

# A SIMPLE AND GENERAL GRAPH NEURAL NETWORK WITH STOCHASTIC MESSAGE PASSING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

1 Graph neural networks (GNNs) are emerging machine learning models on graphs.  
 2 One key property behind the expressiveness of existing GNNs is that the learned  
 3 node representations are permutation-equivariant. Though being a desirable prop-  
 4 erty for certain tasks, however, permutation-equivariance prevents GNNs from  
 5 being proximity-aware, i.e., preserving the walk-based proximities between pairs  
 6 of nodes, which is another critical property for graph analytical tasks. On  
 7 the other hand, some variants of GNNs are proposed to preserve node prox-  
 8 imities, but they fail to maintain permutation-equivariance. How to empower  
 9 GNNs to be proximity-aware while maintaining permutation-equivariance re-  
 10 mains an open problem. In this paper, we propose Stochastic Message Passing  
 11 (SMP), a general and simple GNN to maintain both proximity-awareness and  
 12 permutation-equivariance properties. Specifically, we augment the existing GNNs  
 13 with stochastic node representations learned to preserve node proximities. Though  
 14 seemingly simple, we prove that such a mechanism can enable GNNs to preserve  
 15 node proximities in theory while maintaining permutation-equivariance with cer-  
 16 tain parametrization. Extensive experimental results demonstrate the effectiveness  
 17 and efficiency of SMP for tasks including node classification and link prediction.

## 18 1 INTRODUCTION

19 Graph neural networks (GNNs), as generalizations of neural networks in analyzing graphs, have  
 20 attracted considerable research attention. GNNs have been widely applied to various applications  
 21 such as social recommendation (Ma et al., 2019), physical simulation (Kipf et al., 2018), and protein  
 22 interaction prediction (Zitnik & Leskovec, 2017).

23 One key property of most existing GNNs is permutation-equivariance, i.e., if we randomly permu-  
 24 tate the IDs of nodes while maintaining the graph structure, the representations of nodes in GNNs  
 25 are permuted accordingly. Mathematically, permutation-equivariance reflects one basic symmet-  
 26 ric group of graph structures. Although it is a desirable property for tasks such as node or graph  
 27 classification (Keriven & Peyré, 2019; Maron et al., 2019b), permutation-equivariance also prevents  
 28 GNNs from being proximity-aware, i.e., permutation-equivariant GNNs cannot preserve walk-based  
 29 proximities between nodes such as the shortest distance or high-order proximities (see Theorem 1).

30 Pairwise proximities between nodes are crucial for graph analytical tasks such as link predic-  
 31 tion (Hu et al., 2020; You et al., 2019). To enable a proximity-aware GNN, Position-aware GNN  
 32 (P-GNN) (You et al., 2019)<sup>1</sup> proposes a sophisticated GNN architecture and shows better perfor-  
 33 mance for proximity-aware tasks. But P-GNN needs to explicitly calculate the shortest distance be-  
 34 tween nodes and its computational complexity is unaffordable for large graphs. Moreover, P-GNN  
 35 completely ignores the permutation-equivariance property. Therefore, it cannot produce satisfactory  
 36 results when permutation-equivariance is helpful.

37 In real-world scenarios, both proximity-awareness and permutation-equivariance are indispensable  
 38 properties for GNNs. Firstly, different tasks may require different properties. For example, recom-  
 39 mendation applications usually require the model to be proximity-aware (Konstas et al., 2009) while  
 40 permutation-equivariance is a basic assumption in centrality measurements (Borgatti, 2005). Even

<sup>1</sup>In (You et al., 2019), the authors consider the special case of shortest distance between nodes and name such property as “position-aware”. In this paper, we consider a more general case of any walk-based proximity.

41 for the same task, different datasets may have different requirements on these two properties. Taking  
 42 link prediction as an example, we observe that permutation-equivariant GNNs such as GCN (Kipf &  
 43 Welling, 2017) or GAT (Velickovic et al., 2018) show better results than P-GNN in coauthor graphs,  
 44 but the opposite in biological graphs (please see Section 5.2 for details). Unfortunately, in the current  
 45 GNN frameworks, these two properties are contradicting, as we show in Theorem 1. Whether there  
 46 exists a general GNN to be proximity-aware while maintaining permutation-equivariance remains  
 47 an open problem.

48 In this paper, we propose Stochastic Message Passing (SMP), a general and simple GNN to pre-  
 49 serve both proximity-awareness and permutation-equivariance properties. Specifically, we augment  
 50 the existing GNNs with stochastic node representations learned to preserve proximities. Though  
 51 seemingly simple, we prove that our proposed SMP can enable GNNs to preserve walk-based prox-  
 52 imities in theory (see Theorem 2 and Theorem 3). Meanwhile, SMP is equivalent to a permutation-  
 53 equivariant GNN with certain parametrization and thus is at least as powerful as those GNNs in  
 54 permutation-equivariant tasks (see Remark 1). Therefore, SMP is general and flexible in handling  
 55 both proximity-aware and permutation-equivariant tasks, which is also demonstrated by our exten-  
 56 sive experimental results. Besides, owing to the simple structure, SMP is computationally efficient,  
 57 with a running time roughly the same as those of the most simple GNNs such as SGC (Wu et al.,  
 58 2019) and is at least an order of magnitude faster than P-GNN on large graphs. Ablation studies  
 59 further show that a linear instantiation of SMP is expressive enough as adding extra non-linearities  
 60 does not lift the performance of SMP on the majority of datasets. Our contributions are as follows.

- 61 • We propose SMP, a simple and general GNN to handle both proximity-aware and permutation-  
 62 equivariant graph analytical tasks.
- 63 • We prove that SMP has theoretical guarantees in preserving walk-based proximities and is at  
 64 least as powerful as the existing GNNs in permutation-equivariant tasks.
- 65 • Extensive experimental results demonstrate the effectiveness and efficiency of SMP. We show  
 66 that a linear instantiation of SMP is expressive enough on the majority of datasets.

## 67 2 RELATED WORK

68 We briefly review GNNs and their permutation-equivariance and proximity-awareness property.

69 The earliest GNNs adopt a recursive definition of node states (Scarselli et al., 2008; Gori et al.,  
 70 2005) or a contextual realization (Micheli, 2009). GGS-NNs (Li et al., 2016) replace the recursive  
 71 definition with recurrent neural networks (RNNs). Spectral GCNs (Bruna et al., 2014) defined graph  
 72 convolutions using graph signal processing (Shuman et al., 2013; Ortega et al., 2018) with Cheb-  
 73 Net (Defferrard et al., 2016) and GCN (Kipf & Welling, 2017) approximating the spectral filters us-  
 74 ing a  $K$ -order Chebyshev polynomial and the first-order polynomial, respectively. MPNNs (Gilmer  
 75 et al., 2017), GraphSAGE (Hamilton et al., 2017), and MoNet (Monti et al., 2017) are proposed  
 76 as general frameworks by characterizing GNNs with a message-passing function and an updating  
 77 function. More advanced variants such as GAT (Velickovic et al., 2018), JK-Nets (Xu et al., 2018b),  
 78 GIN (Xu et al., 2018a), and GraphNets (Battaglia et al., 2018) follow these frameworks.

79 Li *et al.* (Li et al., 2018), Xu *et al.* (Xu et al., 2018a), Morris *et al.* (Morris et al., 2019), and  
 80 Maron *et al.* (Maron et al., 2019a) show the connection between GNNs and the Weisfeiler-Lehman  
 81 algorithm (Shervashidze et al., 2011) of graph isomorphism tests, in which permutation-equivariance  
 82 holds a key constraint. Maron *et al.* (Maron et al., 2019b) and Keriven *et al.* (Keriven & Peyré, 2019)  
 83 analyze the permutation-equivariance property of GNNs more theoretically. To date, most of the  
 84 existing GNNs are permutation-equivariant and thus are not proximity-aware. The only exception  
 85 is P-GNN (You et al., 2019), which proposes to capture the positions of nodes using the relative  
 86 distance between the target node and some randomly chosen anchor nodes. However, P-GNN cannot  
 87 satisfy permutation-equivariance and is computationally expensive.

88 Very recently, motivated by enhancing the expressive power of GNNs in graph isomorphism tests  
 89 and distributed computing literature (Angluin, 1980; Linial, 1992; Naor & Stockmeyer, 1995),  
 90 some studies suggest assigning unique node identifiers for GNNs (Loukas, 2020) such as one-hot  
 91 IDs (Murphy et al., 2019) or random numbers (Dasoulas et al., 2019; Sato et al., 2020; Corso et al.,  
 92 2020). For example, Sato *et al.* (Sato et al., 2020) novelly show that random numbers can enhance  
 93 GNNs in tackling two important graph-based NP problems with a theoretical guarantee, namely the

94 minimum dominating set and the maximum matching problem, and Fey *et al.* (Fey *et al.*, 2020) empirically show the effectiveness of random features in the graph matching problem. Our work differs 95 in that we systematically study how to preserve permutation-equivariance and proximity-awareness 96 simultaneously in a simple yet effective framework, which is a new topic different from these existing 97 works. Besides, we theoretically prove that our proposed method can preserve walk-based 98 proximities by using the random projection literature. We also demonstrate the effectiveness of our 99 method on various large-scale benchmarks for both node- and edge-level tasks, while no similar 100 results are reported in the literature. 101

102 The design of our method is also inspired by the random projection literature in dimensionality reduction 103 (Vempala, 2005) and to the best of our knowledge, we are the first to study random projection 104 in the scope of GNNs. More remotely, our definition of node proximities is inspired and inherited 105 from graph kernels (Gärtner *et al.*, 2003; Borgwardt & Kriegel, 2005), network embedding (Perozzi 106 *et al.*, 2014; Grover & Leskovec, 2016), and general studies of graphs (Newman, 2018).

### 107 3 MESSAGE-PASSING GNNs

108 We consider a graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$  where  $\mathcal{V} = \{v_1, \dots, v_N\}$  is the set of  $N = |\mathcal{V}|$  nodes,  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$  109 is the set of  $M = |\mathcal{E}|$  edges, and  $\mathbf{F} \in \mathbb{R}^{N \times d_0}$  is a matrix of  $d_0$  node features. The adjacency matrix 110 is denoted as  $\mathbf{A}$ , where its  $i^{\text{th}}$  row,  $j^{\text{th}}$  column and an element denoted as  $\mathbf{A}_{i,:}$ ,  $\mathbf{A}_{:,j}$ , and  $\mathbf{A}_{i,j}$ , 111 respectively. In this paper, we assume the graph is unweighted and undirected. The neighborhood 112 of node  $v_i$  is denoted as  $\mathcal{N}_i$  and  $\tilde{\mathcal{N}}_i = \mathcal{N}_i \cup \{v_i\}$ .

113 The existing GNNs usually follow a message-passing framework (Gilmer *et al.*, 2017), where the  $l^{\text{th}}$  114 layer adopts a neighborhood aggregation function  $\text{AGG}(\cdot)$  and an updating function  $\text{UPDATE}(\cdot)$ :

$$\mathbf{m}_i^{(l)} = \text{AGG}(\{\mathbf{h}_j^{(l)}, \forall j \in \tilde{\mathcal{N}}_i\}), \mathbf{h}_i^{(l+1)} = \text{UPDATE}([\mathbf{h}_i^{(l)}, \mathbf{m}_i^{(l)}]), \quad (1)$$

115 where  $\mathbf{h}_i^{(l)} \in \mathbb{R}^{d_l}$  is the representation of node  $v_i$  in the  $l^{\text{th}}$  layer,  $d_l$  is the dimensionality, and  $\mathbf{m}_i^{(l)}$  116 are the messages. We also denote  $\mathbf{H}^{(l)} = [\mathbf{h}_1^{(l)}, \dots, \mathbf{h}_N^{(l)}]$  and  $[\cdot, \cdot]$  is the concatenation operation. The 117 node representations are initialized as node features, i.e.,  $\mathbf{H}^{(0)} = \mathbf{F}$ . We denote a GNN following 118 Eq. (1) with  $L$  layers as a parameterized function as follows<sup>2</sup>:

$$\mathbf{H}^{(L)} = \mathcal{F}_{\text{GNN}}(\mathbf{A}, \mathbf{F}; \mathbf{W}), \quad (2)$$

119 where  $\mathbf{H}^{(L)}$  are final node representations learned by the GNN and  $\mathbf{W}$  denotes all the parameters.

120 One key property of the existing GNNs is permutation-equivariance.

121 **Definition 1** (Permutation-equivariance). Consider a graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$  and any permutation 122  $\mathcal{P} : \mathcal{V} \rightarrow \mathcal{V}$  so that  $G' = (\mathcal{V}, \mathcal{E}', \mathbf{F}')$  has an adjacency matrix  $\mathbf{A}' = \mathbf{PAP}^T$  and a feature matrix 123  $\mathbf{F}' = \mathbf{PF}$ , where  $\mathbf{P} \in \{0, 1\}^{N \times N}$  is the permutation matrix corresponding to  $\mathcal{P}$ , i.e.,  $\mathbf{P}_{i,j} = 1$  iff 124  $\mathcal{P}(v_i) = v_j$ . A GNN satisfies permutation-equivariance if the node representations are equivariant 125 with respect to  $\mathcal{P}$ , i.e.,

$$\mathbf{PF}_{\text{GNN}}(\mathbf{A}, \mathbf{F}; \mathbf{W}) = \mathcal{F}_{\text{GNN}}(\mathbf{PAP}^T, \mathbf{PF}; \mathbf{W}). \quad (3)$$

126 It is known that GNNs following Eq. (1) are permutation-equivariant (Maron *et al.*, 2019b).

127 **Definition 2** (Automorphism). A graph  $G$  is said to have (non-trivial) automorphism if there exists 128 a non-identity permutation matrix  $\mathbf{P} \neq \mathbf{I}_N$  so that  $\mathbf{A} = \mathbf{PAP}^T$  and  $\mathbf{F} = \mathbf{PF}$ . We denote the 129 corresponding automorphic node pairs as  $\mathcal{C}_G = \bigcup_{\mathbf{P} \neq \mathbf{I}_N} \{(i, j) | \mathbf{P}_{i,j} \neq 0, i \neq j\}$

130 **Corollary 1.** Using Definition 1 and 2, if a graph has automorphism, a permutation-equivariant 131 GNN will produce identical node representations for automorphic node pairs:

$$\mathbf{h}_i^{(L)} = \mathbf{h}_j^{(L)}, \forall (i, j) \in \mathcal{C}_G. \quad (4)$$

132 Since the node representations are used for downstream tasks, the corollary shows that permutation- 133 equivariant GNNs cannot differentiate automorphic node pairs. A direct consequence of Corollary 1 134 is that permutation-equivariant GNNs cannot preserve walk-based proximities between pairs 135 of nodes. The formal definitions are as follows.

<sup>2</sup>Since the final layer of GNNs is task-specific, e.g., a softmax layer for node classification or a readout layer for graph classification, we only consider the GNN architecture to its last hidden layer.

136 **Definition 3** (Walk-based Proximities). For a given graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$ , we use a matrix  $\mathbf{S} \in$   
 137  $\mathbb{R}^{N \times N}$  to denote walk-based proximities between pairs of nodes defined as:

$$\mathbf{S}_{i,j} = \mathcal{S}(\{v_i \rightsquigarrow v_j\}), \quad (5)$$

138 where  $v_i \rightsquigarrow v_j$  denotes walks from node  $v_i$  to  $v_j$  and  $\mathcal{S}(\cdot)$  is an arbitrary real-valued function. The  
 139 length of a walk-based proximity is the maximum length of all the walks of the proximity.

140 Typical examples of walk-based proximities include the shortest distance (You et al., 2019), the high-  
 141 order proximities (a sum of walks weighted by their lengths) (Zhang et al., 2018), and random walk  
 142 probabilities (Klicpera et al., 2019). Next, we give a definition of preserving walk-based proximities.

143 **Definition 4.** For a given walk-based proximity, a GNN is said to be able to preserve the proximity  
 144 if there exists a decoder function  $\mathcal{F}_{de}(\cdot)$  satisfying that for any graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$ , there exist  
 145 parameters  $\mathbf{W}_G$  so that  $\forall \epsilon > 0$ :

$$\left| \mathbf{S}_{i,j} - \mathcal{F}_{de}(\mathbf{H}_{i,:}^{(L)}, \mathbf{H}_{j,:}^{(L)}) \right| < \epsilon, \quad (6)$$

146 where

$$\mathbf{H}^{(L)} = \mathcal{F}_{GNN}(\mathbf{A}, \mathbf{F}; \mathbf{W}_G). \quad (7)$$

147 Note that we do not constrain the GNN architecture as long as it follows Eq. (1), and the decoder  
 148 function is also arbitrary (but notice that it cannot take the graph structure as inputs). In fact, both  
 149 the GNN and the decoder function can be arbitrarily deep and with sufficient hidden units.

150 **Theorem 1.** The existing permutation-equivariant GNNs cannot preserve any walk-based proximity  
 151 except the trivial solution that all node pairs have the same proximity.<sup>3</sup>

152 The formulation and proof of the theorem are given in Appendix A.1. Since walk-based proximities  
 153 are rather general and widely adopted in graph analytical tasks such as link prediction, the theorem  
 154 shows that the existing permutation-equivariant GNNs cannot handle these tasks well.

## 155 4 THE MODEL

### 156 4.1 A GNN FRAMEWORK USING STOCHASTIC MESSAGE PASSING

157 A major shortcoming of permutation-equivariant GNNs is that they cannot differentiate automorphic  
 158 node pairs. To solve that problem, we need to introduce some mechanism as “symmetry breaking”,  
 159 i.e., to enable GNNs to distinguish these nodes. To achieve this goal, we sample a stochastic matrix  
 160  $\mathbf{E} \in \mathbb{R}^{N \times d}$  where each element follows an i.i.d. normal distribution  $\mathcal{N}(0, 1)$ . The stochastic matrix  
 161 can provide signals in distinguishing the nodes because they are randomly sampled without being  
 162 affected by the graph automorphism. In fact, we can easily calculate that the Euclidean distance  
 163 between two stochastic signals divided by a constant  $\sqrt{2}$  follows a chi distribution  $\chi_d$ :

$$\frac{1}{\sqrt{2}} |\mathbf{E}_{i,:} - \mathbf{E}_{j,:}| \sim \chi_d, \forall i, j. \quad (8)$$

164 When  $d$  is reasonably large, e.g.,  $d > 20$ , the probability of two signals being close is very low.  
 165 Then, inspired by the message-passing framework, we apply a GNN on the stochastic matrix so that  
 166 nodes can exchange information of the stochastic signals:

$$\tilde{\mathbf{E}} = \mathcal{F}_{GNN}(\mathbf{A}, \mathbf{E}; \mathbf{W}). \quad (9)$$

167 We call  $\tilde{\mathbf{E}}$  the stochastic representation of nodes. Using the stochastic matrix and message-passing,  
 168  $\tilde{\mathbf{E}}$  can be used to preserve node proximities (see Theorem 2 and Theorem 3). Then, to let our model  
 169 still be able to utilize node features, we concatenate  $\tilde{\mathbf{E}}$  with the node representations from another  
 170 GNN with node features as inputs:

$$\begin{aligned} \mathbf{H} &= \mathcal{F}_{\text{output}}([\tilde{\mathbf{E}}, \mathbf{H}^{(L)}]) \\ \tilde{\mathbf{E}} &= \mathcal{F}_{GNN}(\mathbf{A}, \mathbf{E}; \mathbf{W}), \mathbf{H}^{(L)} = \mathcal{F}_{GNN'}(\mathbf{A}, \mathbf{F}; \mathbf{W}'), \end{aligned} \quad (10)$$

<sup>3</sup>Proposition 1 in (You et al., 2019) can be regarded as a special case of Theorem 1 using the shortest distance proximity.

171 where  $\mathcal{F}_{\text{output}}(\cdot)$  is an aggregation function such as a linear function or simply the identity mapping.  
 172 In a nutshell, our proposed method augments the existing GNNs with a stochastic representation  
 173 learned by message-passings to differentiate different nodes and preserve node proximities.

174 There is also a delicate choice worthy mentioning, i.e., whether the stochastic matrix  $\mathbf{E}$  is fixed or  
 175 resampled in each epoch. By fixing  $\mathbf{E}$ , the model can learn to memorize the stochastic representation  
 176 and distinguish different nodes, but with the cost of unable to handle nodes not seen during training.  
 177 On the other hand, by resampling  $\mathbf{E}$  in each epoch, the model can have a better generalization  
 178 ability since the model cannot simply remember one specific stochastic matrix. However, the node  
 179 representations are not fixed (but pairwise proximities are preserved; see Theorem 2). In these cases,  
 180  $\tilde{\mathbf{E}}$  is more capable of handling pairwise tasks such as link prediction or pairwise node classification.  
 181 In this paper, we use a fixed  $\mathbf{E}$  for transductive datasets and resample  $\mathbf{E}$  for inductive datasets.

182 **Time Complexity** From Eq.(10), the time complexity of our framework mainly depends on the  
 183 two GNNs in learning the stochastic and permutation-equivariant node representations. In this paper,  
 184 we instantiate these two GNNs using simple message-passing GNNs such as GCN (Kipf & Welling,  
 185 2017) and SGC (Wu et al., 2019) (see Section 4.2 and Section 4.3). Thus, the time complexity of  
 186 our method is the same as these models, which is  $O(M)$ , i.e., linear with respect to the number of  
 187 edges. We also empirically compare the running time of different models in Appendix 5.5. Besides,  
 188 many acceleration schemes for GNNs such as sampling (Chen et al., 2018a;b; Huang et al., 2018)  
 189 or partitioning the graph (Chiang et al., 2019) can be directly applied to our framework.

## 190 4.2 A LINEAR INSTANTIATION

191 Based on the general framework shown in Eq. (10), we attempt to explore its minimum model  
 192 instantiation, i.e., a linear model. Specifically, inspired by Simplified Graph Convolution (SGC) (Wu  
 193 et al., 2019), we adopt a linear message-passing for both GNNs, i.e.,

$$\mathbf{H} = \mathcal{F}_{\text{output}}([\tilde{\mathbf{E}}, \mathbf{H}^{(L)}]) = \mathcal{F}_{\text{output}}([\tilde{\mathbf{A}}^K \mathbf{E}, \tilde{\mathbf{A}}^K \mathbf{F}]), \quad (11)$$

194 where  $\tilde{\mathbf{A}} = (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}} (\mathbf{A} + \mathbf{I}) (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}$  is the normalized graph adjacency matrix with self-loops  
 195 proposed in GCN (Kipf & Welling, 2017) and  $K$  is the number of propagation steps. We also set  
 196  $\mathcal{F}_{\text{output}}(\cdot)$  in Eq. (11) as a linear mapping or identity mapping.

197 Though seemingly simple, we show that such an SMP instantiation possesses a theoretical guarantee  
 198 in preserving the walk-based proximities.

199 **Theorem 2.** *An SMP in Eq. (11) with the message-passing matrix  $\tilde{\mathbf{A}}$  and the number of propagation  
 200 steps  $K$  can preserve the walk-based proximity  $\tilde{\mathbf{A}}^K (\tilde{\mathbf{A}}^K)^T$  with high probability if the dimensional-  
 201 ity of the stochastic matrix  $d$  is sufficiently large, where the superscript  $T$  denotes matrix transpose.  
 202 The theorem is regardless of whether  $\mathbf{E}$  are fixed or resampled.*

203 The mathematical formulation and proof of the theorem are given in Appendix A.2. In addition, we  
 204 show that SMP is equivalent to a permutation-equivariant GNN with certain parametrization.

205 **Remark 1.** *Suppose we adopt  $\mathcal{F}_{\text{output}}(\cdot)$  as a linear function with the output dimensionality the  
 206 same as  $\mathcal{F}_{GNN'}$ . Then, Eq. (10) is equivalent to the permutation-equivariant  $\mathcal{F}_{GNN'}(\mathbf{A}, \mathbf{F}; \mathbf{W}')$  if the  
 207 parameters in  $\mathcal{F}_{\text{output}}(\cdot)$  are all-zeros for  $\tilde{\mathbf{E}}$  and an identity matrix for  $\mathbf{H}^{(L)}$ .*

208 The result is straightforward from the definition. Then, we have the following corollary.

209 **Corollary 2.** *For any task, Eq. (10) with the aforementioned linear  $\mathcal{F}_{\text{output}}(\cdot)$  is at least as powerful  
 210 as the permutation-equivariant  $\mathcal{F}_{GNN'}(\mathbf{A}, \mathbf{F}; \mathbf{W}')$ , i.e., the minimum training loss of using  $\mathbf{H}$  in  
 211 Eq. (10) is equal to or smaller than using  $\mathbf{H}^{(L)} = \mathcal{F}_{GNN'}(\mathbf{A}, \mathbf{F}; \mathbf{W}')$ .*

212 In other words, SMP will not hinder the performance<sup>4</sup> even the tasks are permutation-equivariant  
 213 since the stochastic representations are concatenated with the permutation-equivariant GNNs fol-  
 214 lowed by a linear mapping. In these cases, the linear SMP is equivalent to SGC (Wu et al., 2019).

215 Combining Theorem 2 and Corollary 2, the linear SMP instantiation in Eq. (11) is capable of han-  
 216 dling both proximity-aware and permutation-equivariant tasks.

<sup>4</sup>Similar to previous works such as (Hamilton et al., 2017; Xu et al., 2018a), we only consider the minimum training loss because the optimization landscapes and generalization gaps are difficult to analyze analytically.

## 217 4.3 NON-LINEAR EXTENSIONS

218 One may question whether a more sophisticated variant of Eq. (10) can further improve the expres-  
 219 siveness of SMP. There are three adjustable components in Eq. (10): two GNNs in propagating the  
 220 stochastic matrix and node features, respectively, and an output function. In theory, adopting non-  
 221 linear models as either component is able to enhance the expressiveness of SMP. Indeed, if we use  
 222 a sufficiently expressive GNN in learning  $\tilde{\mathbf{E}}$  instead of linear propagations, we can prove a more  
 223 general version of Theorem 2 as follows.

224 **Theorem 3.** *An SMP variant following Eq.(10) with  $\mathcal{F}_{GNN}(\mathbf{A}, \mathbf{E}; \mathbf{W})$  containing  $L$  layers can  
 225 preserve any length- $L$  walk-based proximity if the message-passing and updating functions in the  
 226 GNN are sufficiently expressive. In this theorem, we also assume the Gaussian random vectors  $\mathbf{E}$   
 227 are rounded to machine precision so that  $\mathbf{E}$  is drawn from a countable subspace of  $\mathbb{R}$ .*

228 The proof of the theorem is given in Appendix A.3. Similarly, we can adopt more advanced methods  
 229 for  $\mathcal{F}_{\text{output}}(\cdot)$  such as gating or attention so that the two GNNs are more properly integrated.

230 Although non-linear extensions of SMP can, in theory, increase the model expressiveness, they also  
 231 take a higher risk of over-fitting due to model complexity, not to mention that the computational cost  
 232 will also increase. In practice, we find in ablation studies that the linear SMP instantiation in Eq. (11)  
 233 works reasonably well on most of the datasets (please refer to Section 5.4 for further details).

## 234 5 EXPERIMENTS

## 235 5.1 EXPERIMENTAL SETUPS

236 **Datasets** We conduct experiments on the following **ten** datasets: two simulation datasets, **Grid**  
 237 and **Communities** (You et al., 2019), a communication dataset **Email** (You et al., 2019), two coau-  
 238 thor networks, **CS** and **Physics** (Shchur et al., 2018), two protein interaction networks, **PPI** (Hamil-  
 239 ton et al., 2017) and **PPA** (Hu et al., 2020), and three GNN benchmarks, **Cora**, **CiteSeer**, and  
 240 **PubMed** (Yang et al., 2016). We only report the results of three benchmarks for the node classi-  
 241 fication task and the results for other tasks are shown in Appendix B due to the page limit. More  
 242 details of the datasets including their statistics are provided in Appendix C.1. These datasets cover  
 243 a wide spectrum of domains, sizes, and with or without node features. Since Email and PPI contain  
 244 more than one graph, we conduct experiments in an *inductive setting* on these two datasets, i.e., the  
 245 training, validation, and testing set are split with respect to different graphs.

246 **Baselines** We adopt two sets of baselines. The first set is permutation-equivariant GNNs including  
 247 GCN (Kipf & Welling, 2017), GAT (Velickovic et al., 2018), and SGC (Wu et al., 2019), which are  
 248 widely adopted GNN architectures. The second set contains P-GNN (You et al., 2019), the only  
 249 proximity-aware GNN to date. We use the P-GNN-F version.

250 In comparing with the baselines, we mainly evaluate two variants of SMP with different  $\mathcal{F}_{\text{output}}(\cdot)$ :  
 251 SMP-Identity, i.e.,  $\mathcal{F}_{\text{output}}(\cdot)$  as an identity mapping, and SMP-Linear, i.e.,  $\mathcal{F}_{\text{output}}(\cdot)$  as a linear  
 252 mapping. Note that both variants adopt linear message-passing functions as SGC. We conduct more  
 253 ablation studies with different SMP variants in Section 5.4.

254 For fair comparisons, we adopt the same architecture and hyper-parameters for all the methods  
 255 (please refer to Appendix C.2 for the details). For datasets without node features, we adopt a con-  
 256 stant vector as the node features. We experiment on two tasks: link prediction and node classifica-  
 257 tion. Additional experiments on graph reconstruction, pairwise node classification, and running time  
 258 comparison are provided in Appendix B. We repeat the experiments 10 times for datasets except for  
 259 PPA and 3 times for PPA, and report the average results.

## 260 5.2 LINK PREDICTION

261 Link prediction aims to predict missing links of a graph. Specifically, we split the edges into 80%-  
 262 10%-10% and use them for training, validation, and testing, respectively. Besides adopting those real  
 263 edges as positive samples, we obtain negative samples by randomly sampling an equal number of  
 264 node pairs that do not have edges. For all the methods, we set a simple classifier:  $\text{Sigmoid}(\mathbf{H}_i^T \mathbf{H}_j)$ ,  
 265 i.e., use the inner product to predict whether a node pair  $(v_i, v_j)$  forms a link, and use AUC (area

Table 2: The results of link prediction tasks measured in AUC (%). The best results and the second-best results for each dataset, respectively, are in bold and underlined.

Model	Grid	Communities	Email	CS	Physics	PPI
SGC	57.6±3.8	51.9±1.6	68.5±7.0	96.5±0.1	<b>96.6±0.1</b>	80.5±0.4
GCN	61.8±3.6	50.3±2.5	67.4±6.9	93.4±0.3	93.8±0.2	78.0±0.4
GAT	61.0±5.5	51.1±1.6	53.5±6.3	93.7±0.9	94.1±0.4	79.3±0.5
PGNN <sup>5</sup>	<u>73.4±6.0</u>	<u>97.8±0.6</u>	70.9±6.4	82.2±0.5	Out of memory	80.8±0.4
SMP-Identity	55.1±4.8	<b>98.0±0.7</b>	<u>72.9±5.1</u>	<u>96.5±0.1</u>	<u>96.5±0.1</u>	<u>81.0±0.2</u>
SMP-Linear	<b>73.6±6.2</b>	97.7±0.5	<b>75.7±5.0</b>	<b>96.7±0.1</b>	96.1±0.1	<b>81.9±0.3</b>

266 under the curve) as the evaluation metric. One exception to the aforementioned setting is that on the  
 267 PPA dataset, we follow the splits and evaluation metric (i.e., Hits@100) provided by the dataset (Hu  
 268 et al., 2020). The results except PPA are shown in Table 2. We make the following observations.

- 269 • Our proposed SMP achieves the best results on five out of the six datasets and is highly compet-  
 270 itive (the second-best result) on the other (Physics). The results demonstrate the effectiveness of  
 271 our proposed method on link prediction tasks. We attribute the strong performance of SMP to its  
 272 capability of maintaining both proximity-awareness and permutation-equivariance properties.
- 273 • On Grid, Communities, Email, and PPI, both SMP and P-GNN outperform the permutation-  
 274 equivariant GNNs, proving the importance of preserving node proximities. Although SMP is  
 275 simpler and more computationally efficient than P-GNN, SMP reports even better results.
- 276 • When node features are available (CS, Physics, and PPI), SGC can outperform GCN and GAT.  
 277 The results re-validate the experiments in SGC (Wu et al., 2019) that the non-linearity in GNNs is  
 278 not necessarily indispensable. Some plausible reasons include that the additional model complex-  
 279 ity brought by non-linear operators makes the models tend to overfit and also difficult to train (see  
 280 Appendix B.6). On those datasets, SMP retains comparable performance on two coauthor graphs  
 281 and shows better performance on PPI, possibly because node features on protein graphs are less  
 282 informative than node features on coauthor graphs for predicting links, and thus preserving graph  
 283 structure is more beneficial on PPI.
- 284 • As Email and PPI are conducted in an inductive setting, i.e., using different graphs for train-  
 285 ing/validation/testing, the results show that SMP can handle inductive tasks as well.

286 The results on PPA are shown in Table 1. SMP  
 287 again outperforms all the baselines, showing that  
 288 it can handle large-scale graphs with millions of  
 289 nodes and edges. PPA is part of a recently re-  
 290 leased Open Graph Benchmark (Hu et al., 2020).  
 291 The superior performance on PPA further demon-  
 292 strates the effectiveness of our proposed method  
 293 in the link prediction task.

### 294 5.3 NODE CLASSIFICATION

295 Next, we conduct experiments of node classifica-  
 296 tion, i.e., predicting the labels of nodes. Since  
 297 we need ground-truths in the evaluation, we only  
 298 adopt datasets with node labels. Specifically, for  
 299 CS and Physics, following (Shchur et al., 2018), we adopt 20/30 labeled nodes per class for train-  
 300 ing/validation and the rest for testing. For Communities, we adjust the number as 5/5/10 labeled  
 301 nodes per class for training/validation/testing. For Cora, CiteSeer, and PubMed, we use the default  
 302 splits that came with the datasets. We do not adopt Email because some graphs in the dataset are too  
 303 small to show stable results and exclude PPI as it is a multi-label dataset.

Table 1: The results of link prediction on the PPA dataset. The best result and the second-best result are in bold and underlined, respectively.

Model	Hits@100
SGC	0.1187±0.0012
GCN	0.1867±0.0132
GraphSAGE	0.1655±0.0240
P-GNN	Out of Memory
Node2vec	0.2226±0.0083
Matrix Factorization	<u>0.3229±0.0094</u>
SMP-Identity	0.2018±0.0148
SMP-Linear	<b><u>0.3582±0.0070</u></b>

<sup>5</sup>The results of PGNN are slightly different compared to the paper because we adopt a more practical and common setting that negative samples in the data are not known apriori but randomly sampled in each epoch.

304 We use a softmax layer on the learned node representations as the classifier and adopt accuracy,  
 305 i.e., how many percentages of nodes are correctly classified, as the evaluation criteria. We omit the  
 306 results of SMP-Identity for this task since the node representations in SMP-Identity have a fixed  
 307 dimensionality that does not match the number of classes.

Table 3: The results of node classification tasks measured by accuracy (%). The best results and the second-best results for each dataset, respectively, are in bold and underlined.

Model	Communities	CS	Physics	Cora	CiteSeer	PubMed
SGC	7.1±2.1	67.2±12.8	92.3±1.6	76.9±0.2	63.6±0.0	74.2±0.1
GCN	<u>7.5±1.2</u>	91.1±0.7	93.1±0.8	81.4±0.5	<b>71.3±0.5</b>	<b>79.3±0.4</b>
GAT	5.0±0.0	<u>90.5±0.5</u>	<b>93.1±0.4</b>	<b>82.9±0.5</b>	<u>71.2±0.6</u>	<u>77.9±0.5</u>
PGNN	5.2±0.5	77.6±7.6	Out of memory	59.2±1.5	55.7±0.9	Out of memory
SMP-Linear	<b>99.9±0.3</b>	<b>91.5±0.8</b>	<u>93.1±0.8</u>	80.9±0.8	68.2±1.0	76.5±0.8

308 The results are shown in Table 3. From the table, we observe that SMP reports nearly perfect results  
 309 on Communities. Since the node labels are generated by graph structures on Communities and there  
 310 are no node features, the model needs to be proximity-aware to handle it well. P-GNN, which shows  
 311 promising results in the link prediction task, also fails miserably here.

312 On the other five graphs, SMP reports highly competitive performance. These graphs are commonly-  
 313 used benchmarks for GNNs. P-GNN, which completely ignores permutation-equivariance, performs  
 314 poorly as expected. In contrast, SMP can manage to recover the permutation-equivariant GNNs  
 315 and avoid being misled, as proven in Remark 1. In fact, SMP even shows better results than its  
 316 counterpart, SGC, indicating that preserving proximities is also helpful for these datasets.

#### 317 5.4 ABLATION STUDIES

318 We conduct ablation studies by comparing different SMP variants, including SMP-Identity, SMP-  
 319 Linear, and the additional three variants as follows:

- 320 • SMP-MLP: we set  $\mathcal{F}_{\text{output}}(\cdot)$  as a fully-connected network with 1 hidden layer.
- 321 • SMP-Linear-GCN<sub>feat</sub>: we set  $\mathcal{F}_{\text{GNN}'}(\mathbf{A}, \mathbf{F}; \mathbf{W}')$  in Eq. (10) to be a GCN (Kipf & Welling,  
 322 2017), i.e., induce non-linearity in message passing for features.  $\mathcal{F}_{\text{output}}(\cdot)$  is still linear.
- 323 • SMP-Linear-GCN<sub>both</sub>: we set both  $\mathcal{F}_{\text{GNN}}(\mathbf{A}, \mathbf{E}; \mathbf{W})$  and  $\mathcal{F}_{\text{GNN}'}(\mathbf{A}, \mathbf{F}; \mathbf{W}')$  to be a GCN (Kipf  
 324 & Welling, 2017), i.e., induce non-linearity in message passing for both features and stochastic  
 325 representations.  $\mathcal{F}_{\text{output}}(\cdot)$  is linear.

326 We show the results for link prediction tasks in Table 4. The results for node classification and  
 327 pairwise node classification, which imply similar conclusions, are provided in Table 10 and Table 11  
 328 in Appendix B.5. We make the following observations.

- 329 • In general, SMP-Linear shows good-enough performance, achieving the best or second-best re-  
 330 sults on six datasets and highly competitive on the other (Communities). SMP-Identity, which  
 331 does not have parameters in the output function, performs slightly worse. The results demon-  
 332 strate the importance of adopting a learnable linear layer in the output function, which is con-  
 333 sistent with Remark 1. SMP-MLP does not lift the performance in general, showing that adding  
 334 extra complexities in  $\mathcal{F}_{\text{output}}(\cdot)$  brings no gain in those datasets.
- 335 • SMP-Linear-GCN<sub>feat</sub> reports the best results on Communities, PPI, and PPA, indicating that  
 336 adding extra non-linearities in propagating node features are helpful for some graphs.
- 337 • SMP-Linear-GCN<sub>both</sub> reports the best results on Grid with a considerable margin. Recall that  
 338 Grid has no node features. The results indicate that inducing non-linearities can help the stochas-  
 339 tic representations capture more proximities, which is more helpful for featureless graphs.

#### 340 5.5 EFFICIENCY COMPARISON

341 To compare the efficiency of different methods quantitatively, we report the running time of different  
 342 methods in Table 5. The results are averaged over 3,000 epochs on an NVIDIA TESLA M40 GPU

Table 4: The ablation study of different SMP variants for the link prediction task. Datasets except PPA are measured by AUC (%) and PPA is measured by Hits@100. The best results and the second-best results for each dataset are in bold and underlined, respectively.

Model	Grid	Communities	Email	CS	Physics	PPI	PPA
SMP-Identity	55.1±4.8	<b>98.0±0.7</b>	72.9±5.1	<u>96.5±0.1</u>	<b>96.5±0.1</b>	81.0±0.2	0.2018±0.0148
SMP-Linear	<u>73.6±6.2</u>	97.7±0.5	<b>75.7±5.0</b>	<b>96.7±0.1</b>	<u>96.1±0.1</u>	81.9±0.3	0.3582±0.0070
SMP-MLP	72.1±4.3	<u>97.8±0.6</u>	62.7±8.1	88.9±0.8	89.2±0.4	80.1±0.3	0.2035±0.0038
SMP-Linear-GCN <sub>feat</sub>	72.8±4.2	<b>98.0±0.4</b>	<u>74.2±3.9</u>	92.9±0.6	94.3±0.2	<b>82.3±1.0</b>	<b>0.4090±0.0087</b>
SMP-Linear-GCN <sub>both</sub>	<b>80.5±3.9</b>	97.3±0.7	<u>73.4±5.5</u>	89.8±2.0	91.7±0.2	79.7±0.3	0.2125±0.0232

Table 5: The average running time (in milliseconds) for each epoch (including both training and testing), on link prediction task.

Model	Grid	Communities	Email	CS	Physics	PPI
SGC	25	28	58	210	651	704
GCN	25	35	75	214	612	784
GAT	36	43	140	258	801	919
PGNN	81	84	206	19,340	Out of Memory	6,521
SMP-Identity	26	37	96	284	751	840
SMP-Linear	28	26	84	212	616	832
SMP-MLP	23	28	83	237	614	831
SMP-Linear-GCN <sub>feat</sub>	23	29	90	231	636	855
SMP-Linear-GCN <sub>both</sub>	34	40	95	228	626	895

343 with 12 GB of memory. The results show that SMP is computationally efficient, i.e., only marginally  
 344 slower than SGC and comparable to GCN. P-GNN is at least an order of magnitude slower except  
 345 for the extremely small graphs such as Grid, Communities, or Email with no more than a thousand  
 346 nodes. In addition, the expensive memory cost makes P-GNN unable to work on large-scale graphs.

## 347 5.6 MORE EXPERIMENTAL RESULTS

348 Besides the aforementioned experiments, we also conduct experiments on the following tasks: graph  
 349 reconstruction (Appendix B.1), pairwise node classification (Appendix B.2), and comparing with  
 350 one-hot IDs (Appendix B.3). Please refer to the Appendix for experimental results and correspond-  
 351 ing analyses.

## 352 6 CONCLUSION

353 In this paper, we propose SMP, a general and simple GNN to maintain both proximity-awareness  
 354 and permutation-equivariance properties. We propose to augment the existing GNNs with stochastic  
 355 node representations learned to preserve node proximities. We prove that SMP can enable GNN to  
 356 preserve node proximities in theory and is equivalent to a permutation-equivariant GNN with certain  
 357 parametrization. Experimental results demonstrate the effectiveness and efficiency of SMP. Ablation  
 358 studies show that a linear SMP instantiation works reasonably well on most of the datasets.

## 359 REFERENCES

- 360 Dana Angluin. Local and global properties in networks of processors. In *Proceedings of the twelfth*  
 361 *annual ACM symposium on Theory of computing*, pp. 82–93, 1980.
- 362 Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi,  
 363 Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al.  
 364 Relational inductive biases, deep learning, and graph networks. *arXiv:1806.01261*, 2018.
- 365 Stephen P Borgatti. Centrality and network flow. *Social networks*, 27(1):55–71, 2005.

- 366 Karsten M Borgwardt and Hans-Peter Kriegel. Shortest-path kernels on graphs. In *Fifth IEEE*  
367 *international conference on data mining (ICDM'05)*, pp. 8–pp. IEEE, 2005.
- 368 Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann Lecun. Spectral networks and locally  
369 connected networks on graphs. In *International Conference on Learning Representations*, 2014.
- 370 Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with  
371 variance reduction. In *International Conference on Machine Learning*, pp. 942–950, 2018a.
- 372 Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks via  
373 importance sampling. In *International Conference on Learning Representations*, 2018b.
- 374 Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An  
375 efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of*  
376 *the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp.  
377 257–266, 2019.
- 378 Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal  
379 neighbourhood aggregation for graph nets. *arXiv preprint arXiv:2004.05718*, 2020.
- 380 George Dasoulas, Ludovic Dos Santos, Kevin Scaman, and Aladin Virmaux. Coloring graph neural  
381 networks for node disambiguation. *arXiv preprint arXiv:1912.06058*, 2019.
- 382 Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on  
383 graphs with fast localized spectral filtering. In *Advances in neural information processing systems*,  
384 pp. 3844–3852, 2016.
- 385 Matthias Fey, Jan E Lenssen, Christopher Morris, Jonathan Masci, and Nils M Kriege. Deep graph  
386 matching consensus. In *International Conference on Learning Representations*, 2020.
- 387 Thomas Gärtner, Peter Flach, and Stefan Wrobel. On graph kernels: Hardness results and efficient  
388 alternatives. In *Learning theory and kernel machines*, pp. 129–143. Springer, 2003.
- 389 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural  
390 message passing for quantum chemistry. In *International Conference on Machine Learning*, pp.  
391 1263–1272, 2017.
- 392 Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains.  
393 In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, volume 2,  
394 pp. 729–734. IEEE, 2005.
- 395 Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings*  
396 *of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*,  
397 pp. 855–864, 2016.
- 398 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.  
399 In *Advances in neural information processing systems*, pp. 1024–1034, 2017.
- 400 Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, YongDong Zhang, and Meng Wang. Lightgcn:  
401 Simplifying and powering graph convolution network for recommendation. In *Proceedings of*  
402 *the 43rd International ACM SIGIR Conference on Research and Development in Information*  
403 *Retrieval*, pp. 639–648, 2020.
- 404 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,  
405 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. *arXiv*  
406 *preprint arXiv:2005.00687*, 2020.
- 407 Wenbing Huang, Tong Zhang, Yu Rong, and Junzhou Huang. Adaptive sampling towards fast graph  
408 representation learning. In *Advances in neural information processing systems*, 2018.
- 409 Nicolas Keriven and Gabriel Peyré. Universal invariant and equivariant graph neural networks. In  
410 *Advances in Neural Information Processing Systems*, pp. 7090–7099, 2019.

- 411 Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational  
412 inference for interacting systems. In *International Conference on Machine Learning*, pp. 2688–  
413 2697, 2018.
- 414 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional net-  
415 works. In *Proceedings of the 6th International Conference on Learning Representations*, 2017.
- 416 Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. Predict then propagate:  
417 Graph neural networks meet personalized pagerank. In *International Conference on Learning  
418 Representations*, 2019.
- 419 Ioannis Konstas, Vassilios Stathopoulos, and Joemon M Jose. On social networks and collaborative  
420 recommendation. In *Proceedings of the 32nd international ACM SIGIR conference on Research  
421 and development in information retrieval*, pp. 195–202, 2009.
- 422 Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for  
423 semi-supervised learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- 424 Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. Gated graph sequence neural  
425 networks. In *International Conference on Learning Representations*, 2016.
- 426 Nathan Linial. Locality in distributed graph algorithms. *SIAM Journal on computing*, 21(1):193–  
427 201, 1992.
- 428 Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *International Con-  
429 ference on Learning Representations*, 2020.
- 430 Jianxin Ma, Chang Zhou, Peng Cui, Hongxia Yang, and Wenwu Zhu. Learning disentangled repre-  
431 sentations for recommendation. In *Advances in Neural Information Processing Systems 32*, pp.  
432 5712–5723. 2019.
- 433 Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph  
434 networks. In *Advances in Neural Information Processing Systems*, pp. 2156–2167, 2019a.
- 435 Haggai Maron, Heli Ben-Hamu, Nadav Shamir, and Yaron Lipman. Invariant and equivariant graph  
436 networks. In *International Conference on Learning Representations*, 2019b.
- 437 Alessio Micheli. Neural network for graphs: A contextual constructive approach. *IEEE Transactions  
438 on Neural Networks*, 20(3):498–511, 2009.
- 439 Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M  
440 Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In  
441 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5115–  
442 5124, 2017.
- 443 Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav  
444 Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks.  
445 In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019.
- 446 Ryan Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Relational pooling  
447 for graph representations. In *International Conference on Machine Learning*, 2019.
- 448 Moni Naor and Larry Stockmeyer. What can be computed locally? *SIAM Journal on Computing*,  
449 24(6):1259–1277, 1995.
- 450 Mark Newman. *Networks*. Oxford university press, 2018.
- 451 Antonio Ortega, Pascal Frossard, Jelena Kovačević, José MF Moura, and Pierre Vandergheynst.  
452 Graph signal processing: Overview, challenges, and applications. *Proceedings of the IEEE*, 106  
453 (5):808–828, 2018.
- 454 Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social repre-  
455 sentations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge  
456 discovery and data mining*, pp. 701–710, 2014.

- 457 Ryoma Sato, Makoto Yamada, and Hisashi Kashima. Random features strengthen graph neural  
458 networks. *arXiv preprint arXiv:2002.03155*, 2020.
- 459 Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.  
460 The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008.
- 461 Kakade Sham and Shakhnarovich Greg. Random projections. CMSC 35900 (Spring 2009) Large  
462 Scale Learning, 2020. URL [https://ttic.uchicago.edu/~gregory/courses/  
463 LargeScaleLearning/lectures/jl.pdf](https://ttic.uchicago.edu/~gregory/courses/LargeScaleLearning/lectures/jl.pdf). [Online; accessed 4-September-2020].
- 464 Oleksandr Shchur, Maximilian Mumme, Aleksandar Bojchevski, and Stephan Günnemann. Pitfalls  
465 of graph neural network evaluation. *Relational Representation Learning Workshop, NeurIPS  
466 2018*, 2018.
- 467 Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M Borg-  
468 wardt. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(Sep):2539–  
469 2561, 2011.
- 470 David I Shuman, Sunil K Narang, Pascal Frossard, Antonio Ortega, and Pierre Vandergheynst. The  
471 emerging field of signal processing on graphs: Extending high-dimensional data analysis to net-  
472 works and other irregular domains. *IEEE signal processing magazine*, 30(3):83–98, 2013.
- 473 Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale  
474 information network embedding. In *Proceedings of the 24th international conference on world  
475 wide web*, pp. 1067–1077, 2015.
- 476 Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua  
477 Bengio. Graph attention networks. In *Proceedings of the 7th International Conference on Learn-  
478 ing Representations*, 2018.
- 479 Santosh S Vempala. *The random projection method*, volume 65. American Mathematical Soc.,  
480 2005.
- 481 Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. In *Proceedings of  
482 the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pp.  
483 1225–1234, 2016.
- 484 Duncan J Watts. Networks, dynamics, and the small-world phenomenon. *American Journal of  
485 sociology*, 105(2):493–527, 1999.
- 486 Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Sim-  
487 plifying graph convolutional networks. In *International Conference on Machine Learning*, pp.  
488 6861–6871, 2019.
- 489 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
490 networks? In *International Conference on Learning Representations*, 2018a.
- 491 Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie  
492 Jegelka. Representation learning on graphs with jumping knowledge networks. In *International  
493 Conference on Machine Learning*, pp. 5453–5462, 2018b.
- 494 Zhilin Yang, William W Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning  
495 with graph embeddings. In *Proceedings of the 33rd International Conference on International  
496 Conference on Machine Learning-Volume 48*, pp. 40–48, 2016.
- 497 Jiaxuan You, Rex Ying, and Jure Leskovec. Position-aware graph neural networks. In *International  
498 Conference on Machine Learning*, pp. 7134–7143, 2019.
- 499 Ziwei Zhang, Peng Cui, Xiao Wang, Jian Pei, Xuanrong Yao, and Wenwu Zhu. Arbitrary-order  
500 proximity preserved network embedding. In *Proceedings of the 24th ACM SIGKDD International  
501 Conference on Knowledge Discovery & Data Mining*, pp. 2778–2786, 2018.
- 502 Marinka Zitnik and Jure Leskovec. Predicting multicellular function through multi-layer tissue  
503 networks. *Bioinformatics*, 33(14):i190–i198, 2017.

## 504 A THEOREMS AND PROOFS

## 505 A.1 THEOREM 1

506 Here we formulate and prove Theorem 1.

507 **Theorem 1.** *For any walk-based proximity function  $\mathcal{S}(\cdot)$ , a permutation-equivariant GNN cannot*  
 508 *preserve  $\mathcal{S}(\cdot)$ , except the trivial solution that all node pairs have the same proximity, i.e.,  $\mathbf{S}_{i,j} =$*   
 509  *$c, \forall i, j$ , where  $c$  is a constant.*

510 *Proof.* We prove the theorem by contradiction. Assume there exists a non-trivial  $\mathcal{S}(\cdot)$  which a  
 511 permutation-equivariant GNN can preserve. Consider any graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$  and denote  $N =$   
 512  $|\mathcal{V}|$ . We can create  $G' = (\mathcal{V}', \mathcal{E}', \mathbf{F}')$  with  $|\mathcal{V}'| = 2N$  so that:

$$\mathcal{E}'_{i,j} = \begin{cases} \mathcal{E}_{i,j} & \text{if } i \leq N, j \leq N \\ \mathcal{E}_{i-N, j-N} & \text{if } i > N, j > N \\ 0 & \text{else} \end{cases}, \quad \mathbf{F}'_{i,:} = \begin{cases} \mathbf{F}_{i,:} & \text{if } i \leq N \\ \mathbf{F}_{i-N,:} & \text{if } i > N \end{cases}. \quad (12)$$

513 Basically, we generate two “copies” of the original graph, one indexing from 1 to  $N$ , and the other  
 514 indexing from  $N + 1$  to  $2N$ . By assumption, there exists a permutation-equivariant GNN which can  
 515 preserve  $\mathcal{S}(\cdot)$  in  $G'$  and we denote the node representations as  $\mathbf{H}'^{(L)} = \mathcal{F}_{\text{GNN}}(\mathbf{A}', \mathbf{F}'; \mathbf{W}_{G'})$ . It is  
 516 easy to see that node  $v'_i$  and  $v'_{i+N}$  in  $G'$  form an automorphic node pair. Using Corollary 1, their  
 517 representations will be identical in any permutation-equivariant GNN, i.e.,

$$\mathbf{H}'_{i,:}^{(L)} = \mathbf{H}'_{i+N,:}^{(L)}, \forall i \leq N. \quad (13)$$

518 Also, note that there exists no walk from the two copies, i.e.  $\{v'_i \rightsquigarrow v'_j\} = \{v'_j \rightsquigarrow v'_i\} = \emptyset, \forall i \leq$   
 519  $N, j > N$ . As a result, for  $\forall i \leq N, j \leq N, \forall \epsilon > 0$ , we have:

$$\begin{aligned} |\mathbf{S}_{i,j} - \mathcal{S}(\emptyset)| &\leq \left| \mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}^{(L)}, \mathbf{H}'_{j,:}^{(L)}) \right| + \left| \mathcal{S}(\emptyset) - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}^{(L)}, \mathbf{H}'_{j,:}^{(L)}) \right| \\ &= \left| \mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}^{(L)}, \mathbf{H}'_{j,:}^{(L)}) \right| + \left| \mathbf{S}_{i,j+N} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}^{(L)}, \mathbf{H}'_{j+N,:}^{(L)}) \right| < 2\epsilon. \end{aligned} \quad (14)$$

520 We can prove the same for  $\forall i > N, j > N$ . The equation naturally holds if  $i \leq N, j > N$  or  
 521  $i > N, j \leq N$  since  $\{v'_i \rightsquigarrow v'_j\} = \emptyset$ . Combining the results, we have  $\forall \epsilon > 0, \forall i, j, |\mathbf{S}_{i,j} - \mathcal{S}(\emptyset)| <$   
 522  $2\epsilon$ . Since  $\epsilon$  can be arbitrarily small, the equation shows that all node pairs have the same proximity  
 523  $c = \mathcal{S}(\emptyset)$ , which leads to a contraction and finishes our proof.  $\square$

524 Notice that in our proof,  $G'$  can be constructed for any graph, so rather than designing one specific  
 525 counter-example, we have shown that there always exists an infinite number of counter-examples by  
 526 constructing automorphisms in the graph.

527 Some may find that our counter-examples in the above proof will lead to multiple connected com-  
 528 ponents. Next, we give an alternative proof maintaining one connected component (assuming the  
 529 original graph is connected) under the assumption that the walk-based proximity is of finite length.

530 *Proof.* Similar to the previous proof, we assume there exists a non-trivial  $\mathcal{S}(\cdot)$  which a permutation-  
 531 equivariant GNN can preserve. Besides, we assume the length of  $\mathcal{S}(\cdot)$  is upper bounded by  $l_{\max}$ ,  
 532 where  $l_{\max}$  is any finite number, i.e.,  $\forall i, j$ ,

$$\mathbf{S}_{i,j} = \mathcal{S}(\{v_i \rightsquigarrow v_j\}) = \mathcal{S}(\{v_i \rightsquigarrow v_j | \text{len}(v_i \rightsquigarrow v_j) \leq l_{\max}\}). \quad (15)$$

533 Then, for a connected graph  $G = (\mathcal{V}, \mathcal{E}, \mathbf{F})$ , we create  $G' = (\mathcal{V}', \mathcal{E}', \mathbf{F}')$  similar to Eq. (12). Specif-  
534 ically, denoting  $\tilde{N} = N + l_{\max}$ , we let  $G'$  have  $3\tilde{N}$  nodes so that:

$$\mathcal{E}'_{i,j} = \begin{cases} \mathcal{E}_{i,j} & \text{if } i, j \leq N \\ 1 & \text{if } N \leq i, j \leq \tilde{N} + 1, |j - i| = 1 \\ \mathcal{E}_{i-\tilde{N}, j-\tilde{N}} & \text{if } \tilde{N} < i, j \leq \tilde{N} + N \\ 1 & \text{if } \tilde{N} + N \leq i, j \leq 2\tilde{N} + 1, |j - i| = 1 \\ \mathcal{E}_{i-2\tilde{N}, j-2\tilde{N}} & \text{if } 2\tilde{N} < i, j \leq 2\tilde{N} + N \\ 1 & \text{if } 2\tilde{N} + N \leq i, j, |j - i| = 1 \\ 1 & \text{if } i = 3\tilde{N}, j = 1 \text{ or } j = 3\tilde{N}, i = 1 \\ 0 & \text{else} \end{cases}, \mathbf{F}'_{i,:} = \begin{cases} \mathbf{F}_{i,:} & \text{if } i \leq N \\ 0 & \text{if } N < i \leq \tilde{N} \\ \mathbf{F}_{i-\tilde{N},:} & \text{if } \tilde{N} < i \leq \tilde{N} + N \\ 0 & \text{if } \tilde{N} + N < i \leq 2\tilde{N} \\ \mathbf{F}_{i-2\tilde{N},:} & \text{if } 2\tilde{N} < i \leq 2\tilde{N} + N \\ 0 & \text{if } 2\tilde{N} + N < i \end{cases} \quad (16)$$

535 Intuitively, we create three ‘‘copies’’ of  $G$  and three ‘‘bridges’’ to connect the copies and thus make  
536  $G'$  also connected. It is also easy to see that nodes  $v'_i, v'_{i+\tilde{N}}$ , and  $v'_{i+2\tilde{N}}$  all form automorphic node  
537 pairs and thus we have:

$$\mathbf{H}'_{i,:} = \mathbf{H}'_{i+\tilde{N},:} = \mathbf{H}'_{i+2\tilde{N},:}, \forall i \leq \tilde{N}. \quad (17)$$

538 Next, we can see that the nodes in  $G'$  are divided into six parts (three copies and three bridges),  
539 which we denote as  $\mathcal{V}'_1 = \{v_1, \dots, v_N\}$ ,  $\mathcal{V}'_2 = \{v_{N+1}, \dots, v_{\tilde{N}}\}$ ,  $\mathcal{V}'_3 = \{v_{\tilde{N}+1}, \dots, v_{\tilde{N}+N}\}$ ,  $\mathcal{V}'_4 =$   
540  $\{v_{\tilde{N}+N+1}, \dots, v_{2\tilde{N}}\}$ ,  $\mathcal{V}'_5 = \{v_{2\tilde{N}+1}, \dots, v_{2\tilde{N}+N}\}$ , and  $\mathcal{V}'_6 = \{v_{2\tilde{N}+N+1}, \dots, v_{3\tilde{N}}\}$ . Since  $\mathcal{V}'_2, \mathcal{V}'_4, \mathcal{V}'_6$   
541 are bridges with length  $l_{\max}$ , any walk crosses these bridges will have a length large than  $l_{\max}$ . For  
542 example, let us focus on  $v_i \in \mathcal{V}'_1$ , i.e.,  $i \leq N$ . If  $v_j$  is in  $\mathcal{V}'_3, \mathcal{V}'_4$ , or  $\mathcal{V}'_5$  (i.e.,  $\tilde{N} < j \leq 2\tilde{N} + N$ ), any  
543 walk  $v_i \rightsquigarrow v_j$  will either pass the bridge  $\mathcal{V}'_2$  or  $\mathcal{V}'_6$  and thus has a length larger than  $l_{\max}$ . As a result,  
544 we have:

$$\mathbf{S}_{i,j} = \mathcal{S}(\{v_i \rightsquigarrow v_j\}) = \mathcal{S}(\{v_i \rightsquigarrow v_j | \text{len}(v_i \rightsquigarrow v_j) \leq l_{\max}\}) = \mathcal{S}(\emptyset). \quad (18)$$

545 If  $v_j \in \mathcal{V}'_1$  or  $v_j \in \mathcal{V}'_2$ , i.e.,  $j \leq \tilde{N}$ , we can use the fact that  $v_j$  and  $v_{j+\tilde{N}}$  forms an automorphic node  
546 pair similar to Eq. (14), i.e.,  $\forall \epsilon > 0$ , we have

$$\begin{aligned} |\mathbf{S}_{i,j} - \mathcal{S}(\emptyset)| &\leq \left| \mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}, \mathbf{H}'_{j,:}) \right| + \left| \mathcal{S}(\emptyset) - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}, \mathbf{H}'_{j,:}) \right| \\ &= \left| \mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}, \mathbf{H}'_{j,:}) \right| + \left| \mathbf{S}_{i,j+\tilde{N}} - \mathcal{F}_{\text{de}}(\mathbf{H}'_{i,:}, \mathbf{H}'_{j+\tilde{N},:}) \right| < 2\epsilon. \end{aligned} \quad (19)$$

547 Similarly, if  $v_j \in \mathcal{V}'_6$ , i.e.,  $2\tilde{N} + N < j$ , we can use the fact that  $v_j$  and  $v_{j-\tilde{N}}$  forms an automorphic  
548 node pair to prove the same inequality. Thus, we prove that if  $i \leq N, \forall \epsilon > 0, \forall j, |\mathbf{S}_{i,j} - \mathcal{S}(\emptyset)| < 2\epsilon$ .  
549 The same proof strategy can be applied to  $i > N$ . Since  $\epsilon$  can be arbitrarily small, the results show  
550 that all node pairs have the same proximity  $\mathcal{S}(\emptyset)$ , which leads to a contraction and finishes our  
551 proof.  $\square$

## 552 A.2 THEOREM 2

553 Here we formulate and prove Theorem 2. Note that some notations and definitions are introduced in  
554 Appendix A.1.

555 **Theorem 2.** *For the walk-based proximity  $\mathbf{S} = \tilde{\mathbf{A}}^K (\tilde{\mathbf{A}}^K)^T$ , SMP can preserve the proximity with  
556 high probability if the dimensionality of the stochastic matrix is sufficiently large, i.e.,  $\forall \epsilon > 0, \forall \delta >$   
557  $0$ , there  $\exists d_0$  so that any  $d > d_0$ :*

$$P(|\mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}_{i,:}, \mathbf{H}_{j,:})| < \epsilon) > 1 - \delta, \quad (20)$$

558 where  $\mathbf{H}$  are the node representation obtained from SMP in Eq. (11). The result holds for any  
559 stochastic matrix and thus is regardless of whether  $\mathbf{E}$  is fixed or resampled during each epoch.

560 *Proof.* Our proof is mostly based on the standard random projection theory. Firstly, since we have  
561 proven in Theorem 1 that the permutation-equivariant representations cannot preserve any walk-  
562 based proximity, here we prove that we can preserve the proximity only using  $\tilde{\mathbf{E}}$ , which can be

563 easily achieved by ignoring  $\mathbf{H}^{(L)}$  in  $\mathcal{F}_{\text{output}}([\tilde{\mathbf{E}}, \mathbf{H}^{(L)}])$ , e.g., if we set  $\mathcal{F}_{\text{output}}$  as a linear function, the  
 564 model can learn to set the corresponding weights for  $\mathbf{H}^{(L)}$  as all-zeros.

565 We set the decoder function as a normalized inner product:

$$\mathcal{F}_{\text{de}}(\mathbf{H}_{i,:}, \mathbf{H}_{j,:}) = \frac{1}{d} \mathbf{H}_{i,:} \mathbf{H}_{j,:}^T. \quad (21)$$

566 Then, denoting  $\mathbf{a}_i = \tilde{\mathbf{A}}_{i,:}^K$  and recalling  $\tilde{\mathbf{E}} = \tilde{\mathbf{A}}^K \mathbf{E}$ , we have:

$$|\mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}_{i,:}, \mathbf{H}_{j,:})| = |\mathbf{a}_i \mathbf{a}_j^T - \frac{1}{d} \tilde{\mathbf{E}}_{i,:} \tilde{\mathbf{E}}_{j,:}^T| = |\mathbf{a}_i \mathbf{a}_j^T - \mathbf{a}_i \frac{1}{d} \mathbf{E} \mathbf{E}^T \mathbf{a}_j^T|. \quad (22)$$

567 Since  $\mathbf{E}$  is a Gaussian random matrix, from the Johnson-Lindenstrauss lemma (Vempala, 2005) (in  
 568 the inner product preservation forum, e.g., see Corollary 2.1 and its proof in (Sham & Greg, 2020)),  
 569  $\forall 0 < \epsilon' < \frac{1}{2}$ , we have:

$$P\left(|\mathbf{a}_i \mathbf{a}_j^T - \mathbf{a}_i \frac{1}{d} \mathbf{E} \mathbf{E}^T \mathbf{a}_j^T| \leq \frac{\epsilon'}{2} (\|\mathbf{a}_i\| + \|\mathbf{a}_j\|)\right) > 1 - 4e^{-\frac{(\epsilon'^2 - \epsilon'^3)d}{4}}. \quad (23)$$

570 By setting  $\epsilon' = \frac{\epsilon}{\max_i \|\mathbf{a}_i\|}$ , we have  $\epsilon > \frac{\epsilon'}{2} (\|\mathbf{a}_i\| + \|\mathbf{a}_j\|)$  and:

$$P(|\mathbf{S}_{i,j} - \mathcal{F}_{\text{de}}(\mathbf{H}_{i,:}, \mathbf{H}_{j,:})| < \epsilon) > 1 - 4e^{-\frac{(\frac{\epsilon}{\max_i \|\mathbf{a}_i\|})^2 - \frac{\epsilon}{\max_i \|\mathbf{a}_i\|})^3 d}{4}}, \quad (24)$$

571 which leads to the theorem by solving and setting  $d_0$  as follows:

$$4e^{-\frac{(\frac{\epsilon}{\max_i \|\mathbf{a}_i\|})^2 - \frac{\epsilon}{\max_i \|\mathbf{a}_i\|})^3 d_0}{4}} = \delta \Rightarrow d_0 = \frac{4 \log \frac{4}{\delta} (\max_i \|\mathbf{a}_i\|)^3}{\epsilon^2 \max_i \|\mathbf{a}_i\| - \epsilon^3}. \quad (25)$$

572 □

### 573 A.3 THEOREM 3

574 Here we formulate and prove Theorem 3. Note that some notations and definitions are introduced in  
 575 Appendix A.1.

**Theorem 3.** *For any length- $L$  walk-based proximity, i.e.,*

$$\mathbf{S}_{i,j} = \mathcal{S}(\{v_i \rightsquigarrow v_j\}) = \mathcal{S}(\{v_i \rightsquigarrow v_j | \text{len}(v_i \rightsquigarrow v_j) \leq L\}),$$

576 *where  $\text{len}(\cdot)$  is the length of a walk, there exists an SMP variant in Eq. (10) with  $\mathcal{F}_{\text{GNN}}(\mathbf{A}, \mathbf{E}; \mathbf{W})$*   
 577 *containing  $L$  layers (including the input layer) to preserve that proximity if the following conditions*  
 578 *hold: (1) The stochastic matrix  $\mathbf{E}$  contains unique signals for different nodes, i.e.  $\mathbf{E}_{i,:} \neq \mathbf{E}_{j,:}, \forall i \neq$*   
 579  *$j$ . (2) The message-passing and updating functions in learning  $\tilde{\mathbf{E}}$  are bijective. (3) The decoder*  
 580 *function  $\mathcal{F}_{\text{de}}(\cdot)$  also takes  $\mathbf{E}$  as inputs and is universal approximation.*

581 *Proof.* Similar as Theorem 2, we only utilize  $\tilde{\mathbf{E}}$  during our proof. We use  $\mathbf{e}_i^{(l)}, 0 \leq l < L$  to  
 582 denote the node representations in the  $l^{\text{th}}$  layer of  $\mathcal{F}_{\text{GNN}}(\mathbf{A}, \mathbf{E}; \mathbf{W})$ , i.e.,  $\mathbf{e}_i^{(0)} = \mathbf{E}_{i,:}$  and  $\mathbf{e}_i^{(L-1)} =$   
 583  $\tilde{\mathbf{E}}_{i,:}$ . Our proof strategy is to show that the stochastic node representations can remember all the  
 584 information about the walks.

585 Firstly, as the message-passing and updating function are bijective by assumption, we can recover  
 586 from the node representations in each layer all their neighborhood representations in the previous  
 587 layer. Specifically, there exist  $\mathcal{F}^{(l)}(\cdot), 1 \leq l < L$  such that:

$$\mathcal{F}^{(l)}(\mathbf{e}_i^{(l)}) = \left[ \mathbf{e}_i^{(l-1)}, \left\{ \mathbf{e}_j^{(l-1)}, j \in \mathcal{N}_i \right\} \right]^6. \quad (26)$$

588 For notation conveniences, we split the function into two parts, one for the node itself and the other  
 589 for its neighbors:

$$\begin{aligned} \mathcal{F}_{\text{self}}^{(l)}(\mathbf{e}_i^{(l)}) &= \mathbf{e}_i^{(l-1)}, \\ \mathcal{F}_{\text{neighbor}}^{(l)}(\mathbf{e}_i^{(l)}) &= \left\{ \mathbf{e}_j^{(l-1)}, j \in \mathcal{N}_i \right\}. \end{aligned} \quad (27)$$

<sup>6</sup>To let  $\mathcal{F}^{(l)}(\cdot)$  output a set with arbitrary lengths, we can adopt sequence-based models such an LSTM.

590 For the first function, if we successively apply such functions from the  $l^{th}$  layer to the input layer, we  
 591 can recover the input features of the GNN, i.e.,  $\mathbf{E}$ . Since the stochastic matrix  $\mathbf{E}$  contains a unique  
 592 signal for different nodes, we can decode the node ID from  $\mathbf{e}_i^{(0)}$ , i.e., there exists  $\mathcal{F}_{\text{self}}^{(0)}(\mathbf{e}_i^{(0)}; \mathbf{E}) = i$ .  
 593 For brevity, we denote applying such  $l + 1$  functions to get the node ID as

$$\mathcal{F}_{\text{self}}^{(0:l)}(\mathbf{e}_i^{(l)}) = \mathcal{F}_{\text{self}}^{(0)}\left(\mathcal{F}_{\text{self}}^{(1)}\left(\dots\left(\mathcal{F}_{\text{self}}^{(l)}\left(\mathbf{e}_i^{(l)}\right)\right)\right)\right); \mathbf{E} = i. \quad (28)$$

594 For the second function, we can apply  $\mathcal{F}_{\text{neighbor}}^{(l-1)}$  to the decoded vector set so that we can recover their  
 595 neighborhood representations in the  $(l - 2)^{th}$  layer, etc.

596 Next, we show that for  $\mathbf{e}_j^{(l-1)}$ , there exists a length- $l$  walk  $v_i \rightsquigarrow v_j = (v_{a_1}, v_{a_2}, \dots, v_{a_l})$ , where  
 597  $v_{a_1} = v_i, v_{a_l} = v_j$  if and only if  $\mathcal{F}_{\text{self}}^{(0:l-1)}(\mathbf{e}_j^{(l-1)}) = a_l = j$  and there exists  $\mathbf{e}^{(l-2)}, \dots, \mathbf{e}^{(0)}$  such  
 598 that:

$$\begin{aligned} \mathbf{e}^{(l-2)} &\in \mathcal{F}_{\text{neighbor}}^{(l-1)}(\mathbf{e}_j^{(l-1)}), \mathcal{F}_{\text{self}}^{(0:l-2)}(\mathbf{e}^{(l-2)}) = a_{l-1}, \\ \mathbf{e}^{(l-3)} &\in \mathcal{F}_{\text{neighbor}}^{(l-2)}(\mathbf{e}^{(l-2)}), \mathcal{F}_{\text{self}}^{(0:l-3)}(\mathbf{e}^{(l-3)}) = a_{l-2}, \\ &\dots \\ \mathbf{e}^{(0)} &\in \mathcal{F}_{\text{neighbor}}^{(1)}(\mathbf{e}^{(1)}), \mathcal{F}_{\text{self}}^{(0:0)}(\mathbf{e}^{(0)}) = a_1 = i. \end{aligned} \quad (29)$$

599 This result is easily verified as:

$$\begin{aligned} (v_{a_1}, v_{a_2}, \dots, v_{a_l}) \text{ is a walk} &\Leftrightarrow \mathcal{E}_{a_i, a_j} = \mathcal{E}_{a_j, a_i} = 1 \Leftrightarrow a_i \in \mathcal{N}_{a_{i+1}}, \forall 1 \leq i < l \\ &\Leftrightarrow \exists \mathbf{e}^{(i-1)} \in \mathcal{F}_{\text{neighbor}}^{(i)}(\mathbf{e}^{(i)}), \mathcal{F}_{\text{self}}^{(0:i-1)}(\mathbf{e}^{(i-1)}) = a_i, \forall 1 \leq i < l. \end{aligned} \quad (30)$$

600 Note that all the information is encoded in  $\tilde{\mathbf{E}}$ , i.e., we can decode  $\{v_i \rightsquigarrow v_j | \text{len}(v_i \rightsquigarrow v_j) \leq L\}$  from  
 601  $\mathbf{e}_j^{(L-1)}$  by successively applying  $\mathcal{F}_{\text{self}}^{(l)}(\cdot), \mathcal{F}_{\text{neighbor}}^{(l)}(\cdot)$ . We can also apply  $\mathcal{F}_{\text{self}}^{(0:L-1)}$  to  $\mathbf{e}_i^{(L-1)}$  to get  
 602 the start node ID  $i$ . Putting it together, we have:

$$\mathcal{F}(\mathbf{e}_j^{(L-1)}, \mathbf{e}_i^{(L-1)}) = \{v_i \rightsquigarrow v_j | \text{len}(v_i \rightsquigarrow v_j) \leq L\}, \quad (31)$$

603 where  $\mathcal{F}(\cdot)$  is composed of  $\mathcal{F}_{\text{self}}^{(l)}(\cdot), 0 \leq l < L$  and  $\mathcal{F}_{\text{neighbor}}^{(l)}(\cdot), 1 \leq l < L$ . Applying the proximity  
 604 function  $\mathcal{S}(\cdot)$ , we have:

$$\mathcal{S}\left(\mathcal{F}(\mathbf{e}_j^{(L-1)}, \mathbf{e}_i^{(L-1)})\right) = \mathbf{S}_{i,j}. \quad (32)$$

605 We finish the proof by setting the real decoder function  $\mathcal{F}_{\text{de}}(\cdot)$  to arbitrarily approximate this desired  
 606 function  $\mathcal{S}(\mathcal{F}(\cdot, \cdot))$  under the universal approximation assumption.  $\square$

## 607 B ADDITIONAL EXPERIMENTAL RESULTS

### 608 B.1 GRAPH RECONSTRUCTION

609 To verify that our proposed SMP can indeed preserve node proximities, we conduct experiments of  
 610 graph reconstruction (Wang et al., 2016), i.e., using the node representations learned by GNNs to  
 611 reconstruct the edges of the graph. Graph reconstruction corresponds to the first-order proximity  
 612 between nodes, i.e., whether two nodes directly have a connection, which is the most straight-  
 613 forward node proximity (Tang et al., 2015). Specifically, following Section 5.2, we adopt the inner  
 614 product classifier  $\text{Sigmoid}(\mathbf{H}_i^T \mathbf{H}_j)$  and use AUC as the evaluation metric. To control the impact  
 615 of node features (i.e., since many graphs exhibit assortative mixing, even models only using node  
 616 features can reconstruct the edges to a certain extent), we do not use node features for all the models.

617 We report the results in Table 6. The results show that SMP greatly outperforms permutation-  
 618 equivariant GNNs such as GCN and GAT in graph reconstruction, clearly demonstrating that SMP  
 619 can better preserve node proximities. PGNN shows highly competitive results as SMP. However,  
 620 similar to other tasks, the intensive memory usage makes PGNN unable to handle medium-scale  
 621 graphs such as Physics and PubMed.

Table 6: The results of graph reconstruction measured in AUC (%). The best and the second-best results for each dataset, respectively, are in bold and underlined. OOM represents out of memory.

Model	Grid	Communities	Email	CS	Physics	PPI	Cora	CiteSeer	PubMed
SGC	74.8±0.4	65.4±1.6	71.6±0.3	66.7±0.1	66.2±0.0	76.3±0.2	56.7±9.7	58.5±0.1	71.9±0.1
GCN	73.0±0.3	63.7±1.2	72.5±0.4	75.5±0.4	76.8±0.4	79.2±0.4	68.2±3.9	69.8±8.0	77.2±2.1
GAT	59.6±1.2	52.9±1.1	56.9±1.9	57.0±1.4	59.1±0.7	61.1±1.9	57.8±1.0	63.2±1.5	58.8±0.8
PGNN	<b>99.4±0.1</b>	<u>97.7±0.1</u>	<u>85.6±0.8</u>	<b>97.2±0.6</b>	OOM	<u>85.2±0.6</u>	<b>98.1±0.6</b>	<b>99.7±0.1</b>	OOM
SMP-Identity	99.2±0.1	97.5±0.1	80.0±0.3	77.1±2.3	73.7±0.3	79.5±0.2	89.7±5.7	97.1±0.8	77.0±0.1
SMP-Linear	99.1±0.1	<b>97.8±0.1</b>	<b>86.7±0.2</b>	<u>96.3±0.2</u>	<b>95.5±0.2</b>	<b>85.5±0.1</b>	<u>96.3±0.1</u>	<u>98.2±0.1</u>	<b>95.8±0.2</b>

## 622 B.2 PAIRWISE NODE CLASSIFICATION

623 Besides standard node classification experiments reported in Section 5.3, we follow (You et al.,  
624 2019) and experiment on pairwise node classification, i.e., predicting whether two nodes have the  
625 same label. Compared with standard node classification, pairwise node classification focuses more  
626 on the relations between nodes and thus requires the model to be proximity-aware to perform well.

627 Similar to link prediction, we split the positive samples (i.e., node pairs with the same label) into an  
628 80%-10%-10% training-validation-testing set with an equal number of randomly sampled negative  
629 pairs. For large graphs, since the possible positive samples are intractable (i.e.  $O(N^2)$ ), we use a  
630 random subset. Since we also need node labels as the ground-truth, we only conduct pairwise node  
631 classification on datasets when node labels are available. We also exclude the results of PPI since  
632 the dataset is multi-label and cannot be used in a pairwise setting (You et al., 2019). Similar to  
633 Section 5.2, we adopt a simple inner product classifier and use AUC as the evaluation metric.

634 The results are shown in Table 7. We observe consistent results as link prediction in Section 5.2, i.e.,  
635 SMP reports the best results on four datasets and the second-best results on the other three datasets.  
636 These results again verify that SMP can effectively preserve and utilize node proximities when  
637 needed while retaining comparable performance when the tasks are more permutation-equivariant  
638 like, e.g., on CS and Physics.

Table 7: The results of pairwise node classification tasks measured in AUC (%). The best results and the second-best results for each dataset, respectively, are in bold and underlined.

Model	Communities	Email	CS	Physics	Cora	CiteSeer	PubMed
SGC	67.4±2.4	56.3±5.4	<b>99.8±0.0</b>	<u>99.6±0.0</u>	<u>99.2±0.3</u>	<b>95.5±0.7</b>	92.3±0.3
GCN	64.9±2.3	55.0±5.7	96.8±0.7	<b>99.7±0.1</b>	97.7±0.6	92.9±1.2	<b>94.8±0.4</b>
GAT	52.5±1.3	47.7±2.7	95.2±0.6	96.3±0.2	91.6±0.7	73.6±2.7	87.1±0.2
PGNN	<u>98.6±0.5</u>	<u>63.3±5.5</u>	90.0±0.5	Out of memory	85.5±1.2	49.8±1.8	Out of memory
SMP-Identity	<b>98.8±0.5</b>	56.9±4.1	99.7±0.0	<u>99.6±0.0</u>	<u>99.2±0.2</u>	95.2±1.1	91.9±0.3
SMP-Linear	<b>98.8±0.5</b>	<b>74.5±4.1</b>	<b>99.8±0.0</b>	<u>99.6±0.0</u>	<b>99.3±0.3</b>	<u>95.3±0.4</u>	<u>93.4±0.2</u>

## 639 B.3 COMPARISON WITH USING IDs

640 We further compare SMP with augmenting GNNs using a one-hot encoding of node IDs, i.e., the  
641 identity matrix. Intuitively, since the IDs of nodes are unique, such a method does not suffer from the  
642 automorphism problem and should also enable GNNs to preserve node proximities. However, theo-  
643 retically speaking, using such a one-hot encoding has two major problems. Firstly, the dimensionality  
644 of the identity matrix is  $N \times N$ , and thus the number of parameters in the first message-passing  
645 layer is also on the order of  $O(N)$ . Therefore, the method will inevitably be computationally ex-  
646 pensive and may not be scalable to large-scale graphs. A large number of parameters will also more  
647 likely lead to the overfitting problem. Secondly, the node IDs are not transferable across different  
648 graphs, i.e., the node  $v_1$  in one graph and the node  $v_1$  in another graph do not necessarily share a  
649 similar meaning. But as the parameters in the message-passings depend on the node IDs (since they  
650 are input features), such a mechanism cannot handle inductive tasks well.<sup>7</sup>

<sup>7</sup>One may question whether SMP is transferable across different graphs since the stochastic features are independently drawn. Empirically, we find that SMP reports reasonably well results on inductive datasets such

Table 8: The results of comparing SMP with using one-hot IDs in GCNs. OOM represents out of memory. — represents the task is unavailable.

Task	Model	Grid	Communities	Email	CS	Physics	PPI	Cora	CiteSeer	PubMed
Link Prediction	GCN <sub>onehot</sub>	91.5±2.1	98.3±0.7	71.2±3.5	93.1±1.3	OOM	78.6±0.3	86.8±1.5	81.7±1.1	89.4±0.5
	SMP-Linear	73.6±6.2	97.7±0.5	75.7±5.0	96.7±0.1	96.1±0.1	81.9±0.3	92.7±0.7	92.6±1.0	95.4±0.2
Pairwise Node Classification	GCN <sub>onehot</sub>	—	98.9±0.5	67.3±5.6	97.6±0.2	OOM	—	98.2±0.3	94.4±1.2	98.9±0.1
	SMP-Linear	—	98.8±0.5	74.5±4.1	99.8±0.0	99.6±0.0	—	99.3±0.3	95.3±0.4	93.4±0.2
Node Classification	GCN <sub>onehot</sub>	—	99.6±1.0	—	86.9±1.5	OOM	—	77.6±1.1	57.7±5.8	74.9±0.6
	SMP-Linear	—	99.9±0.3	—	91.5±0.8	93.1±0.8	—	80.9±0.8	68.2±1.0	76.5±0.8

651 We also empirically compare such a method with SMP and report the results in Table 8. The results  
652 show that SMP-Linear outperforms GCN<sub>onehot</sub> in most cases. Besides, GCN<sub>onehot</sub> fails to handle  
653 Physics, which is only a medium-scale graph, due to the heavy memory usage. One surprising result  
654 is that GCN<sub>onehot</sub> outperforms SMP-Linear on Grid, the simulated graph where nodes are placed on  
655 a  $20 \times 20$  grid. A plausible reason is that since the edges in Grid follow a specific rule, using a  
656 one-hot encoding gives GCN<sub>onehot</sub> enough flexibility to learn and remember the rules, and the model  
657 does not overfit because the graph has a rather small scale.

#### 658 B.4 ADDITIONAL LINK PREDICTION RESULTS

659 We further report the results of link prediction on three GNN benchmarks: Cora, CiteSeer, and  
660 PubMed. The results are shown in Table 9. The results show similar trends as other datasets pre-  
661 sented in Section 5.2, i.e., SMP reports comparable results as other permutation-equivariant GNNs  
662 while PGNN fails to handle the task well.

Table 9: The results of the link prediction task measured in AUC (%). The best results and the second-best results for each dataset, respectively, are in bold and underlined.

Model	Cora	CiteSeer	PubMed
SGC	<b>93.6±0.6</b>	<b>94.7±0.8</b>	<b>95.8±0.2</b>
GCN	90.6±1.0	78.2±1.7	92.4±0.9
GAT	88.5±1.2	87.8±1.0	89.2±0.8
PGNN	75.4±2.3	70.6±1.1	Out of memory
SMP-Identity	<u>93.0±0.6</u>	<u>92.9±0.5</u>	94.5±0.3
SMP-Linear	<u>92.7±0.7</u>	92.6±1.0	<u>95.4±0.2</u>
SMP-MLP	82.8±0.9	80.7±1.1	88.0±0.6
SMP-Linear-GCN <sub>feat</sub>	86.7±1.4	81.1±1.4	90.5±0.6
SMP-Linear-GCN <sub>both</sub>	80.1±2.5	80.0±2.0	81.1±2.0

#### 663 B.5 ADDITIONAL ABLATION STUDIES

664 We report the ablation study results for the node classification task and pairwise node classification  
665 task in Table 10 and Table 11, respectively. The results again show that SMP-Linear generally  
666 achieves good-enough results on the majority of the datasets and adding non-linearities does not  
667 necessarily lift the performance of SMP.

668 We also compare whether the stochastic signals  $\mathbf{E}$  are fixed or not during different training epochs  
669 for our proposed SMP. For brevity, we only report the results for the link prediction task in Table 12.  
670 The results show that fixing  $\mathbf{E}$  usually leads to better results on transductive datasets (recall that  
671 datasets except Email and PPI are transductive) and resampling  $\mathbf{E}$  leads to better results on inductive  
672 datasets in general. The results are consistent with our analysis in Section 4.1.

as Email and PPI. One plausible reason is that since the proximities of nodes are preserved even the random features per se are different (see Theorem 2), all subsequent parameters based on proximities can be transferred.

Table 10: The ablation study of different SMP variants for the node classification task. The best results and the second-best results are in bold and underlined, respectively.

Model	Communities	CS	Physics	Cora	CiteSeer	PubMed
SMP-Linear	99.9±0.3	<b>91.5±0.8</b>	<b>93.1±0.8</b>	<b>80.9±0.8</b>	<b>68.2±1.0</b>	76.5±0.8
SMP-MLP	<b>100.0±0.2</b>	90.1±0.5	92.3±0.8	79.3±0.8	67.0±1.5	76.8±0.9
SMP-Linear-GCN <sub>feat</sub>	<b>100.0±0.0</b>	89.8±0.7	<u>92.9±0.8</u>	78.9±1.2	<u>67.8±0.6</u>	<b>77.3±0.6</b>
SMP-Linear-GCN <sub>both</sub>	<b>100.0±0.2</b>	77.4±4.2	87.1±3.5	69.2±2.5	49.8±3.1	68.1±4.1

Table 11: The ablation study of different SMP variants for the pairwise node classification task. The best results and the second-best results are in bold and underlined, respectively.

Model	Communities	Email	CS	Physics	Cora	CiteSeer	PubMed
SMP-Identity	98.8±0.5	56.9±4.1	99.7±0.0	<b>99.6±0.0</b>	<u>99.2±0.2</u>	95.2±1.1	91.9±0.3
SMP-Linear	<u>98.8±0.5</u>	<b>74.5±4.1</b>	<b>99.8±0.0</b>	<b>99.6±0.0</b>	<b>99.3±0.3</b>	<b>95.3±0.4</b>	93.4±0.2
SMP-MLP	98.7±0.3	<u>65.4±6.3</u>	94.3±0.6	97.6±0.4	90.3±3.0	67.7±13.7	93.4±0.4
SMP-Linear-GCN <sub>feat</sub>	<b>99.0±0.4</b>	60.2±9.3	95.6±0.7	98.3±0.7	96.1±0.7	88.8±1.6	<b>94.8±0.2</b>
SMP-Linear-GCN <sub>both</sub>	<u>98.8±0.4</u>	61.6±6.0	95.2±0.7	97.8±0.8	94.3±1.9	83.5±3.9	<u>94.1±0.7</u>

## 673 B.6 COMPARISON OF PERMUTATION-EQUIVARIANT GNNs FOR LINK PREDICTION

674 To investigate the performance of linear and non-linear variants of permutation-equivariant GNNs  
675 for the link prediction task, we additionally report both the training accuracies and the testing accu-  
676 racies of SGC, GCN, and GAT in Table 13. Notice that to ensure a fair comparison, we do not adopt  
677 the early stopping strategy here so that different models have the same number of training epochs  
678 (otherwise, if a model tends to overfit, the early stopping strategy will terminate the training process  
679 when the number of training epochs is small and result in a spurious underfitting phenomena).

680 The results show that non-linear variants of GNNs (GCN and GAT) are more likely to overfit, i.e.,  
681 the margins between the training accuracies and the testing accuracies are usually larger, than the  
682 linear variant SGC. Besides, though possessing extra model expressiveness, non-linear GNNs are  
683 also difficult to train, i.e., the training accuracies of GCN and GAT are not necessarily higher than  
684 SGC. The results are consistent with the literature Wu et al. (2019); He et al. (2020).

## 685 C EXPERIMENTAL DETAILS FOR REPRODUCIBILITY

### 686 C.1 DATASETS

- 687 • **Grid** (You et al., 2019): A simulated 2D grid graph with size  $20 \times 20$  and no node feature.
- 688 • **Communities** (You et al., 2019): A simulated caveman graph (Watts, 1999) composed of 20  
689 communities with each community containing 20 nodes. The graph is perturbed by rewiring 1%  
690 edges randomly. It has no node feature and the label of each node indicates which community  
691 the node belongs to.
- 692 • **Email**<sup>8</sup> (You et al., 2019): Seven real-world email communication graphs. Each graph has six  
693 communities and each node has an integer label indicating the community the node belongs to.
- 694 • **Coauthor Networks**<sup>9</sup> (Shchur et al., 2018): Two networks from Microsoft academic graph in  
695 *CS* and *Physics* with their nodes representing authors and edges representing co-authorships  
696 between authors. The node features are embeddings of the paper keywords of the authors.
- 697 • **PPI**<sup>8</sup> (Hamilton et al., 2017): 24 protein-protein interaction networks. Each node has a 50-  
698 dimensional feature vector.
- 699 • **PPA**<sup>10</sup> (Hu et al., 2020): A network representing biological associations between proteins from  
700 58 different species. The node features are one-hot vectors of the species that the proteins are  
701 taken from.

<sup>8</sup><https://github.com/JiaxuanYou/P-GNN/tree/master/data>

<sup>9</sup><https://github.com/shchur/gnn-benchmark/tree/master/data/npz/>

<sup>10</sup><https://snap.stanford.edu/ogb/data/linkproppred/ppassoc.zip>

Table 12: The results of comparing whether the stochastic signals  $\mathbf{E}$  are fixed or not during different training epochs for the link prediction task. The better of the two results are in bold.

Model	$\mathbf{E}$	Grid	Communities	CS	Physics	Email	PPI
SMP-Identity	Fixed	55.1±4.8	<b>98.0±0.7</b>	<b>96.5±0.1</b>	96.5±0.1	<b>75.9±3.9</b>	80.4±0.4
	Not Fixed	<b>55.2±4.1</b>	97.6±0.7	96.4±0.1	96.5±0.1	72.9±5.1	<b>81.0±0.2</b>
SMP-Linear	Fixed	<b>73.6±6.2</b>	<b>97.7±0.5</b>	<b>96.7±0.1</b>	96.1±0.1	71.3±3.9	71.5±0.7
	Not Fixed	64.4±2.9	97.4±0.1	96.2±0.1	96.1±0.1	<b>75.7±5.0</b>	<b>81.9±0.3</b>

Table 13: The results of SGC, GCN, and GAT for the link prediction task. Both the training accuracies and the testing accuracies are reported.

Method	Results	Grid	Communities	Email	CS	Physics	PPI	Cora	CiteSeer	PubMed
SGC	Train	52.4±0.5	50.4±0.5	68.4±1.0	98.3±0.0	97.7±0.0	84.9±0.4	99.5±0.0	99.9±0.0	98.9±0.0
	Test	51.6±2.9	49.2±1.6	67.0±9.3	96.5±0.1	96.5±0.1	78.8±0.7	92.5±0.7	94.0±0.5	95.6±0.3
GCN	Train	51.3±1.0	50.0±0.6	68.9±4.0	95.3±0.3	95.1±0.3	76.1±0.6	96.5±1.0	77.3±1.3	93.9±1.6
	Test	54.6±4.3	49.2±1.1	66.0±3.1	93.2±0.3	93.8±0.3	74.9±0.6	89.2±0.3	76.1±2.5	90.6±1.2
GAT	Train	47.5±1.4	49.6±0.3	52.2±3.6	95.7±0.9	96.4±0.5	82.3±0.3	97.2±0.5	99.3±0.1	98.3±0.1
	Test	50.8±6.0	50.8±2.1	47.6±4.4	91.0±1.1	93.9±0.3	78.5±0.4	75.7±1.7	81.0±0.7	83.7±1.3

702 • **Cora, CiteSeer, PubMed**<sup>11</sup> (Yang et al., 2016): Three citation graphs where nodes correspond  
703 to papers and edges correspond to citations between papers. The node features are bag-of-words  
704 and the node labels are the ground truth topics of the papers.

705 We summarize the statistics of datasets in Table 14.

## 706 C.2 HYPER-PARAMETERS

707 We use the following hyper-parameters:

- 708 • All datasets except PPA: we uniformly set the number of layers for all the methods as 2, i.e., 2  
709 message-passing steps, and set the dimensionality of hidden layers as 32, i.e.,  $\mathbf{H}^{(l)} \in \mathbb{R}^{N \times 32}$ ,  
710 for all  $1 \leq l \leq L$  (for GAT, we use 4 heads with each head containing 8 units). We use Adam  
711 optimizer with an initial learning rate of 0.01 and decay the learning rate by 0.1 at epoch 200.  
712 The weight decay is  $5e-4$ . We train the model for 1,000 epochs and evaluate the model every 5  
713 epochs. We adopt an early-stopping strategy by reporting the testing performance at the epoch  
714 which achieves the best validation performance. For SMP, the dimensionality of the stochastic  
715 matrix is  $d = 32$ . For P-GNN, we use the P-GNN-F version, which uses the truncated 2-hop  
716 shortest path distance instead of the exact shortest distance.
- 717 • PPA: as suggested in the original paper (Hu et al., 2020), we set the number of GNN layers  
718 as 3 with each layer containing 256 hidden units and add a three-layer MLP after taking the  
719 Hadamard product between pair-wise node embeddings as the predictor, i.e.,  $\text{MLP}(\mathbf{H}_i \odot \mathbf{H}_j)$ .  
720 We use Adam optimizer with an initial learning rate of 0.01. We set the number of epochs for  
721 training as 40, evaluate the results on validation sets every epoch, and report the testing results  
722 using the model with the best validation performance. We also found that the dataset had issues  
723 with exploding gradients and adopt a gradient clipping strategy by limiting the maximum p2-  
724 norm of gradients as 1.0. The dimensionality of the stochastic matrix in SMP is  $d = 64$ .

## 725 C.3 HARDWARE AND SOFTWARE CONFIGURATIONS

726 All experiments are conducted on a server with the following configurations.

- 727 • Operating System: Ubuntu 18.04.1 LTS
- 728 • CPU: Intel(R) Xeon(R) CPU E5-2699 v4 @ 2.20GHz
- 729 • GPU: NVIDIA TESLA M40 with 12 GB of memory

<sup>11</sup><https://github.com/kimiyoung/planetoid/tree/master/data>

Table 14: The statistics of the datasets. For datasets with more than one graph, #Nodes and #Edges are summed over all the graphs and the experiments are conducted in an inductive setting.

Dataset	#Graphs	#Nodes	#Edges	#Features	#Classes
Grid	1	400	760	-	-
Communities	1	400	3,800	-	20
Email	7	1,005	25,571	-	42
CS	1	18,333	81,894	6,805	15
Physics	1	34,493	247,962	8,415	5
PPI	24	56,944	818,716	50	-
PPA	1	576,289	30,326,273	58	-
Cora	1	2,708	5,429	1,433	7
CiteSeer	1	3,327	4,732	3,703	6
PubMed	1	19,717	44,338	500	3

730 • Software: Python 3.6.8, PyTorch 1.4.0, PyTorch Geometric 1.4.3, NumPy 1.18.1, Cuda 10.1