 1). However, all of them except Role2Vec [20] are *transductive* in the sense that they directly
- 33 generate node embedding from the entire graph (see details in Sec. 2). When unseen nodes are added
- to the graph, the representations of the nodes in the entire graph need to be updated. Therefore, trans-
- <sup>35</sup> ductive embedding methods are not useful for analyzing evolving graphs, where the graph structure
- 36 changes over time and unseen nodes constantly appear, such as social networks on Facebook and

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.

posts on Twitter. Moreover, existing structural embedding approaches are computationally prohibitive 37 in large-scale graphs. For example, struc2vec [18] has a time complexity of  $O(|\mathcal{V}|^3)$ . Role2Vec has a 38 time complexity of  $O(|\mathcal{V}|^2)$ . SEGK [23] has a time complexity of  $O(|\mathcal{V}||\mathcal{E}|)$ . GraphWave's space 39 complexity is  $O(|\mathcal{V}|^2)$ . 40 To overcome these challenges, we propose InSuRE, an Inductive Structural Role Embedding approach 41 to embed the structural roles of the nodes in a graph. InSuRE leverages a local diffusion kernel to 42 characterize the local neighborhood topology of the nodes. The local diffusion kernel is characterized 43 by a local neighborhood radius (i.e., the neighborhood hop number) and a diffusion scale, whose 44 best configurations are identified by two variance-based approaches. We demonstrate that a local 45 diffusion kernel, with the optimized neighborhood radius and diffusion scale, is as effective as a 46 global diffusion kernel in terms of capturing structural roles, and makes InSuRE both inductive and 47 computationally efficient. We also prove that GNN coupled with InSuRE's structural embeddings 48

- <sup>49</sup> provides more expressive power. The contribution of our work is *three-fold*.
- We develop InSuRE, an inductive structural role embedding method. InSuRE extracts node structural roles via a local diffusion kernel (Def. 3.2) consisting of two parameters, a local neighborhood radius and a diffusion scale. The optimized local diffusion kernel embeds the structural roles of unseen nodes efficiently.
- InSuRE's time complexity is linear with the number of edges (Sec. 3.5), which is lower than state-of-the-art structural embedding methods, making it applicable to large-scale graphs.
- We demonstrate that using InSuRE's structural embedding as node features improves the expressive power of MP-GNNs theoretically (Sec. 3.6) and empirically (Sec. 4.5).

# 58 2 Related work

**Structural role node embedding.** Most of the existing structural role embedding methods follow two 59 major steps to embed nodes — first construct a node-feature matrix by extracting nodes' structural 60 features, and then embed the nodes based on either feature-based matrix factorization or feature-based 61 random walk [26, 27]. As one example of feature-based matrix factorization, RolX [17, 28] factorizes 62 the matrix with node topological features to assign nodes with a mixed-membership across the 63 predetermined structural roles. HONE [21] and SEGK [23] also embeds nodes' structural roles based 64 on factorizing a set of motif-based matrices and the structural feature matrix generated by graph 65 kernels, respectively. Another line of structural role embedding methods performs feature-based 66 random walk using the skip-gram model [29, 30]. Struc2vec [18] embeds node structural roles by 67 performing biased random walks on a multilayer graph, with each layer capturing the structural 68 similarities of the neighboring hops between the nodes. RiWalk [31] performs random walks on 69 a role-identification graph constructed by graph kernels. DRNE [22] utilizes a normalized LSTM 70 model recursively to embed nodes with regular equivalence. Struc2gauss [25] leverages Gaussian 71 embedding to embed nodes based on node structural features. However, all of the aforementioned 72 methods require either manual feature engineering or explicit predefined functions to extract nodes' 73 structural features. Moreover, they are all transductive. Recently, GraphWave [19] proposes to extract 74 structural features based on a predefined diffusion process [32, 33, 34] on the graph. It interprets heat 75 diffusion on graphs in the context of graph wavelet transform [35], treats the node diffusion patterns 76 as random variables, and embeds them using a characteristic function [36]. However, GraphWave is 77 78 transductive due to the nature of spectral methods. Graph spectrum changes completely when new nodes are added, thus changing the predefined diffusion process. By contrast, InSuRE characterizes a 79 local diffusion process, and thus can efficiently embed the structural roles of the unseen nodes. 80

Inductive node embedding. While most embedding methods are transductive, a few inductive 81 approaches have been proposed recently. Planetoid [37] is the first inductive embedding approach 82 while it does not use any structural features. GraphSAGE [11] embeds nodes inductively by aggre-83 gating information from nodes' neighborhoods, while it is proximity-preserving. Recently, both 84 SPINE [24] and Role2Vec [20] embed node structural identities inductively. They first generate 85 structural features and perform attributed random walks to embed nodes. While SPINE approximates 86 the Rooted-PageRank score from the local structure to ensure inductiveness, it still relates to local 87 proximity. Role2Vec embeds unseen nodes by mapping them to pre-trained node types and embeds 88 them using the learned skip-gram model. However, it becomes inappropriate when the unseen nodes 89 do not match the pre-trained node types. [38] also propose an inductive embedding method for 90 temporal graphs based on causal anonymous walks. 91

# 92 **3** Methods

Since a node's structural role is mostly reflected by its local neighborhood, InSuRE (Algo. 1)
extracts node structural roles via a *local diffusion kernel* (Def. 3.2), characterized by two parameters,
namely *a local neighborhood radius* and *a diffusion scale*. Section 3.1 introduces basic notations and
definitions. Section 3.2 introduces the embedding procedure given two specific parameters. Section
3.3 and 3.4 describes the optimization of the local neighborhood radius and the diffusion scale,
respectively. Section 3.5 discusses InSuRE's time complexity. Section 3.6 proves that InSuRE's
structural embedding improves GNN's expressive power.

### 100 3.1 Preliminaries

Suppose that  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  is an undirected graph, where  $\mathcal{V} = [n]$  and  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  denote the node set and the edge set respectively. Each node  $u \in \mathcal{V}$  has a node feature vector  $\mathbf{x}_u \in \mathbb{R}^d$ . Let  $\mathbf{A}$  be the adjacency matrix and  $\mathbf{D}$  be the degree matrix. The random walk adjacency matrix is  $\mathbf{A}_{rw} = \mathbf{D}^{-1}\mathbf{A}$ . Only a subgraph is selected in training to accelerate the parameter optimization, with a sampling rate  $\alpha$ .  $T = \{t_1, \ldots, t_d\}$  is a set of evenly spaced sampling time points used to sample an empirical characteristic function.

**Regular equivalence.** We follow Everett and Borgatti [39] and Rossi et al. [40] and consider two regularly equivalent nodes to have the same structural role. Formally, let  $r_u$  be the structural role of node u and  $\mathcal{N}_u$  denote the set of u's neighbors. Two nodes u and v are regularly equivalent iff they are connected to regularly equivalent nodes (i.e.,  $r_u = r_v \iff \{r(s) | s \in \mathcal{N}_u\} = \{r(t) | t \in \mathcal{N}_v\}, \forall u, v \in \mathcal{V}$ ).

Graph Neural Networks (GNNs). GNNs learn a vector representation of node v,  $\mathbf{m}_v$ , within the message passing framework. During each message-passing iteration of a messgae-passing GNN (MP-GNN),  $\mathbf{m}_v$  is updated as

$$\mathbf{m}_{u}^{(k+1)} = \text{UPDATE}(\mathbf{m}_{u}^{(k)}, \text{AGG}(\{\mathbf{m}_{v}^{(k)}, \forall v \in \mathcal{N}(u)\})), \tag{1}$$

where  $\mathbf{m}_{u}^{0} = \mathbf{x}_{u} \in \mathbb{R}^{d}, \forall u \in \mathcal{V}$ . UPDATE is a differentiable function (e.g., ReLU) and AGG is a differentiable permutation invariant function (e.g., summation).

Weisfeiler-Lehman (WL) test. Testing for graph isomorphism is to declare whether two graphs have identical graph structure while only differing in the node ordering in the adjacency matrix. No general polynomial time algorithms are known for the problem [41]. The Weisfeiler-Lehman (WL) test of graph isomorphism [42] is an effective and efficient algorithm for approximate isomorphism testing [43]. Its simplest form, commonly known as the 1-WL, iteratively aggregates the labels of nodes and their neighborhoods, and hashes the aggregated labels into unique new labels. The algorithm declares two graphs to be isomorphic iff the labels of the nodes between the two graphs are identical.

**GNNs and** 1-**WL test.** Despite the success of GNNs, recent works have proved that the expressive power of MP-GNNs is upper bounded by the 1-WL test as follows.

**Theorem 3.1.** [44, 45]. Consider a K-layer MP-GNN with each layer following the form of Eq. 1. Suppose the initial input node feature is discrete, i.e.,  $\mathbf{m}_u^0 = \mathbf{x}_u \in \mathbb{Z}^d, \forall u \in \mathcal{V}$ . Then  $\mathbf{m}_u^K \neq \mathbf{m}_v^K$ only if nodes u and v have different labels after K iterations of the WL algorithm.

# 128 **3.2** Structural role embedding from local heat diffusion

This section introduces the embedding 129 procedure given the local neighborhood 130 radius and the diffusion scale. The pro-131 cedure has two steps — the first step is 132 to extract each node's structural infor-133 mation through a local diffusion process, 134 and the second step is to embed their dif-135 fusion pattern. 136

# 137 Step 1. Generating diffusion pattern

InSuRE starts with modeling a diffusion
process on a graph to extract each node's
structural information. The global diffu-

141 sion kernel is defined as follows.

### Algorithm 1 InSuRE

- Input: G={V, E}, initial diffusion scale s<sub>0</sub>, sampling rate α, standard deviation threshold η, sampling time point set T
- Output: Optimal local neighborhood radius k<sup>\*</sup><sub>η</sub>, optimal scale s<sup>\*</sup>, node embedding φ<sub>i</sub> for v<sub>i</sub> ∈ V
- 3: Compute transition matrix  $P=D^{-1/2}AD^{-1/2}$
- 4: Generate subgraph  $\mathcal{G}_{sub}$  and transition matrix  $\mathbf{P}_{sub}$
- 5: Optimize local neighborhood radius  $k_n^*$  (Sec. 3.3)
- 6: Learn the optimal diffusion scale  $s^*$  (Sec. 3.4)
- 7: Compute optimal diffusion pattern  $H^*$  (Eq. 2)
- 8: Generate embedding  $\phi_a$  for  $v_a \in \mathcal{V}$  (Eq. 3)

**Definition 3.1.** (Global diffusion kernel) A diffusion kernel operator is defined as  $H_s = e^{-s\mathbf{L}_{rw}}$ . As the diffusion kernel operator reformulated as a continuous-time random walk [32], the global diffusion kernel,  $H_s$ , is defined as  $H_s = e^{-s} \sum_{k=0}^{\infty} \frac{s^k}{k!} \mathbf{P}^k$ , where  $\mathbf{P}$  is the random walk adjacency matrix  $\mathbf{A}_{rw}$ ,  $\mathbf{L}_{rw} = \mathbf{I} - \mathbf{D}^{-1} \mathbf{A}$  denotes the normalized Laplacian matrix, k is the steps of random walk, and s is the diffusion scale.

The global diffusion kernel describes a flow of heat energy at a diffusion scale *s* over the entire graph [46]. Each node's neighborhood structural information is characterized by the energy received from its neighborhoods, namely the diffusion pattern. Directly computing the global heat kernel is both time and space expensive. However, it is easily verified that  $H_s$ 's coefficient series (i.e.,  $\left\{\frac{e^s s^k}{k!}\right\}_{k=0,\ldots,\infty}$ , s > 0) converges rapidly and the random walk on a graph converges to a stationary distribution as  $k \to \infty$ . The local diffusion kernel is therefore defined as the sum of  $H_s$ 's first k' terms.

**Definition 3.2.** (Local diffusion kernel) A local diffusion kernel  $\tilde{H}(s,k') = e^{-s} \sum_{k=0}^{k'} \frac{(s)^k}{k!} \mathbf{P}^k$ . is characterized by two parameters, namely a local neighborhood radius k' and a diffusion scale s. The

local neighborhood radius k' limits a diffusion process to neighboring nodes within k' hops, and the diffusion scale s controls the heat propagation.

157 With the optimized local neighborhood radius  $k_{\eta}^{*}$  and the diffusion scale  $s^{*}$ , the optimal local diffusion

kernel is  $\tilde{H}(s^*, k^*_{\eta}) = e^{-s^*} \sum_{k=0}^{k^*_{\eta}} \frac{(s^*)^k}{k!} P^k$ . We concatenate the local diffusion kernel within  $k^*_{\eta}$ -hop neighborhoods and obtain the final optimal diffusion pattern as

$$H^* = [\tilde{H}(s^*, 1)^T, \dots, \tilde{H}(s^*, k^*_\eta)^T]^T.$$
(2)

The yielded diffusion pattern for a node u is denoted as  $h_u = H^* \delta_u$ , where  $\delta_u = \mathbb{I}_u$  denotes the onehot vector of node u. Since the local diffusion kernel can be easily computed given  $k_{\eta}^*$  and  $s^*$ , InSuRE performs well in an inductive learning setting (see Sec. 4.2 and 4.4). The local diffusion calculation also makes InSuRE computationally efficient (see Sec. 3.5). Note that we use the symmetric random walk adjacency matrix  $\mathbf{A}_{sym} = \mathbf{D}^{-1/2} \mathbf{A} \mathbf{D}^{-1/2}$ , which yields better performance in experiments.

### 165 Step 2. Embedding structural roles

Given nodes u and v have the same structural role, the neighborhood subgraphs rooted at u and 166 v are isomorphic. As a node's neighborhood topology is encoded in the yielded local diffusion 167 pattern, nodes u and v also have "isomorphic" local diffusion pattern (i.e.,  $\exists \pi : \mathcal{V} \rightarrow \mathcal{V}, \pi(h_u) = h_v$ ). 168 However, directly solving the diffusion pattern matching problem is computationally prohibitive [47]. 169 To compare the patterns more efficiently, we regard each node u's diffusion pattern  $h_u$  as a random 170 variable following a distribution, with  $h_{mu}$  treated as if they were sampled observations. Then its 171 empirical characteristic function is used to characterize the distribution. Specifically, for a node u, the characteristic function of  $h_u$  is defined as  $z_u(t) = \mathbb{E}[e^{ith_u}]$ . It is a Fourier transform of the probability distribution  $h_a$ , and thus fully describes  $h_a$ . Based on the calculated coefficients of node u's diffusion 172 173 174 pattern, the empirical characteristic distribution of  $h_u$  is defined as  $z_u(t) = \frac{1}{n} \sum_{m=1}^{k_{\eta}^* \cdot n} e^{ith_{mu}}$ . By 175 Euler's formula (i.e.,  $e^{ix} = \cos x + i \sin x$ ), we divide  $z_u(t)$  into real and imaginary parts. Sampling 176 the characteristic function at d points in set T, node u's embedding is generated by concatenating all 177 the values [19], 178

$$\phi_{u} = \left[ \operatorname{Re}\left( z_{u}\left( t_{k} \right) \right), \operatorname{Im}\left( z_{u}\left( t_{k} \right) \right) \right]_{t_{k} \in T}, \tag{3}$$

179 where 
$$\operatorname{Re}(z_u(t_k)) = \frac{1}{k_{\eta}^{*,n}} \sum_{m=1}^{k_{\eta}^{*,n}} \cos(t_k h_{mu})$$
, and  $\operatorname{Im}(z_u(t_k)) = \frac{1}{k_{\eta}^{*,n}} \sum_{m=1}^{k_{\eta}^{*,n}} \sin(t_k h_{mu})$ 

For unseen nodes: InSuRE generates the induced  $k_{\eta}^*$ -hop neighborhood subgraph centered at each unseen node, and applies Eq. 2 to embed their structural roles.

### 182 **3.3 Optimizing local neighborhood radius** k

The local neighborhood radius k' controls the range of heat diffusion. Too small or too large k' leads to limited or redundant diffusion, and cannot extract the proper structural information. As k' also denotes the number of random walk steps, the random walk converges to a stationary distribution when k' goes to infinity [48]. The stationary distribution is identical for all nodes in a graph, and cannot differentiate the diffusion pattern between nodes. Thus, optimizing k' is related to choosing a proper random walk convergence level. Since the initial transition matrix  $\mathbf{P}^0$  is usually far from

convergence, we derive the optimal  $k_{\eta}^*$  via r(k'), the ratio between the average standard deviation 189 of row vectors in the k'-step random walk matrix  $\mathbf{P}^{k'}$  and the average standard deviation of row 190 vectors in the initial random walk matrix  $\mathbf{P}^0$ , namely  $r(k') = \frac{\operatorname{std}(\mathbf{P}^{k'})}{\operatorname{std}(\mathbf{P}^0)} = \frac{\sum \operatorname{std}(p_i^{k'})}{\sum \operatorname{std}(p_0^{k'})}$ , where  $p_i^{k'} \in \mathbb{R}^n$ 191 denotes the *i*-th row of  $\mathbf{P}^{k'}$ ,  $i=0, 1, \ldots, \infty$ . Given a predefined threshold  $\eta$ , the random walk with k'192 steps is not informative on identifying nodes' diffusion pattern if  $r(k') < \eta$ . Let F(k') denote the 193 194

set that contains all the values of k' satisfying  $r(k') \ge \eta$  (i.e.,  $F(k') = \{k' | r(k') \ge \eta\}$ ), and  $k_{\eta}^*$  is computed as  $k_{\eta}^* = \max_{k'} F(k')$ . The relationship between k and the random walk convergence level is formally derived as follows (with proof in the supplementary material). 195 196

Theorem 3.2. Consider a random walk on a connected graph. For every initial probability distribu-197 tion  $p_0$  and every  $k \ge 0$ , we have  $||p_k - \pi|| \le \sqrt{\frac{\max_u d_u}{\min_u d_u}} \lambda_2^k$ , where  $p_k$  denotes the distribution after k steps,  $\pi$  denotes the stationary distribution,  $d_u$  denotes node u's degree, and  $\lambda_2$  denotes the second 198

199 largest eigenvalue of the transition matrix **P**. 200

**Corollary 3.3.** Given a ratio threshold  $\eta$ , we have  $r(k) = \eta \iff k \leq \log_{\lambda_2}(\frac{\eta_k^2 \operatorname{Var}(\mathbf{P}^0)}{d})$ , where  $d = \frac{\max_u d_u}{\min_u d_u}$ ,  $\operatorname{Var}(\mathbf{P}^0) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n \|\mathbf{P}_{ij}^0 - p_i^0\|_2^2$ , and  $p_i^0 = \frac{1}{n} \sum_{j=1}^n \mathbf{P}_{ij}^0$ ,  $i, j \in \mathcal{V}$ . 201 202

Note that the convergence of random walk using transition matrix  $\mathbf{P}=\mathbf{A}_{rw}$  is equivalent to the convergence of random walk using  $\mathbf{P}=\mathbf{A}_{sym}$ , as  $\mathbf{A}_{sym}^{k}=\mathbf{D}^{-1/2}\mathbf{A}\mathbf{A}_{rw}^{k-1}\mathbf{D}^{-1/2}$ . Since  $\lambda_{2} < 1$ , Corollary 3.3 indicates that the value of  $\eta$  is proportional to the selection of  $k_{\eta}^{*}$ . We refer to  $\eta$  as a 203 204 205 hyper-parameter in our method, with a larger  $\eta$  corresponding to global node structural roles, and a 206 smaller  $\eta$  corresponding to local node structural roles. 207

#### **Optimizing diffusion scale** s 3.4 208

The scale parameter s determines the heat propagation time. Too small or too large s leads to limited 209 diffusion or identical diffusion. For a given  $k_{\eta}^*$ , the scale parameter s is optimized to best identify 210 nodes' structural roles. Thus, we optimize s by maximizing the covariance between nodes' diffusion 211 pattern, which is denoted as  $cov(h_i, h_j) = \frac{1}{n} \sum_{a} \sum_{b} h_{ai} h_{bj} - \bar{h}_i \bar{h}_j$ , where  $\bar{h}_i$  and  $\bar{h}_j$  denote the means 212 of  $h_i$  and  $h_j$ . However, directly computing the covariance has a time complexity  $O(n^4)$ . To reduce 213 computation, the covariance between the sample mean of  $h_i$ ,  $i \in \mathcal{V}$ , is maximized with a loss function 214

$$L(s) = -n \sum_{i=1}^{n} \sum_{j=1}^{n} \operatorname{cov}(\bar{h}_i, \bar{h}_j) = -(\sum_{i=1}^{n} \bar{h}_i^2 - n\bar{h}^2),$$
(4)

where  $h = \frac{1}{n} \sum_{i} \bar{h}_{i}$ . The simplification is based on the idea that covariance between the diffusion 215 pattern of the nodes can be estimated by the variance between their sample means, with details in the 216 supplementary material. To accelerate the optimization, only a subset of the nodes are considered 217 during training. We first randomly select an initial node, traversing the graph with a breadth-first 218 search (BFS) strategy and sampling nodes with the sampling rate  $\alpha$ , until  $\alpha * |\mathcal{V}|$  nodes are collected. 219 This BFS strategy preserves nodes' structure since nodes with more neighbors still have larger degrees 220 after sampling. 221

#### Time complexity 3.5 222

Since InSuRE utilizes a subset of nodes during training, its time complexity stems from computing 223 the optimal local diffusion pattern (line 7 in Algo. 1), which requires multiplying transition matrices of different orders (i.e.,  $\mathbf{P}^k \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$ ,  $k=0, 1, \dots, k_{\eta}^*$ ). Sparse matrix multiplication is used to reduce 224 225 computation. Let  $\tau(\mathbf{P})$  denote the number of the non-zero elements in  $\mathbf{P}$ , and the expected number 226 of multiplication operations between  $\mathbf{P}^k$  and  $\mathbf{P}$  is  $\tau(\mathbf{P}^k)\tau(\mathbf{P})/|\mathcal{V}|$  [49]. For large sparse graphs, 227  $|\mathcal{E}|=c|\mathcal{V}|$ , where c is a constant. The *total time complexity* for embedding existing nodes is  $O(|\mathcal{E}|)$ . 228 For a set of unseen nodes  $\mathcal{V}' = \{v'_1, \dots, v'_m\}$ , subgraphs  $\mathcal{G}'_i = \{\mathcal{V}'_i, \mathcal{E}'_i\}$ , centered at these nodes are generated. The time complexity for *embedding unseen nodes* is  $O(|\mathcal{V}'|^2)$ , where  $\mathcal{V}'$  denotes the 229 230 average size of  $\mathcal{G}'_i$  for each node  $v_i \in \mathcal{V}'$ . 231

#### 3.6 Improving a GNN's expressive power with InSuRE's structural embedding 232

MP-GNN's limited expressive power is inherited from the 1-WL test which does not capture distance 233 information between the nodes. Recently, Li et al. [50] proposed DE-GNN to improve MP-GNN's 234

expressive power by encoding distance information as extra node features. Due to InSuRE's efficiency and effectiveness in embedding nodes' structural roles, we couple MP-GNN with InSuRE's structural embedding, termed as InSuRE-GNN, and theoretically demonstrate that InSuRE-GNN increases the expressive power of MP-GNNs by proving that InSuRE-GNN is more powerful than DE-GNN with a simple aggregation function. We introduce definitions and theorems below, and save the detailed derivations in the supplementary material.

**Definition 3.3.** (*DE-GNN*) [50]. Given a subset of nodes  $S \subseteq \mathcal{V}$  (e.g., a set including node u and its neighbors), the distance encoding (*DE*) for node u is

$$\zeta(u \mid S) = AGG(\{\zeta(u \mid v) \mid v \in S\}),\tag{5}$$

where  $\zeta$  and AGG are differentiable permutation invariant functions.  $\zeta(u \mid v)$  chracterizes a certain distance between nodes v and u. For instance,

$$\zeta(u \mid v) = g\left(\ell_{uv}\right), \ell_{uv} = \left(1, \left(\mathbf{A}_{rw}\right)_{uv}, \dots, \left(\mathbf{A}_{rw}^{k}\right)_{uv}\right).$$
(6)

245 *DE-GNN couples MP-GNN with DE as extra features. DE-GNN is called proper if UPDATE, AGG in* 246 *Eq. 1, AGG in Eq. 5, and g in Eq. 6 are all injective mapping if the input features are all countable.* 

Li et al. [50] formally prove that DE-GNN is more powerful than MP-GNN in distinguishing 247 structural identities in regular graphs without node attributes. Nevertheless, the simple aggregation 248 function AGG:  $\mathbb{R}^{|S|} \times \mathbb{R}^{|S|} \to \mathbb{R}$  used in DE-GNN satisfies the injective assumption only if |S|=1, 249 where DE contains limited distance information from itself. When |S|>2, AGG becomes surjective, 250 and nodes with different distance information may be mapped together. Thus, AGG has limited 251 capability to generate distinguishable node structural features. By contrast, InSuRE uses the empirical 252 characteristic function to embed the distribution of the distance information between nodes (i.e., the 253 diffusion pattern of nodes). Note that InSuRE's diffusion pattern generation is equivalent to set g 254 in Eq. 6 as heat kernel operator and the obtained diffusion pattern containing distance information 255 between nodes. As the empirical characteristic function uniquely characterizes the distribution of 256 node diffusion pattern, it is capable of distinguishing non-isomorphic node diffusion patterns, which 257 is formally stated as follows. 258

**Theorem 3.4.** The empirical characteristic function of nodes u and v are the same if and only if their diffusion pattern  $h_u$ ,  $h_v$  are isomorphic (i.e.,  $\exists \pi : \mathcal{V} \to \mathcal{V}, \pi(h_u) = h_v$ ).

Based on Theorem 3.4, we conclude as follows.

**Corollary 3.5.** InSuRE-GNN is more expressive than DE-GNN with a simple set aggregation function.

263

# **264 4 Experiments**

We evaluate InSuRE's performance in five experiments with different purposes. The first experiment (Sec. 4.1) investigates the effectiveness of a local diffusion kernel comparing to a global diffusion kernel in structural role embedding. Then, we use a simulated barbell graph to examine InSuRE's structural embedding in both transductive and inductive settings (Sec. 4.2). We further test InSuRE's embedding in transductive (Sec. 4.3) and inductive (Sec. 4.4) node classification tasks. Finally, we evaluate InSuRE-GNN's expressive power in node classification tasks (Sec. 4.5).

Baselines. We compare InSuRE with struc2vec [18], GraphWave [19], and Role2Vec [20] in terms 271 of their embedding quality and inductive ability. We also compare InSuRE with RolX [17] in the 272 simulated experiments and node2vec [7] in the transductive node classification. In the experiments for 273 evaluating InSuRE's expressive power, seven baselines are chosen. The first three methods, Struc2vec-274 GNN, GraphWave-GNN, Role2Vec-GNN, denote MP-GNN coupled with Struc2vec, GraphWave, 275 Role2Vec node embeddings. Regarding DE-GNN, we choose its landing probability (LP) variant, as 276 it achieved the best performance on the node classification task reported in the original paper. The 277 remaining three baselines, GCN [51], GraphSAGE [11], and GIN [44] are representative MP-GNNs, 278 and node degrees are used as the initial node features. Parameters are tuned for all the baselines. 279

### 280 4.1 Comparing local and global diffusion kernels in structural role embedding

To investigate a local diffusion kernel's capability of identifying nodes' structural roles, we compare the performance of a local diffusion kernel with different neighborhood radii (k = 1, 2, 3, and 4) and the global diffusion kernel used in GraphWave on a simulated barbell graph. The simulated barbell graph is a symmetric, unweighted graph with two 70-vertex cliques, connected by a path of seven nodes (Fig. 2). It has six structural equivalent classes indicated by different colors. Fig. 3 shows the 2D-PCA projection of the embedding obtained by a local diffusion kernel with neighborhood radii k = 1, 2, 3, and 4, and the global diffusion kernel (i.e.,  $k = \infty$ ) used by GraphWave. The optimal local neighborhood radius  $k_{\eta}^*$  learned by InSuRE is 4, which produces a embedding nearly the same as that from a



Figure 2: Barbell graph with node colors

global diffusion kernel. The result demonstrates that a denoting structural roles. local diffusion kernel is as effective as a global diffusion kernel in extracting nodes' structural roles.



Figure 3: 2D-PCA projections of the embedding obtained by a local diffusion kernel with neighborhood radius k = 1, 2, 3, 4, and the global diffusion kernel (i.e.,  $k = \infty$ ) used by GraphWave.



Figure 4: 2D-PCA projections of node structural embedding from different methods in (A) transductive setting and (B) inductive setting.

293

292

### **4.2** Comparing structural role embedding in both transductive and inductive settings

We use the same barbell graph and evaluate all the approaches in both transductive and inductive 295 settings. In the transductive setting, InSuRE and baselines are applied to the entire graph, and the 296 structural embedding from different approaches is visualized and compared. In the inductive setting, 297 we randomly remove one yellow clique node and its edges, and regard it as an unseen node. All 298 approaches are applied to the remainder graph to generate their embedding. The unseen node together 299 with its neighborhood is then fed to all the approaches to generate its embedding. A successful 300 inductive embedding approach should project the unseen node and its structural equivalent class 301 in the remainder graph to the same place. The structural embedding from different approaches is 302 visualized in Fig. 4. The colors of the nodes represent their ground-truth structural roles. In the 303 transductive embedding task (Fig. 4A), InSuRE, GraphWave, and struc2vec correctly identify all 304 of the six structural roles in the graph, whereas RolX only identifies three structural roles, with 305 purple/red/brown nodes projected together, green/yellow nodes projected together, and blue-green 306 nodes projected as the third role. Role2Vec fails to classify the structural roles between green and 307 yellow nodes as well as between purple and red nodes. Moreover, InSuRE and GraphWave accurately 308 capture the similarity between different structural roles — yellow and green nodes with similar roles 309 are projected closer, and the other four similar roles are projected closer. In the inductive experiment 310 (Fig. 4B), only InSuRE and Role2Vec correctly classify the unseen node into the right class. However, 311 InSuRE still identifies all of the six structural roles whereas Role2Vec fails to classify the structural 312 roles between nodes, especially for the classes other than yellow. 313

MEASURE	GRAPH	INSURE	STRUC2VEC	GRAPHWAVE	ROLE2VEC	NODE2VEC	
	Brazil-Air	$\textbf{82.59} \pm \textbf{8.93}$	$76.67\pm8.77$	$78.52\pm7.73$	$43.33\pm10.74$	$43.70\pm10.86$	
ACCURACY (%)	EUROPE-AIR	$\textbf{58.00} \pm \textbf{4.00}$	$57.63 \pm 4.09$	$52.37 \pm 5.63$	$32.00\pm 6.35$	$38.12\pm4.55$	
	USA-AIR	$61.51 \pm 2.52$	$60.42 \pm 2.21$	$56.17 \pm 3.18$	$48.57\pm3.43$	$52.98 \pm 2.84$	
	ACTOR	$45.56\pm0.69$	$\textbf{46.17} \pm \textbf{0.99}$	$45.45 \pm 1.26$	$34.50\pm9.31$	$35.73\pm0.99$	
	Brazil-Air	1.05	23.68	0.34	12.29	6.01	
RUN TIME	EUROPE-AIR	3.21	155.02	3.54	52.30	23.85	
(SECOND)	USA-AIR	15.63	874.71	18.11	143.54	84.92	
	ACTOR	453.28	10490.72	559.52	2324.80	460.55	
	Brazil-Air	5.15	8.89	19.46	15.18	17.79	
PEAK MEMORY	EUROPE-AIR	48.83	9.59	174.88	101.27	89.97	
USAGE (MB)	USA-AIR	305.99	11.66	1482.80	286.20	271.77	
	ACTOR	6378.05	28.95	58457.95	5711.69	1460.35	

TABLE 1: ACCURACY, RUN TIME, AND PEAK MEMORY USAGE OF DIFFERENT METHODS ON TRANSDUCTIVE NODE CLASSIFICATION.

 TABLE 2: PERFORMANCE ON INDUCTIVE NODE

 CLASSIFICATION WITH FACEBOOK SOCIAL CIRCLE

 NETWORK

TABLE 3: PERFORMANCE ON IN-<br/>DUCTIVE NODE CLASSIFICATION WITH<br/>TWITTER SOCIAL CIRCLE NETWORK

MEASURE	CLASSIFIER	INSURE	ROLE2VEC	GRAPHWAVE	STRUC2VEC	I WITTER SOCIAL CIRCLE NETWORK
						MEASURE\CLASSIFIER KNN NAIVE BAYES
MEAN	KNN	0.998	0.996	0.932	0.996	
ACCURACY	NAIVE BAYES	1.000	0.996	0.787	0.996	MEAN ACCURACY 0.9895 0.9840
E SCOPE	KNN	0.667	0.000	0.00	0.000	F <sub>1</sub> -score 0.3313 0.5068
F1-SCORE	NAIVE BAYES	1.000	0.000	0.027	0.000	

### 314 4.3 Testing structural embedding in a transductive node classification task

We evaluate the structural embedding from different approaches in transductive node classification 315 tasks under the assumption that a good structural embedding should correctly identify nodes' structural 316 labels. Four real datasets, Brazil-Airports, Europe-Airports, USA-Airports [18], and Actor co-317 occurrence [52], are chosen as their labels indicate the structural roles. The first three datasets are 318 flight traffic networks, where nodes correspond to airports and edges indicate the existence of non-stop 319 flights. Each node is assigned a label according to its level of activity. In the Actor co-occurrence 320 network dataset, nodes represent actors and edges indicate the co-occurrence on the same Wikipedia 321 page. Each node is assigned with a label according to its influence level. For all datasets, we use 80%322 and 20% dataset splitting for training and testing. All the embedding methods are applied to learn 323 each node's embedding. We train a L2-regularized logistic regression model based on embeddings of 324 training nodes. Accuracy score is used to evaluate the performance of methods, run time and peak 325 memory usage are also recorded. 326

Results (averaged from 10 replicates with different random seeds) are summarized in Table 1. InSuRE 327 outperforms the baseline methods at almost all datasets, which indicates that InSuRE is capable 328 of embedding node structural roles effectively. Especially, InSuRE outperforms GraphWave in the 329 three airport traffic networks, and achieves slightly better performance in the Actor co-occurrence 330 network. This demonstrates local diffusion kernel (used by InSuRE) is more powerful than the global 331 diffusion kernel (used by GraphWave) in identifying node structural roles in the real-world datasets, 332 as the global diffusion kernel may be interfered with by noisy information. Struc2vec has the best 333 performance among all the baseline methods, which implies the structural kernel is effective to extract 334 structural information. Role2Vec and node2vec do not perform well in the node classification task. 335 Regarding run time and peak memory usage, InSuRE is the fastest among all the methods and has 336 comparable peak memory usage with Role2Vec and node2vec. Although struc2vec uses the least 337 memory, it takes too much time. GraphWave has a comparable run time with InSuRE, while it 338 consumes too much space. Therefore, only InSuRE achieves a decent node classification accuracy in 339 340 an effective manner for embedding large-scale graphs.

### **4.4** Using structural embedding for an inductive node classification task

In an inductive node classification task, we are given a graph  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}\)$ , where each node has its ground-truth structural label. To obtain the test set, we randomly remove 25% of the nodes in  $\mathcal{G}$  and their edges, and use  $\mathcal{V}_{test}$  to denote the set of removed nodes. The remainder graph and node set are denoted as  $\mathcal{G}_{train}$  and  $\mathcal{V}_{train}$ . For each node  $v_i \in \mathcal{V}_{test}$ , we construct its subgraph  $\mathcal{G}'_i$  based on its neighborhood. Then we merge all the subgraphs  $\mathcal{G}'_i$  to one graph called  $\mathcal{G}_{test}$ . For transductive methods (i.e., GraphWave, struc2vec), we apply them on  $\mathcal{G}_{train}$  and  $\mathcal{G}_{test}$  to embed nodes in  $\mathcal{V}_{train}$  and  $\mathcal{V}_{test}$  separately. For inductive methods (i.e., InSuRE and Role2Vec), we learn

TABLE 4: PERFORMANCE OF DIFFERENT GNN MODELS ON NODE CLASSIFICATION TASK (MEASURE IN ACCURACY SCORE (%)).

GRAPH	INSURE-GNN	STRUC2VEC-GNN	GRAPHWAVE-GNN	ROLE2VEC-GNN	DE-GNN	GCN	GRAPHSAGE	GIN
BRAZIL-AIR Europe-Air USA-Air	$\begin{array}{c} 83.85 \pm 7.26 \\ 65.00 \pm 4.33 \\ 66.98 \pm 2.02 \end{array}$	$\begin{array}{c} 72.11 \pm 7.41 \\ 61.25 \pm 8.00 \\ 62.61 \pm 2.36 \end{array}$	$\begin{array}{c} 76.15 \pm 9.39 \\ 61.00 \pm 3.74 \\ 63.10 \pm 3.71 \end{array}$	$\begin{array}{c} 65.39 \pm 11.54 \\ 59.50 \pm 7.56 \\ 61.09 \pm 4.54 \end{array}$	$\begin{array}{c} 75.46 \pm 6.64 \\ 61.25 \pm 3.75 \\ 62.77 \pm 3.68 \end{array}$	$\begin{array}{c} 73.85 \pm 9.23 \\ 59.25 \pm 5.24 \\ 58.65 \pm 3.85 \end{array}$	$\begin{array}{c} 73.08 \pm 11.53 \\ 59.50 \pm 5.78 \\ 50.85 \end{array}$	$\begin{array}{c} 73.24 \pm 7.19 \\ 60.50 \pm 5.45 \\ 62.68 \pm 2.67 \end{array}$

parameters of the embedding function based on  $\mathcal{G}_{train}$ , and embed nodes in  $\mathcal{V}_{train}$  and  $\mathcal{V}_{test}$  using 349 the learned model. Based on the obtained embedding, we predict the label of each node in  $\mathcal{V}_{test}$  with 350 two classifiers — one k-nearest neighbor classifier and one naive Bayesian classifier. A successful 351 inductive embedding should accurately predict the ground-truth node label based on the trained model. 352 Mean accuracy and  $F_1$ -score are used to measure their performance. We use the Facebook and Twitter 353 social circle network datasets from Stanford Network Analysis Project database. Both of the graphs 354 are unweighted. Facebook network consists of 4,039 nodes and 88,234 edges, and Twitter network 355 contains 81,306 nodes and 1,768, 149 edges. Nodes in the two networks represent individuals, and 356 edges represent the pairwise friendship between two individuals. Facebook network includes ten 357 ego-networks, and Twitter network contains 973 ego-networks. Each ego-network includes one ego 358 node and different numbers of alter nodes. We assign labels to every node in the network according 359 to their structural identities in the community, namely, "ego" or "alter". 360

Results on Facebook and Twitter networks are summarized in Tables 2 and 3. Note that for Twitter 361 network, we could not run GraphWave on it with a workstation equipped with 64 GB memory due to 362 its high space complexity, and we could not obtain the result of Role2Vec and struc2vec after running 363 for one day due to their high time complexity. Regarding Facebook network, all methods appear to 364 be accurate due to the imbalance between the two classes (i.e., only ten "ego" nodes in the 4,039365 nodes). Six of the "ego" nodes are in the train graph and the other four are in the test graph. InSuRE 366 achieves both high accuracy and  $F_1$ -scores, while Role2Vec and struc2vec yield high accuracy but 367 poor  $F_1$ -scores, since they fail to identify any of the "ego" nodes in the test set. This demonstrates 368 that InSuRE correctly identifies the structural roles of the unseen nodes even though training samples 369 are very few. Regarding Twitter network, InSuRE yields high mean accuracy and decent  $F_1$ -score 370 with a 1-hour run time, which indicates its capability of embedding nodes in large-scale graphs. 371

### 372 4.5 Evaluating MP-GNN coupled with structural embeddings on node classification task

The expressive power of different GNN models is evaluated on the node classification task under the assumption that a GNN with higher expressive power yields better performance in classifying the nodes' structural roles. The three flight traffic network datasets, Brazil-Airports, Europe-Airports, USA-Airports, are chosen, with the detailed introduction in Sec. 4.3. For all the datasets, 80% of data are used for training, 10% of data are used for validation, and the remaining 10% are used for testing. Accuracy score is used to evaluate the performance of each method.

The results averaged from 10 replicates with different random seeds are summarized in Table 379 4. InSuRE-GNN outperforms the baseline models in all datasets, which indicates that it is more 380 powerful than DE-GNN with a simple aggregation function. DE-GNN performs better than MP-GNN 381 baselines, which demonstrates its expressive power in structural representations. GIN achieves the 382 best performance among the three MP-GNN baselines, which is consistent with the theory in Xu 383 et al. [44]. Although struc2vec-GNN and GraphWave-GNN achieve competitive performance with 384 DE-GNN, it is either time-consuming or space-consuming for the two methods to generate node 385 embeddings. By contrast, InSuRE's efficiency makes InSuRE-GNN applicable to large-scale graphs. 386

# 387 5 Discussion

In conclusion, InSuRE innovatively uses a local diffusion kernel to capture node structural roles,
 which is proven to be effective and efficient in inductive node embedding. Theoretical and empirical
 results suggest that InSuRE's embedding as node features increases MP-GNN's expressive power.

Limitations of our work In the diffusion scale parameter optimization, how different choices of sampling points influence structural embedding needs further investigation. Two alternative methods can be considered in future study. One is to measure the pairwise distribution gap with distribution discrepancy measures such as maximum mean discrepancy and Wasserstein distance. The other one is to embed distributions with the first *l* cumulants of empirical cumulant-generating function. However, the former one may be space-consuming, whereas the latter may be time-consuming.

397 Negative societal impacts of our work Our work has no potential negative societal impacts.

# **398** References

- Elizabeth A Leicht, Petter Holme, and Mark EJ Newman. "Vertex similarity in networks". In:
   *Physical Review E* 73.2 (2006), p. 026120.
- [2] Laura A Zager and George C Verghese. "Graph similarity scoring and matching". In: *Applied mathematics letters* 21.1 (2008), pp. 86–94.
- [3] Ruoming Jin, Victor E Lee, and Hui Hong. "Axiomatic ranking of network role similarity". In:
   *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2011, pp. 922–930.
- [4] Ronald S Burt. *Structural holes: The social structure of competition*. Harvard university press, 2009.
- <sup>408</sup> [5] Daniel L Sussman et al. "A consistent adjacency spectral embedding for stochastic blockmodel <sup>409</sup> graphs". In: *Journal of the American Statistical Association* 107.499 (2012), pp. 1119–1128.
- [6] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. "Deepwalk: Online learning of social
   representations". In: *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2014, pp. 701–710.
- [7] Aditya Grover and Jure Leskovec. "node2vec: Scalable feature learning for networks". In:
   *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2016, pp. 855–864.
- [8] Jian Tang et al. "Line: Large-scale information network embedding". In: *Proceedings of the* 24th international conference on world wide web. 2015, pp. 1067–1077.
- <sup>418</sup> [9] Shaosheng Cao, Wei Lu, and Qiongkai Xu. "Deep neural networks for learning graph repre-<sup>419</sup> sentations". In: *Thirtieth AAAI conference on artificial intelligence*. 2016.
- [10] Daixin Wang, Peng Cui, and Wenwu Zhu. "Structural deep network embedding". In: *Proceed- ings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2016, pp. 1225–1234.
- <sup>423</sup> [11] Will Hamilton, Zhitao Ying, and Jure Leskovec. "Inductive representation learning on large <sup>424</sup> graphs". In: *Advances in Neural Information Processing Systems*. 2017, pp. 1024–1034.
- Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. "metapath2vec: Scalable representation learning for heterogeneous networks". In: *Proceedings of the 23rd ACM SIGKDD international conference on knowledge discovery and data mining*. 2017, pp. 135–144.
- [13] Giang Hoang Nguyen et al. "Continuous-time dynamic network embeddings". In: *Companion Proceedings of the The Web Conference 2018*. 2018, pp. 969–976.
- [14] Aynaz Taheri, Kevin Gimpel, and Tanya Berger-Wolf. "Learning to represent the evolution of
   dynamic graphs with recurrent models". In: *Companion Proceedings of The 2019 World Wide* Web Conference. 2019, pp. 301–307.
- [15] Sandro Cavallari et al. "Embedding both finite and infinite communities on graphs [application notes]". In: *IEEE Computational Intelligence Magazine* 14.3 (2019), pp. 39–50.
- [16] Di Jin et al. "Latent network summarization: Bridging network embedding and summarization".
  In: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. 2019, pp. 987–997.
- [17] Keith Henderson et al. "Rolx: structural role extraction & mining in large graphs". In: *Proceed- ings of the 18th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM. 2012, pp. 1231–1239.
- [18] Leonardo FR Ribeiro, Pedro HP Saverese, and Daniel R Figueiredo. "struc2vec: Learning node representations from structural identity". In: *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM. 2017, pp. 385–394.
- [19] Claire Donnat et al. "Learning structural node embeddings via diffusion wavelets". In: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM. 2018, pp. 1320–1329.
- <sup>447</sup> [20] Nesreen K Ahmed et al. "Learning role-based graph embeddings". In: *arXiv preprint* <sup>448</sup> *arXiv:1802.02896* (2018).
- [21] Ryan A. Rossi et al. "A Structural Graph Representation Learning Framework". In: *Proceedings of the 13th International Conference on Web Search and Data Mining*. WSDM
  '20. Houston, TX, USA: Association for Computing Machinery, 2020, pp. 483–491. ISBN: 9781450368223. DOI: 10.1145/3336191.3371843.

- Ke Tu et al. "Deep recursive network embedding with regular equivalence". In: *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*.
   2018, pp. 2357–2366.
- [23] Giannis Nikolentzos and Michalis Vazirgiannis. "Learning structural node representations
   using graph kernels". In: *IEEE Transactions on Knowledge and Data Engineering* (2019).
- Iunliang Guo, Linli Xu, and Jingchang Liu. "Spine: structural identity preserved inductive
   network embedding". In: *arXiv preprint arXiv:1802.03984* (2018).
- Yulong Pei et al. "struc2gauss: Structural role preserving network embedding via Gaussian
   embedding". In: *Data Mining and Knowledge Discovery* 34 (2020), pp. 1072–1103.
- 462 [26] Ryan A Rossi et al. "From community to role-based graph embeddings". In: *arXiv preprint* 463 *arXiv:1908.08572* (2019).
- <sup>464</sup> [27] Ryan A Rossi and Nesreen K Ahmed. "Role discovery in networks". In: *IEEE Transactions* on *Knowledge and Data Engineering* 27.4 (2014), pp. 1112–1131.
- Keith Henderson et al. "It's who you know: graph mining using recursive structural features".
   In: *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. 2011, pp. 663–671.
- [29] Tomas Mikolov et al. "Efficient estimation of word representations in vector space". In: *arXiv preprint arXiv:1301.3781* (2013).
- <sup>471</sup> [30] Tomas Mikolov et al. "Distributed representations of words and phrases and their composition-<sup>472</sup> ality". In: *Advances in Neural Information Processing Systems*. 2013, pp. 3111–3119.
- [31] Xuewei Ma et al. "RiWalk: fast structural node embedding via role identification". In: 2019 *IEEE International Conference on Data Mining (ICDM)*. IEEE. 2019, pp. 478–487.
- [32] Fan Chung. "The heat kernel as the pagerank of a graph". In: *Proceedings of the National Academy of Sciences* 104.50 (2007), pp. 19735–19740.
- [33] Ronald R Coifman and Stéphane Lafon. "Diffusion maps". In: *Applied and Computational Harmonic Analysis* 21.1 (2006), pp. 5–30.
- [34] Risi Imre Kondor and John Lafferty. "Diffusion kernels on graphs and other discrete structures".
  In: *Proceedings of the 19th international conference on machine learning*. Vol. 2002. 2002,
  pp. 315–22.
- [35] David K Hammond, Pierre Vandergheynst, and Rémi Gribonval. "Wavelets on graphs via spectral graph theory". In: *Applied and Computational Harmonic Analysis* 30.2 (2011), pp. 129–150.
- 485 [36] Eugene Lukacs. "Characteristic functions". In: (1970).
- [37] Zhilin Yang, William W Cohen, and Ruslan Salakhutdinov. "Revisiting semi-supervised
   learning with graph embeddings". In: *arXiv preprint arXiv:1603.08861* (2016).
- [38] Yanbang Wang et al. "Inductive Representation Learning in Temporal Networks via Causal
   Anonymous Walks". In: *arXiv preprint arXiv:2101.05974* (2021).
- [39] Martin G Everett and Stephen P Borgatti. "Regular equivalence: General theory". In: *Journal* of Mathematical Sociology 19.1 (1994), pp. 29–52.
- [40] Ryan A Rossi et al. "On proximity and structural role-based embeddings in networks: Miscon ceptions, techniques, and applications". In: *ACM Transactions on Knowledge Discovery from Data (TKDD)* 14.5 (2020), pp. 1–37.
- [41] László Babai. "Graph isomorphism in quasipolynomial time". In: *Proceedings of the forty- eighth annual ACM symposium on Theory of Computing*. 2016, pp. 684–697.
- <sup>497</sup> [42] Boris Weisfeiler and Andrei Leman. "The reduction of a graph to canonical form and the <sup>498</sup> algebra which appears therein". In: *NTI, Series* 2.9 (1968), pp. 12–16.
- [43] László Babai and Ludik Kucera. "Canonical labelling of graphs in linear average time". In:
   20th Annual Symposium on Foundations of Computer Science (sfcs 1979). IEEE. 1979, pp. 39–
   46.
- [44] Keyulu Xu et al. "How powerful are graph neural networks?" In: *arXiv preprint arXiv:1810.00826* (2018).
- [45] Christopher Morris et al. "Weisfeiler and leman go neural: Higher-order graph neural networks".
   In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 33. 01. 2019, pp. 4602–4609.

- [46] David I Shuman et al. "The emerging field of signal processing on graphs: Extending high dimensional data analysis to networks and other irregular domains". In: *IEEE signal processing magazine* 30.3 (2013), pp. 83–98.
- [47] Chang Wang et al. "Manifold alignment". In: *Manifold Learning: Theory and Applications*.
   CRC Press, 2011.
- [48] László Lovász et al. "Random walks on graphs: A survey". In: *Combinatorics, Paul erdos is eighty* 2.1 (1993), pp. 1–46.
- [49] Keivan Borna and Sohrab Fard. "A note on the multiplication of sparse matrices". In: *Open Computer Science* 4.1 (2014), pp. 1–11.
- [50] Pan Li et al. "Distance encoding: Design provably more powerful neural networks for graph
   representation learning". In: *arXiv preprint arXiv:2009.00142* (2020).
- <sup>518</sup> [51] Thomas N Kipf and Max Welling. "Semi-supervised classification with graph convolutional <sup>519</sup> networks". In: *arXiv preprint arXiv:1609.02907* (2016).
- Lei Tang and Huan Liu. "Relational learning via latent social dimensions". In: *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*.
   2009, pp. 817–826.
- 523 Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? [Yes] See Section ??.
- Did you include the license to the code and datasets? [No] The code and the data are proprietary.
- Did you include the license to the code and datasets? [N/A]

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

535 1. For all authors...

536 537	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
538	(b) Did you describe the limitations of your work? [Yes]
539	(c) Did you deserve are minimations of your work? [Yes]
540 541	<ul><li>(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]</li></ul>
542	2. If you are including theoretical results
543	(a) Did you state the full set of assumptions of all theoretical results? [Yes]
544	(b) Did you include complete proofs of all theoretical results? [Yes]
545	3. If you ran experiments
546 547	(a) Did you include the code, data, and instructions needed to reproduce the main experi- mental results (either in the supplemental material or as a URL)?
548 549	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes]
550 551	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes]
552 553	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
554	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
555	(a) If your work uses existing assets, did you cite the creators? [Yes]

556	(b) Did you mention the license of the assets? [Yes]
557	(c) Did you include any new assets either in the supplemental material or as a URL? [No]
558	We did not use any these kinds of data
559	(d) Did you discuss whether and how consent was obtained from people whose data you're
560	using/curating? [No] We did not use any these kinds of data
561	(e) Did you discuss whether the data you are using/curating contains personally identifiable
562	information or offensive content? [No] We did not use any these kinds of data
563	5. If you used crowdsourcing or conducted research with human subjects
564	(a) Did you include the full text of instructions given to participants and screenshots, if
565	applicable? [N/A]
566	(b) Did you describe any potential participant risks, with links to Institutional Review
567	Board (IRB) approvals, if applicable? [N/A]
568	(c) Did you include the estimated hourly wage paid to participants and the total amount
569	spent on participant compensation? [N/A]