

# REVISITING RANDOM WALKS FOR LEARNING ON GRAPHS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We revisit a recent model class for machine learning on graphs, where a random walk on a graph produces a machine-readable record, and this record is processed by a deep neural network to directly make vertex-level or graph-level predictions. We refer to these stochastic machines as **random walk neural networks (RWNNs)**, and through principled analysis, show that we can design them to be isomorphism invariant while capable of universal approximation of graph functions in probability. A useful finding is that almost any kind of record of random walk guarantees probabilistic invariance as long as the vertices are anonymized. This enables us, for example, to record random walks in plain text and adopt a language model to read these text records to solve graph tasks. We further establish a parallelism to message passing neural networks using tools from Markov chain theory, and show that over-smoothing in message passing is alleviated by construction in RWNNs, while over-squashing manifests as probabilistic under-reaching. We empirically demonstrate RWNNs on a range of problems, verifying our theoretical analysis and demonstrating the use of language models for separating strongly regular graphs where the 3-WL test fails, and transductive classification on arXiv citation network.

## 1 INTRODUCTION

Message-passing neural networks (MPNNs) are a popular class of neural networks on graphs where each vertex keeps a feature vector and updates it by propagating messages over neighbors (Battaglia et al., 2018). MPNNs have achieved success, one reason being their respect for the natural symmetries of graph learning problems, i.e., invariance to graph isomorphism (Chen et al., 2019). On the other hand, MPNNs in their basic form can be viewed as implementing color updates of the Weisfeiler-Lehman (1-WL) graph isomorphism test, and thus their expressive power is not stronger (Xu et al., 2019). Also, their inner working is often tied to the topology of the input graph, which is related to over-smoothing, over-squashing, and under-reaching of features under mixing (Giraldo et al., 2023).

In this work, we revisit an alternative direction for learning on graphs, where a random walk on a graph produces a machine-readable record, and this record is processed by a deep neural network that directly makes graph-level or vertex-level predictions. We refer to these stochastic machines as **random walk neural networks (RWNNs)**. Pioneering works were done in this direction, often motivated by the compatibility of random walks with sequence learning methods. Examples include DeepWalk (Perozzi et al., 2014) and node2vec (Grover & Leskovec, 2016) which use skip-gram on walks, AWE (Ivanov & Burnaev, 2018) which uses embeddings of anonymized walks (Micali & Zhu, 2016), and CRaWI (Tönshoff et al., 2023) which uses 1D CNNs on sliding window descriptions of walks. These methods were proven powerful in various contexts, such as graph isomorphism learning that require expressive powers surpassing 1-WL test (Tönshoff et al., 2023; Martinkus et al., 2023).

Despite the promises, unlike MPNNs, we currently have not reached a good principled understanding of RWNNs. First, in context of geometric DL, it is unclear how we can systematically incorporate the symmetries of graph learning problems into RWNNs, as their neural networks are not necessarily isomorphism invariant. Second, the upper bound of their expressive power, along with the requirements on each component to reach it, is not clearly known. Third, whether and how the over-smoothing and over-squashing problems may occur in RWNNs are not well understood. Addressing such questions is important, as they allow us to better understand the strengths and limits of the existing methods, and potentially design enhanced RWNNs by short-circuiting the search in their large design space.

Table 1: An overview of prior methods in context of RWNNs, and new components originating from our analysis. NB means non-backtracking; MDLR means minimum degree local rule (Section 2.1).

Method	Random walk	Recording function $q$	Reader NN $f_\theta$
DeepWalk	Uniform	Identity	Skip-gram
node2vec	Second-order $pq$ -walk	Identity	Skip-gram
AWE	Uniform	Anonymization	Embedding table
CRaWI	Uniform + NB	Sliding window	1D CNN
RW-AgentNet	Uniform	Neighborhoods	RNN
WalkLM	Uniform	Text attributes	Language model
Ours (Sections 2.1, 2.2, 3)	MDLR (+ NB) (2.1) (+ restarts) (2.2)	Anonymization (+ named neighborhoods) (2.1)	Any universal (3.1)

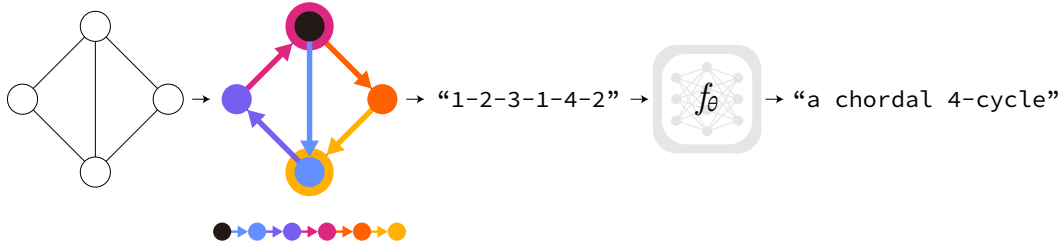


Figure 1: An RWNN that reads text record using a language model.

**Contributions** The main goal of our work is establishing a principled understanding of RWNNs with a focus on the aforementioned aspects. Our starting point is a new formalization of RWNNs which decouples random walks, their records, and neural networks that process them. This abstracts a wide range of design choices including prior methods, as in Table 1. The key idea of our analysis is then to understand RWNNs through a combination of two recent perspectives in geometric DL: probabilistic notions of invariance and expressive power (Bloem-Reddy & Teh, 2020; Abboud et al., 2021), and invariant projection of non-invariant neural networks (Puny et al., 2022; Dym et al., 2024).

This idea allows us to impose probabilistic invariance on RWNNs even if their neural networks lack symmetry, by instead requiring (probabilistic) invariance conditions on random walks and their records. This provides a justification for anonymized recording of walks (Micali & Zhu, 2016) and our novel extension of it using named neighborhoods. As long as probabilistic invariance holds, this also enables recording walks in plain text and adopting a language model to process them (Figure 1).

We then upper-bound the expressive power of RWNNs as universal approximation in probability which surpasses the WL hierarchy of MPNNs, under the condition that random walk records the graph of interest with a high probability. This establishes a useful link to cover times in Markov chains, providing guidance on designing random walks to minimize or bound the cover times. From this we motivate novel adaptation of minimum degree local rule (MDLR) walks (David & Feige, 2018) for RWNNs, and introduce restarts when working on a large and possibly infinite graph. We also provide a justification for the widespread use of non-backtracking (Tönshoff et al., 2023).

Continuing the link to Markov chain theory, we further analyze RWNNs and establish a parallelism to a linearized model of MPNNs. From this, we show that over-smoothing in MPNNs is inherently avoided in RWNNs, while over-squashing manifests as probabilistic under-reaching. This eliminates the typical trade-off between over-smoothing and over-squashing in MPNNs (Giraldo et al., 2023; Nguyen et al., 2023), and allows an RWNN to focus on overcoming under-reaching by scaling the walk length or using rapidly-mixing random walks such as non-backtracking walks.

We empirically demonstrate RWNNs on several graph problems. On synthetic setups, optionally with a small 1-layer transformer, we verify our claims on cover times, over-smoothing, and over-squashing. Then, we demonstrate adapting a language model (DeBERTa (He et al., 2021) to solve the challenging task of isomorphism learning of strongly regular graphs with perfect accuracy whereas the 3-WL test fails, improving over the previously known best result of RWNNs. We further suggest that our approach can turn transductive classification problem on a large graph into an in-context learning problem by recording labeled vertices during random walks. To demonstrate this, we apply Llama 3 (meta llama, 2024) model on transductive classification on arXiv citation network with 170k vertices, and show that it can outperform a range of MPNNs as well as zero- and few-shot baselines. Our experiments show the utility of our approach in analysis and development of RWNNs.

## 2 RANDOM WALK NEURAL NETWORKS

We start our discussion by formalizing random walk neural networks (RWNNs). We define an RWNN as a randomized function  $X_\theta(\cdot)$  that takes a graph  $G$  as input and outputs a random variable  $X_\theta(G)$  on the output space  $\mathbb{R}^d$ .<sup>1</sup> In case of vertex-level tasks, it is additionally queried with an input vertex  $v$  which gives a random variable  $X_\theta(G, v)$ . An RWNN consists of the following components:

1. **Random walk algorithm** that produces  $l$  steps of vertex transitions  $v_0 \rightarrow \dots \rightarrow v_l$  on the input graph  $G$ . If an input vertex  $v$  is given, we fix the starting vertex by  $v_0 = v$ .
2. **Recording function**  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$  that produces a machine-readable record  $\mathbf{z}$  of the random walk. It may access the graph  $G$  to record auxiliary information such as attributes.
3. **Reader neural network**  $f_\theta : \mathbf{z} \mapsto \hat{\mathbf{y}}$  that processes record  $\mathbf{z}$  and outputs a prediction  $\hat{\mathbf{y}}$  in  $\mathbb{R}^d$ . It is the only trainable component and is not restricted to specific architectures.

For graph-level prediction, we sample  $\hat{\mathbf{y}} \sim X_\theta(G)$  by running a random walk on  $G$  and processing its record using the reader NN  $f_\theta$ . For vertex-level prediction  $\hat{\mathbf{y}} \sim X_\theta(G, v)$ , we query the input vertex  $v$  by simply starting the walk from it. In practice, we ensemble several sampled predictions e.g. by averaging, which can be understood as Monte Carlo estimation of the mean predictor  $(\cdot) \mapsto \mathbb{E}[X_\theta(\cdot)]$ .

We now analyze each component for graph-level tasks on finite graphs, and then vertex-level tasks on possibly infinite graphs. The latter simulates the problem of scaling to large graphs such as in transductive classification. As our definition captures many known designs (Table 1), we discuss them as well. We leave pseudocode in Appendix A.2 and leave notations and proofs in Appendix A.5.

### 2.1 GRAPH-LEVEL TASKS

Let  $\mathbb{G}$  be the class of undirected, connected, and simple graphs<sup>2</sup>. Let  $n \geq 1$  and  $\mathbb{G}_n$  be the collection of graphs in  $\mathbb{G}$  with at most  $n$  vertices. Our goal is to model a graph-level function  $\phi : \mathbb{G}_n \rightarrow \mathbb{R}^d$  using an RWNN  $X_\theta(\cdot)$ . Since  $\phi$  is a graph function, it is reasonable to assume isomorphism invariance:

$$\phi(G) = \phi(H), \quad \forall G \simeq H, \quad (1)$$

Incorporating the invariance structure to our model class  $X_\theta(\cdot)$  would offer generalization benefit. As  $X_\theta(\cdot)$  is randomized, we accept the probabilistic notion of invariance (Bloem-Reddy & Teh, 2020):

$$X_\theta(G) \stackrel{d}{=} X_\theta(H), \quad \forall G \simeq H. \quad (2)$$

A justification is that, if  $X_\theta(\cdot)$  is invariant in probability, its mean predictor  $(\cdot) \mapsto \mathbb{E}[X_\theta(\cdot)]$  would be an invariant function. For more in-depth discussion, please see Section 4. We now claim that we can achieve probabilistic invariance of  $X_\theta(\cdot)$  by properly choosing the random walk algorithm and recording function while not imposing any constraint on the reader NN  $f_\theta$ .

**Proposition 2.1.**  *$X_\theta(\cdot)$  is invariant in probability, if its random walk algorithm is invariant in probability and its recording function is invariant.*

In other words, even if  $f_\theta$  lacks symmetry, invariant random walk and recording function provably converts it into an invariant random variable  $X_\theta(\cdot)$ . This is an extension of invariant projection operators on functions (Puny et al., 2022; Dym et al., 2024) to a more general and probabilistic setup.

**Random walk algorithm** A random walk algorithm is invariant in probability if it satisfies:

$$\pi(v_0) \rightarrow \dots \rightarrow \pi(v_l) \stackrel{d}{=} u_0 \rightarrow \dots \rightarrow u_l, \quad \forall G \stackrel{\pi}{\simeq} H, \quad (3)$$

where  $v_{[\cdot]}$  is a random walk on  $G$ ,  $u_{[\cdot]}$  is a random walk on  $H$ , and  $\pi : V(G) \rightarrow V(H)$  specifies the isomorphism from  $G$  to  $H$ . It turns out that many random walk algorithms in literature are already invariant. To see this, let us write the probability of walking from a vertex  $u$  to its neighbor  $x \in N(u)$ :

$$\text{Prob}[v_t = x | v_{t-1} = u] := \frac{c_G(u, x)}{\sum_{y \in N(u)} c_G(u, y)}, \quad (4)$$

where the function  $c_G : E(G) \rightarrow \mathbb{R}_+$  assigns positive weights (conductances) to edges. If we set  $c_G(\cdot) = 1$ , we recover uniform random walk used in DeepWalk (Perozzi et al., 2014). We show:

<sup>1</sup>While any output type such as text is possible (Figure 1), we explain with vector output for simplicity.

<sup>2</sup>We assume this for simplicity but extending to directed or attributed graphs is possible.

**Proposition 2.2.** *The random walk in Equation 4 is invariant in probability if its conductance  $c_{[\cdot]}(\cdot)$  is invariant:*

$$c_G(u, v) = c_H(\pi(u), \pi(v)), \quad \forall G \stackrel{\pi}{\simeq} H. \quad (5)$$

*It includes constant conductance, and any choice that only uses degrees of endpoints  $\deg(u)$ ,  $\deg(v)$ .*

We favor using degrees since it is cheap while potentially improving the behaviors of walks. Among many instantiations, we find that the following conductance function called the minimum degree local rule (MDLR) (Abdullah et al., 2015; David & Feige, 2018) is particularly useful:

$$c_G(u, v) := \frac{1}{\min[\deg(u), \deg(v)]}. \quad (6)$$

MDLR is special as it has  $O(n^2)$  vertex cover time, i.e., expected time of visiting all  $n$  vertices of a graph, optimal among first-order random walks (Equation 4) that use degrees of endpoints.

A common practice is to add non-backtracking property that enforces  $v_{t+1} \neq v_{t-1}$  (Tönshoff et al., 2023), and we also find this beneficial. In Appendix A.1 we extend Equation 3 and Proposition 2.2 to second-order walks that include non-backtracking and node2vec walks (Grover & Leskovec, 2016).

**Recording function** A recording function  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$  takes a random walk and produces a machine-readable record  $\mathbf{z}$ . We let  $q(\cdot, G)$  have access to the graph  $G$  the walk is taking place. This allows recording auxiliary information such as vertex or edge attributes. A recording function is invariant if it satisfies the following for any given random walk  $v_{[\cdot]}$  on  $G$ :

$$q(v_0 \rightarrow \dots \rightarrow v_l, G) = q(\pi(v_0) \rightarrow \dots \rightarrow \pi(v_l), H), \quad \forall G \stackrel{\pi}{\simeq} H. \quad (7)$$

Invariance requires that  $q(\cdot, G)$  produces the same record  $\mathbf{z}$  regardless of re-indexing of vertices of  $G$  into  $H$ . For this, we have to be careful in how we represent each vertex in a walk  $v_0 \rightarrow \dots \rightarrow v_l$  as a machine-readable value, and which auxiliary information we record from  $G$ . For example, identity recording used in DeepWalk (Perozzi et al., 2014) and node2vec (Grover & Leskovec, 2016) are not invariant as the result is sensitive to vertex re-indexing. We highlight two choices which are invariant:

- **Anonymization.** We name each vertex in a walk with a unique integer, starting from  $v_0 \mapsto 1$  and incrementing it based on their order of discovery. For instance, a random walk  $a \rightarrow b \rightarrow c \rightarrow a$  translates to a sequence  $1 \rightarrow 2 \rightarrow 3 \rightarrow 1$ .
- **Anonymization + named neighbors.** While applying anonymization for each vertex  $v$  in a walk, we record its neighbors  $u \in N(v)$  if they are already named but the edge  $v \rightarrow u$  has not been recorded yet. For instance, a walk  $a \rightarrow b \rightarrow c \rightarrow d$  on a fully-connected graph translates to a sequence  $1 \rightarrow 2 \rightarrow 3\langle 1 \rangle \rightarrow 4\langle 1, 2 \rangle$ , where  $\langle \cdot \rangle$  represents the named neighbors.

The pseudocode of both algorithms can be found in Appendix A.2. We now show the following:

**Proposition 2.3.** *A recording function  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$  that uses anonymization, optionally with named neighbors, is invariant.*

While anonymization was originally motivated by privacy concerns in (Micali & Zhu, 2016) (also see Appendix A.6), Proposition 2.3 offers a new justification based on invariance. Our design of recording named neighbors is novel, and is inspired by sublinear algorithms that probe previously discovered neighbors in a walk (Dasgupta et al., 2014). In our context, it is useful since whenever a walk visits a set of vertices  $S \subseteq V(G)$  it automatically records the entire induced subgraph  $G[S]$ . As a result, to record all edges of a graph, a walk only has to visit all vertices. While traversing all edges, i.e. edge cover time, is  $O(n^3)$  (Zuckerman, 1991) in general, it is possible to choose a walk algorithm that takes only  $O(n^2)$  time to visit all  $n$  vertices. MDLR in Equation 6 exactly achieves this.

**Reader neural network** A reader neural network  $f_\theta : \mathbf{z} \mapsto \hat{\mathbf{y}}$  processes the record  $\mathbf{z}$  of the random walk and outputs a prediction  $\hat{\mathbf{y}}$  in  $\mathbb{R}^d$ . As in Proposition 2.1, there is no invariance constraint imposed on  $f_\theta$ , and any neural network that accepts the recorded walks and has a sufficient expressive power can be used (we will make this precise in Section 3.1). This is in contrast to MPNNs where invariance is hard-coded in feature mixing operations. Also, our record  $\mathbf{z}$  can take any format, such as a matrix, byte sequence, or plain text, as long as  $f_\theta$  accepts it. Thus, it is possible (while not required) to choose the record to be plain text, such as "1-2-3-1", and choose  $f_\theta$  to be a pre-trained language model. This offers expressive power (Yun et al., 2020) and has a potential benefit of knowledge transfer from language domain (Lu et al., 2022; Rothermel et al., 2021).

## 2.2 VERTEX-LEVEL TASKS

We now consider vertex-level tasks. In case of finite graphs, we may simply frame a vertex-level task  $(G, v) \mapsto \mathbf{y}$  as a graph-level task  $G' \mapsto \mathbf{y}$  where  $G'$  is  $G$  with its vertex  $v$  marked. Then, we can solve  $G' \mapsto \mathbf{y}$  by querying an RWNN  $X_\theta(G, v)$  to start its walk at  $v_0 = v$ .

A more interesting case is when  $G$  is an infinite graph that is only locally finite, i.e. has finite degrees. This simulates problems on a large graph such as transductive classification. In this case, we may assume that our target function depends on finite local structures (Tahmasebi et al., 2023).

Let  $r \geq 1$ , and  $\mathbb{B}_r := \{B_r(v)\}$  be the collection of local balls in  $G$  of radius  $r$  centered at  $v \in V(G)$ . We would like to model a vertex-level function  $\phi : \mathbb{B}_r \rightarrow \mathbb{R}^d$  on  $G$  using an RWNN  $X_\theta(\cdot)$  by querying the vertex of interest,  $X_\theta(G, v)$ . We assume that  $\phi$  is isomorphism invariant:

$$\phi(B_r(v)) = \phi(B_r(u)), \quad \forall B_r(v) \simeq B_r(u). \quad (8)$$

The probabilistic invariance of  $X_\theta(\cdot)$  is defined as follows:

$$X_\theta(G, v) \stackrel{d}{=} X_\theta(G, u), \quad \forall B_r(v) \simeq B_r(u). \quad (9)$$

We can achieve probabilistic invariance by choosing the walk algorithm and recording function similar to the graph-level case<sup>3</sup>. As a modification, we query  $X_\theta(G, v)$  with the starting vertex  $v_0 = v$  of the random walk. Anonymization informs  $v$  to the reader NN by always naming it as 1.

Then, we make a key choice of localizing random walks with restarts. That is, we reset a walk to its starting vertex  $v_0$  either with a probability  $\alpha \in (0, 1)$  at each step or periodically every  $k$  steps. Restarting walks tend to stay more around starting vertex  $v_0$ , and were used to implement locality bias in personalized PageRank algorithm for search (Page et al., 1999). This localizing effect is crucial in our context since a walk can drift away from  $B_r(v_0)$  before recording all necessary information, which may take an infinite time to return as  $G$  is infinite (Janson & Peres, 2012). Restarts make the return to  $v_0$  mandatory, ensuring that  $B_r(v_0)$  can be recorded in a finite expected time. We show:

**Theorem 2.4.** *For a uniform random walk on an infinite graph  $G$  starting at  $v$ , the vertex and edge cover times of the finite local ball  $B_r(v)$  are not always finitely bounded.*

**Theorem 2.5.** *In Theorem 2.4, if the random walk restarts at  $v$  with any nonzero probability  $\alpha$  or any period  $k \geq r + 1$ , the vertex and edge cover times of  $B_r(v)$  are always finite.*

## 3 ANALYSIS

In Section 2, we have described the design of RWNNs, primarily relying on the principle of (probabilistic) invariance. In this section, we provide in-depth analysis on their expressive power and relations to issues in MPNNs such as over-smoothing, over-squashing, and under-reaching.

### 3.1 EXPRESSIVE POWER

Intuitively, if the records of random walks contain enough information such that the structures of interest, e.g. graph  $G$  or local ball  $B_r(v)$ , can be fully recovered, a powerful reader NN such as an MLP (Hornik et al., 1989) or a transformer (Yun et al., 2020) on these records would be able to approximate any function of interest. Our analysis formalizes this intuition.

We first consider using an RWNN  $X_\theta(\cdot)$  to universally approximate graph-level functions  $\phi(\cdot)$  in probability, as defined in Abboud et al. (2021).

**Definition 3.1.**  *$X_\theta(\cdot)$  is a universal approximator of graph-level functions in probability if, for all invariant functions  $\phi : \mathbb{G}_n \rightarrow \mathbb{R}$  for a given  $n \geq 1$ , and  $\forall \epsilon, \delta > 0$ , there exist choices of length  $l$  of the random walk and network parameters  $\theta$  such that the following holds:*

$$\text{Prob}[|\phi(G) - X_\theta(G)| < \epsilon] > 1 - \delta, \quad \forall G \in \mathbb{G}_n. \quad (10)$$

We show that, if the random walk is long enough and the reader NN  $f_\theta$  is universal, an RWNN  $X_\theta(\cdot)$  is capable of graph-level universal approximation. The length of the walk  $l$  controls the confidence  $> 1 - \delta$ , with edge cover time  $C_E(G)$  or vertex cover time  $C_V(G)$  playing a central role.

<sup>3</sup>This is assuming that the random walks are localized in  $B_r(v)$  and  $B_r(u)$ , e.g. with restarts.

**Theorem 3.2.** An RWNN  $X_\theta(\cdot)$  with a sufficiently powerful  $f_\theta$  is a universal approximator of graph-level functions in probability (Definition 3.1) if it satisfies either of the below:

- It uses anonymization to record random walks of lengths  $l > C_E(G)/\delta$ .
- It uses anonymization + named neighbors to record walks of lengths  $l > C_V(G)/\delta$ .

While both cover times are  $O(n^3)$  for uniform random walks (Aleliunas et al., 1979; Zuckerman, 1991), we can use MDLR in Equation 6 to achieve an  $O(n^2)$  vertex cover time, in conjunction with named neighborhood recording. While universality can be in principle achieved with uniform random walks, our design reduces the worst-case length  $l$  required for the desired reliability  $> 1 - \delta$ .

We now show an analogous result for universal approximation of vertex-level functions.

**Definition 3.3.**  $X_\theta(\cdot)$  is a universal approximator of vertex-level functions in probability if, for all invariant functions  $\phi : \mathbb{B}_r \rightarrow \mathbb{R}$  for a given  $r \geq 1$ , and  $\forall \epsilon, \delta > 0$ , there exist choices of length  $l$  and restart probability  $\alpha$  or period  $k$  of the random walk and network parameters  $\theta$  such that:

$$\text{Prob}[|\phi(B_r(v)) - X_\theta(G, v)| < \epsilon] > 1 - \delta, \quad \forall B_r(v) \in \mathbb{B}_r. \quad (11)$$

We show the following, using local vertex cover time  $C_V(B_r(v))$  and edge cover time  $C_E(B_r(v))$ :

**Theorem 3.4.** An RWNN  $X_\theta(\cdot)$  with a sufficiently powerful  $f_\theta$  and any nonzero restart probability  $\alpha$  or restart period  $k \geq r + 1$  is a universal approximator of vertex-level functions in probability (Definition 3.3) if it satisfies either of the below for all  $B_r(v) \in \mathbb{B}_r$ :

- It uses anonymization to record random walks of lengths  $l > C_E(B_r(v))/\delta$ .
- It uses anonymization + named neighbors to record walks of lengths  $l > C_V(B_r(v))/\delta$ .

Since restarts are required to finitely bound the local cover times on an infinite graph (Section 2.2), non-restarting walks cannot support vertex-level universality in general, and our design is obligatory.

### 3.2 OVER-SMOOTHING, OVER-SQUASHING, AND UNDER-REACHING

Often, in MPNNs, each layer operates by passing features over edges and mixing them using weights deduced from e.g. adjacency matrix. This ties them to the topology of the input graph, and a range of prior work has shown how this relates to the well-known issues of over-smoothing, over-squashing, and under-reaching. We connect these results and RWNNs to verify if similar issues may take place.

Let  $G$  be a connected non-bipartite graph with row-normalized adjacency matrix  $P$ . We consider a linearized MPNN, where the vertex features  $\mathbf{h}^{(0)}$  are initialized as some probability vector  $\mathbf{x}$ , and updated by  $\mathbf{h}^{(t+1)} = \mathbf{h}^{(t)}P$ . This simplification is often useful in understanding the aforementioned issues (Giraldo et al., 2023; Zhao & Akoglu, 2019). Specifically, in this model:

- Over-smoothing happens as the features exponentially converge to a stationary vector  $\mathbf{h}^{(l)} \rightarrow \pi$  as  $l \rightarrow \infty$ , smoothing out the input  $\mathbf{x}$  (Giraldo et al., 2023).
- Over-squashing and under-reaching occur when a feature  $\mathbf{h}_u^{(l)}$  becomes insensitive to distant input  $\mathbf{x}_v$ . While under-reaching refers to insufficient depth  $l < \text{diam}(G)$  (Barceló et al., 2020), over-squashing refers to features getting overly compressed at bottlenecks of  $G$ , even with sufficient depth  $l$ . The latter is described by the Jacobian  $|\partial \mathbf{h}_u^{(l)} / \partial \mathbf{x}_v| \leq [\sum_{t=0}^l P^t]_{uv}$ ,<sup>4</sup> as the bound often decays exponentially with  $l$  (Topping et al., 2022; Black et al., 2023).

What do these results tell us about RWNNs? We can see that, while  $P$  drives feature mixing in the message passing schema, it can be also interpreted as the transition probability matrix of uniform random walk where  $P_{uv}$  is the probability of walking from  $u$  to  $v$ . This parallelism motivates us to design an analogous, simplified RWNN and study its behavior.

We consider a simple RWNN that runs a uniform random walk  $v_0 \rightarrow \dots \rightarrow v_l$ , reads the record  $\mathbf{x}_{v_0} \rightarrow \dots \rightarrow \mathbf{x}_{v_l}$  by averaging, and outputs it as  $\mathbf{h}^{(l)}$ . Like linear MPNN, the model involves  $l$  steps

<sup>4</sup>This bound is obtained by applying Lemma 3.2 of Black et al. (2023) to our linearized MPNN.

of time evolution through  $P$ . However, while MPNN uses  $P$  to process features, this model uses  $P$  only to obtain a record of input, with feature processing decoupled. We show that in this model, over-smoothing does not occur as in simple MPNNs<sup>5</sup>:

**Theorem 3.5.** *The simple RWNN outputs  $\mathbf{h}^{(l)} \rightarrow \mathbf{x}^\top \boldsymbol{\pi}$  as  $l \rightarrow \infty$ .*

Even if the time evolution through  $P$  happens fast (i.e.  $P$  is rapidly mixing) or with many steps  $l$ , the model is resistant to over-smoothing as the input  $\mathbf{x}$  always affects the output. In fact, if  $P$  is rapidly mixing, we may expect improved behaviors based on Section 3.1 as the cover times could reduce.

On the other hand, we show that over-squashing manifests as probabilistic under-reaching:

**Theorem 3.6.** *Let  $\mathbf{h}_u^{(l)}$  be output of the simple RWNN queried with  $u$ . Then:*

$$\mathbb{E} \left[ \left\| \frac{\partial \mathbf{h}_u^{(l)}}{\partial \mathbf{x}_v} \right\| \right] = \frac{1}{l+1} \left[ \sum_{t=0}^l P^t \right]_{uv} \rightarrow \boldsymbol{\pi}_v \quad \text{as } l \rightarrow \infty. \quad (12)$$

The equation shows that the feature Jacobians are bounded by the sum of powers of  $P$ , same as in simple MPNN. Both models are subject to over-squashing phenomenon that is similarly formalized, but manifests through different mechanisms. While in message passing the term is related to over-compression of features at bottlenecks (Topping et al., 2022), in RWNNs it is related to exponentially decaying probability of reaching a distant vertex  $v$ , i.e. probabilistic under-reaching.

In many MPNNs, it is understood that the topology of the input graph inevitably induces a trade-off between over-smoothing and over-squashing (Nguyen et al., 2023; Giraldo et al., 2023). Our results suggest that RWNNs avoid the trade-off, and we can focus on overcoming under-reaching e.g. with long or rapidly-mixing walks, while not worrying much about over-smoothing. Design choices such as MDLR (Equation 6) and non-backtracking (Alon et al., 2007) can be understood as achieving this.

## 4 RELATED WORK

We briefly review the related work. An extended discussion can be found in Appendix A.6.

**Random walks for learning on graphs** In graph learning, random walks have received interest due to compatibility with sequence learning methods (Table 1). DeepWalk (Perozzi et al., 2014) and node2vec (Grover & Leskovec, 2016) used skip-gram models on walks. CRaWl (Tönshoff et al., 2023) used 1D CNN on sliding window-based walk records, with expressive power bounded by the window size. We discuss the relation between CRaWl and our approach in depth in Appendix A.6. AgentNet (Martinkus et al., 2023) learns agents that walk on a graph while recurrently updating their features. While optimizing the walk strategy, this may trade off speed as training requires recurrent roll-out of the network. Our method allows pairing simple and fast walkers, such as MDLR, with parallelizable networks such as transformers. WalkLM (Tan et al., 2023) proposed a fine-tuning method for language models on walks on text-attributed graphs. By utilizing anonymization, our approach is able to process graphs even if no text attribute is given. Lastly, a concurrent work (Wang & Cho, 2024) arrived at a similar use of anonymization combined with RNNs.

**Probabilistic invariant neural networks** Whenever a learning problem is compatible with symmetry, incorporating the associated invariance structure to the hypothesis class often leads to generalization benefit (Bronstein et al., 2021; Elesedy, 2022; 2023). This is also the case for probabilistic invariant NNs (Lyle et al., 2020; Bloem-Reddy & Teh, 2020), which includes our approach. Probabilistic invariant NNs have recently gained interest due to their potential of achieving higher expressive powers compared to deterministic counterparts (Cotta et al., 2023; Sieradzki et al., 2022). In graph learning, this is often achieved with stochastic symmetry breaking between vertices using randomized features (Loukas, 2020; Puny et al., 2020; Abboud et al., 2021; Kim et al., 2022), vertex orderings (Murphy et al., 2019; Kim et al., 2023), or dropout (Papp et al., 2021). Our approach can be understood as using random walk as a symmetry-breaking mechanism for probabilistic invariance, which provides an additional benefit of natural compatibility with sequence learning methods.

<sup>5</sup>While we show not forgetting  $\mathbf{x}$  for brevity, we may extend to initial vertex  $v_0$  using its anonymization as 1.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

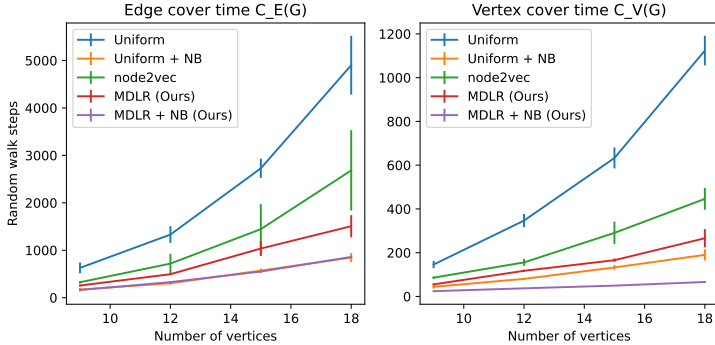


Figure 2: Cover times of random walks on varying sizes of lollipop graphs. NB is non-backtracking.

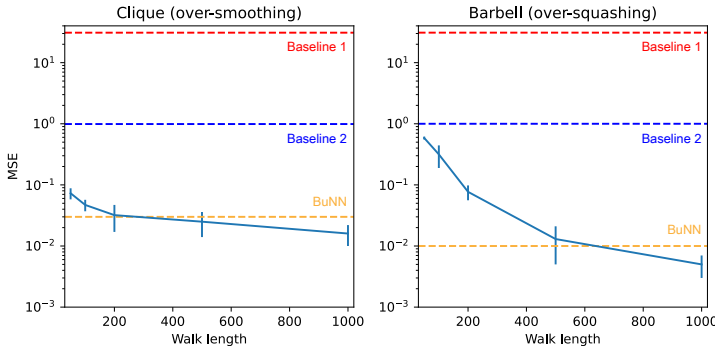


Figure 3: Over-smoothing and over-squashing (MSE ↓) for various walk lengths  $l$ .

Table 2: Over-smoothing and over-squashing results in comparison to baselines from Bamberger et al. (2024). We report aggregated mean squared error (MSE ↓) for 4 randomized runs.

Method	Clique (over-smoothing)	Barbell (over-squashing)
Baseline 1	$30.94 \pm 0.42$	$30.97 \pm 0.42$
Baseline 2	$0.99 \pm 0.08$	$1.00 \pm 0.07$
MLP	$1.10 \pm 0.08$	$1.08 \pm 0.07$
GCN	$29.65 \pm 0.34$	$1.05 \pm 0.08$
SAGE	$0.86 \pm 0.10$	$0.90 \pm 0.29$
GAT	$20.97 \pm 0.40$	$1.07 \pm 0.09$
NSD	$0.08 \pm 0.02$	$1.09 \pm 0.15$
BuNN	$0.03 \pm 0.01$	$0.01 \pm 0.07$
RWNN-transformer (Ours)	<b><math>0.016 \pm 0.006</math></b>	<b><math>0.005 \pm 0.002</math></b>

## 5 EXPERIMENTS

We perform a series of experiments to demonstrate RWNNs. We implement random walk algorithms in C++ based on Chenebaux (2020), which produces good throughput (<0.1 seconds for 10k steps) even without GPU acceleration in (Tönshoff et al., 2023). Pseudocode is given in Appendix A.2. We implement our training and inference pipelines in PyTorch. Supplementary experiments including substructure counting, more real-world datasets, and link prediction, are in Appendix A.4.

### 5.1 SYNTHETIC EXPERIMENTS

On synthetic setups, we first verify our claims on cover times in Section 2.1, focusing on the utility of MDLR (Equation 6), non-backtracking, and neighborhood recording. We measure the edge cover times  $C_E(G)$  and vertex cover times  $C_V(G)$  of random walk algorithms discussed in the main text, on varying sizes of lollipop graphs  $G$  (Feige, 1995) which are commonly used to establish upper bounds of cover times. The results are in Figure 2. We find that (1) MDLR achieves a significant speed-up compared to uniform and node2vec walks, (2) adding non-backtracking significantly improves the behavior of base first-order walks in all cases, and (3) edge cover times are often much larger compared to vertex cover times, strongly justifying the use of neighborhood recording (Theorem 3.2).



Table 3: Isomorphism learning results. We report accuracy of classifying training data rounded to the first decimal point. \*, †, ‡, §, and ◊ are from Papp et al. (2021), Murphy et al. (2019), Zhao et al. (2022b), Martinkus et al. (2023), and Alvarez-Gonzalez et al. (2024), respectively.

Method	CSL	SR16	SR25
	1, 2-WL fails	3-WL fails	3-WL fails
GIN	10.0%†	50.0%§	6.7%‡
PPGN	100.0%§	50.0%§	6.7%‡
GIN-AK+	-	-	6.7%‡
PPGN-AK+	-	-	100.0%‡
GIN+ELENE	-	-	6.7%◊
GIN+ELENE-L (ED)	-	-	100.0%◊
GIN-RNI	16.0%*	-	-
GIN-RP	37.6%†	-	-
GIN-Dropout	82.0%*	-	-
RW-AgentNet	100.0%§	50.5%§	-
AgentNet	100.0%§	100.0%§	6.7%
CRaWl	100.0%§	100.0%§	46.6%
RWNN-DeBERTa (Ours)	100.0%	100.0%	100.0%

Then, we verify our claims in Section 3.2 using Clique and Barbell datasets of Bamberger et al. (2024) that explicitly test for over-smoothing and over-squashing, respectively. Our RWNN uses MDLR walks with non-backtracking, and the reader NN is a 1-layer transformer encoder with width 128, matching the baselines. We precisely follow the experimental procedure of Bamberger et al. (2024). Figure 3 shows the results for varying walk lengths  $l$ . The model performs well on Clique overall, while Barbell requires scaling the walk length. This agrees with our claims in Section 3.2 that RWNNs in general avoid over-smoothing, but over-squashing manifests as probabilistic under-reaching (and is therefore mitigated by sufficiently long walks). With  $l = 1000$ , our model in Table 2 outperforms all baselines including BuNN which is specialized for mitigating over-smoothing and over-squashing.

## 5.2 GRAPH ISOMORPHISM LEARNING

We now test the claims on expressive power in Section 3.1 using three challenging datasets where the task is recognizing the isomorphism type of input graph where certain WL test fails.

The Circular Skip Links (CSL) graphs dataset (Murphy et al., 2019) contains 10 non-isomorphic regular graphs with 41 vertices of degree 4. Distinguishing these graphs requires computing lengths of skip links, and 1-WL and 2-WL tests fail. The  $4 \times 4$  rook’s and Shrikhande graphs (Alvarez-Gonzalez et al., 2024), which we call SR16, are a pair of strongly regular graphs with 16 vertices of degree 6. Distinguishing the pair requires detecting 4-cliques, and 3-WL test fails. The SR25 dataset (Balcilar et al., 2021) contains 15 strongly regular graphs with 25 vertices of degree 12, on which 3-WL test fails. Examples of the considered graphs and random walks on them can be found in Appendix A.2.

Our RWNN uses MDLR random walk with non-backtracking, and recording function with anonymization and named neighbors that produces plain text (Algorithm 3). We use pre-trained DeBERTa-base language model as the reader NN to leverage its capacity, and fine-tune it with cross-entropy loss for at most 100k steps using AdamW optimizer with  $2e-5$  learning rate and 0.01 weight decay on  $8 \times$  RTX 3090 GPUs. We truncate the input to 512 tokens which is the maximum allowed by memory constraint. We use batch size 256, and accumulate gradients for 8 steps for SR25. At test time, we ensemble 4 predictions of the network by averaging classification logits.

The results are in Table 3. MPNNs that align with certain WL test fail when asked to solve harder problems, e.g. GIN aligns with 1-WL and fails in CSL. A limited set of state-of-the-art neural networks solve SR25, but at the cost of introducing specialized structural features (Alvarez-Gonzalez et al., 2024). Alternative approaches based on stochastic symmetry-breaking e.g. random node identification, often fail on CSL although universal approximation is possible in theory (Abboud et al., 2021), possibly due to learning difficulties. For algorithms based on walks, AgentNet and CRaWl solve CSL and SR16, while failing on SR25. This is because learning a policy to walk in AgentNet can be challenging in complex tasks especially at the early stage of training, and the expressiveness of CRaWl is limited by the receptive field of the 1D CNN. Our approach based on DeBERTa language model overcomes the problems, demonstrating as the first RWNN that solves SR25.

In Appendix A.3, we further provide visualizations of learned attentions by mapping attention weights on text records of walks to input graphs. We find that the models often focus on sparse, connected substructures, which presumably provide discriminative information on the isomorphism types.

Table 4: Test accuracy on ogbn-arxiv. † denotes using validation labels for label propagation or in-context learning following Huang et al. (2020). For Llama 3, we ensemble 5 predictions by voting.

Method	Accuracy
MLP	55.50%
node2vec	70.07%
GCN	71.74%
GraphSAGE	71.49%
GAT	73.91%
RevGAT	74.26%
Label propagation	68.50%
C&S	72.62%
C&S†	74.02%
Llama2-13b zero-shot	44.23%
GPT-3.5 zero-shot	73.50%
DeBERTa, fine-tuned	74.13%
Llama3-8b zero-shot	50.20%
Llama3-8b one-shot	52.49%
Llama3-8b one-shot†	52.54%
Llama3-70b zero-shot	65.33%
Llama3-70b one-shot†	67.57%
RWNN-Llama3-8b (Ours)	71.10%
RWNN-Llama3-8b (Ours)†	73.11%
RWNN-Llama3-70b (Ours)†	<b>74.75%</b>

### 5.3 REAL-WORLD TRANSDUCTIVE CLASSIFICATION

Previous results considered tasks on undirected, unattributed, and relatively small graphs. We now show a preliminary result that RWNN based on Llama 3 (meta llama, 2024) language models can solve real-world problem on a large graph with directed edges and textual attributes. We use ogbn-arxiv (Hu et al., 2020), a citation network of 169,343 arXiv papers with title and abstract attributes. The task is transductive classification into 40 areas such as "CS.AI" using a set of labeled vertices.

We consider two representative types of baselines. The first reads title and abstract of each vertex using a language model and solves vertex-wise classification problem, ignoring graph structure. The second initializes vertex features as language model embeddings and trains an MPNN, at the risk of over-compressing the text. To take the best of both worlds, we design the recording function so that, not only it does basic operations such as anonymization, it records a complete information of the local subgraph including title and abstract, edge directions, and notably, **labels in case of labeled vertices** (Appendix A.2) (Sato, 2024). The resulting record naturally includes a number of input-label pairs of the classification task at hand, implying that we can frame transductive classification problem as a simple in-context learning problem (Brown et al., 2020). This allows training-free application of Llama 3 (meta llama, 2024) language model for transductive classification.

The results are in Table 4. Our models based on frozen Llama 3 perform competitively against a range of previous MPNNs on text embeddings, as well as outperforming language models that perform vertex-wise predictions ignoring graph structures such as GPT-3.5 and fine-tuned DeBERTa. Especially, our model largely outperforms one-shot baselines, which are given 40 randomly chosen labeled examples (one per class). This is surprising as our model observes fewer vertices, 29.17 in average, due to other recorded information. This is presumably since the record produced by our algorithm informs useful graph structure to the language model to quickly learn the task in-context compared to randomly chosen shots. In Appendix A.3, we visualize the attention weights, verifying that the model makes use of the graph structure recorded by random walks to make predictions.

As a final note, our approach is related to label propagation algorithms (Zhu & Ghahramani, 2002; Zhu, 2005; Grady, 2006) for transductive classification, which makes predictions by running random walks and probing the distribution of visited labeled vertices. The difference is, in our approach, the reader NN i.e. language model can appropriately use other information such as attributes, as well as do meaningful non-linear processing rather than simply probing the input. As shown in Table 4, our approach outperforms label propagation, verifying our intuition.

## 6 CONCLUSION

We contributed a principled understanding of RWNNs, where random walks on graphs are recorded and processed by NNs to make predictions. We analyzed invariance, expressive power, and information propagation, showing that we can design RWNNs to be invariant and universal without constraining its neural network. Experiments support the utility of our approach.

## REFERENCES

- 540  
541  
542 Ralph Abboud, İsmail İlkan Ceylan, Martin Grohe, and Thomas Lukasiewicz. The surprising power  
543 of graph neural networks with random node initialization. In *IJCAI*, 2021. (pages 2, 5, 7, 9)
- 544  
545 Mohammed Abdullah. The cover time of random walks on graphs. *arXiv*, 2012. (page 46)
- 546  
547 Mohammed Amin Abdullah, Colin Cooper, and Moez Draief. Speeding up cover time of sparse  
graphs using local knowledge. In *IWOCA*, 2015. (pages 4, 46)
- 548  
549 Romas Aleliunas, Richard M. Karp, Richard J. Lipton, László Lovász, and Charles Rackoff. Random  
550 walks, universal traversal sequences, and the complexity of maze problems. In *Annual Symposium  
on Foundations of Computer Science*, 1979. (pages 6, 42, 46)
- 551  
552 Noga Alon, Itai Benjamini, Eyal Lubetzky, and Sasha Sodin. Non-backtracking random walks mix  
553 faster. *Communications in Contemporary Mathematics*, 2007. (pages 7, 18, 46)
- 554  
555 Nurudin Alvarez-Gonzalez, Andreas Kaltenbrunner, and Vicenç Gómez. Improving subgraph-GNNs  
556 via edge-level ego-network encodings. *Transactions on Machine Learning Research*, 2024. ISSN  
2835-8856. (page 9)
- 557  
558 Francesca Arrigo, Desmond J. Higham, and Vanni Noferini. Non-backtracking pagerank. *J. Sci.  
559 Comput.*, 2019. (page 46)
- 560  
561 Muhammet Balcilar, Pierre Héroux, Benoit Gaüzère, Pascal Vasseur, Sébastien Adam, and Paul  
Honeine. Breaking the limits of message passing graph neural networks. In *ICML*, 2021. (page 9)
- 562  
563 Jacob Bamberger, Federico Barbero, Xiaowen Dong, and Michael Bronstein. Bundle neural networks  
564 for message diffusion on graphs. *arXiv*, 2024. (pages 8, 9)
- 565  
566 Pablo Barceló, Egor V. Kostylev, Mikaël Monet, Jorge Pérez, Juan L. Reutter, and Juan Pablo Silva.  
The logical expressiveness of graph neural networks. In *ICLR*, 2020. (pages 6, 47)
- 567  
568 Peter W. Battaglia, Jessica B. Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinícius Flores  
569 Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner,  
Çaglar Gülçehre, H. Francis Song, Andrew J. Ballard, Justin Gilmer, George E. Dahl, Ashish  
570 Vaswani, Kelsey R. Allen, Charles Nash, Victoria Langston, Chris Dyer, Nicolas Heess, Daan  
Wierstra, Pushmeet Kohli, Matthew M. Botvinick, Oriol Vinyals, Yujia Li, and Razvan Pascanu.  
571 Relational inductive biases, deep learning, and graph networks. *arXiv*, 2018. (page 1)
- 572  
573 Anna Ben-Hamou, Roberto I. Oliveira, and Yuval Peres. Estimating graph parameters via random  
574 walks with restarts. In *Annual ACM-SIAM Symposium on Discrete Algorithms*, 2018. (page 46)
- 575  
576 Suman K Bera and C Seshadhri. How to count triangles, without seeing the whole graph. In *KDD*,  
577 2020. (page 47)
- 578  
579 Mitchell Black, Zhengchao Wan, Amir Nayyeri, and Yusu Wang. Understanding oversquashing in  
580 gnn through the lens of effective resistance. In *ICML*, 2023. (pages 6, 47)
- 581  
582 Benjamin Bloem-Reddy and Yee Whye Teh. Probabilistic symmetries and invariant neural networks.  
*J. Mach. Learn. Res.*, 2020. (pages 2, 3, 7)
- 583  
584 Michael M. Bronstein, Joan Bruna, Taco Cohen, and Petar Velickovic. Geometric deep learning:  
585 Grids, groups, graphs, geodesics, and gauges. *arXiv*, 2021. (page 7)
- 586  
587 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,  
Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
588 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,  
Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott  
589 Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya  
Sutskever, and Dario Amodei. Language models are few-shot learners. In *NeurIPS*, 2020. (page  
590 10)
- 591  
592  
593 Eric Richard Bussian. *Bounding the edge cover time of random walks on graphs*. Georgia Institute of  
Technology, 1996. (page 46)

- 594 Lei Cai, Jundong Li, Jie Wang, and Shuiwang Ji. Line graph neural networks for link prediction.  
595 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021. (page 33)  
596
- 597 carnage. Expected hitting time for simple random walk from origin to point (x,y) in 2d-integer-  
598 grid. MathOverflow, 2012. URL <https://mathoverflow.net/questions/112470>.  
599 [Online:] <https://mathoverflow.net/questions/112470>. (page 38)
- 600 Ashok K. Chandra, Prabhakar Raghavan, Walter L. Ruzzo, Roman Smolensky, and Prasoos Tiwari.  
601 The electrical resistance of a graph captures its commute and cover times. *Comput. Complex.*,  
602 1997. (page 47)  
603
- 604 Runjin Chen, Tong Zhao, Ajay Jaiswal, Neil Shah, and Zhangyang Wang. Llaga: Large language and  
605 graph assistant. *arXiv*, 2024a. (page 32)
- 606 Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Longlora:  
607 Efficient fine-tuning of long-context large language models. *arXiv*, 2023a. (page 47)  
608
- 609 Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph  
610 isomorphism testing and function approximation with gnns. In *NeurIPS*, 2019. (page 1)
- 611 Zhengdao Chen, Lei Chen, Soledad Villar, and Joan Bruna. Can graph neural networks count  
612 substructures? In *NeurIPS*, 2020. (page 33)  
613
- 614 Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei  
615 Yin, Wenqi Fan, Hui Liu, and Jiliang Tang. Exploring the potential of large language models  
616 (llms) in learning on graphs. *SIGKDD Explor.*, 2023b. (page 47)
- 617 Zhikai Chen, Haitao Mao, Jingzhe Liu, Yu Song, Bingheng Li, Wei Jin, Bahare Fatemi, Anton  
618 Tsitsulin, Bryan Perozzi, Hui Liu, and Jiliang Tang. Text-space graph foundation models: Compre-  
619 hensive benchmarks and new insights. *arXiv*, 2024b. (page 32)  
620
- 621 Maixent Chenebaux. graph-walker: Fastest random walks generator on networkx graphs. <https://github.com/kerighan/graph-walker>, 2020. (page 8)  
622
- 623 Flavio Chiericetti, Anirban Dasgupta, Ravi Kumar, Silvio Lattanzi, and Tamás Sarlós. On sampling  
624 nodes in a network. In *WWW*, 2016. (page 46)  
625
- 626 Don Coppersmith, Uriel Feige, and James B. Shearer. Random walks on regular and irregular graphs.  
627 *SIAM J. Discret. Math.*, 1996. (page 46)
- 628 Leonardo Cotta, Gal Yehuda, Assaf Schuster, and Chris J. Maddison. Probabilistic invariant learning  
629 with randomized linear classifiers. In *NeurIPS*, 2023. (page 7)  
630
- 631 George Cybenko. Approximation by superpositions of a sigmoidal function. *Math. Control. Signals*  
632 *Syst.*, 1989.
- 633 Anirban Dasgupta, Ravi Kumar, and Tamas Sarlos. On estimating the average degree. In *WWW*,  
634 2014. (pages 4, 46)  
635
- 636 Roe David and Uriel Feige. Random walks with the minimum degree local rule have  $o(n^2)$  cover  
637 time. *SIAM J. Comput.*, 2018. (pages 2, 4, 46)
- 638 Artur Back de Luca and Kimon Fountoulakis. Simulation of graph algorithms with looped transform-  
639 ers. *arXiv*, 2024. (page 47)  
640
- 641 Grégoire Delétang, Anian Ruoss, Jordi Grau-Moya, Tim Genewein, Li Kevin Wenliang, Elliot Catt,  
642 Chris Cundy, Marcus Hutter, Shane Legg, Joel Veness, and Pedro A. Ortega. Neural networks and  
643 the chomsky hierarchy. In *ICLR*, 2023. (page 47)
- 644 Karel Devriendt and Renaud Lambiotte. Discrete curvature on graphs from the effective resistance.  
645 *Journal of Physics: Complexity*, 2022. (page 47)  
646
- 647 Jian Ding, James R Lee, and Yuval Peres. Cover times, blanket times, and majorizing measures. In  
*Annual ACM symposium on Theory of computing*, 2011. (page 46)

- 648 Peter G Doyle and J Laurie Snell. *Random walks and electric networks*. American Mathematical  
649 Soc., 1984. (page 47)
- 650
- 651 Ioana Dumitriu, Prasad Tetali, and Peter Winkler. On playing golf with two balls. *SIAM J. Discret.*  
652 *Math.*, 2003. (pages 38, 46)
- 653 Vijay Prakash Dwivedi, Ladislav Rampásek, Michael Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu,  
654 and Dominique Beaini. Long range graph benchmark. In *NeurIPS*, 2022. (page 32)
- 655
- 656 Nadav Dym, Hannah Lawrence, and Jonathan W. Siegel. Equivariant frames and the impossibility of  
657 continuous canonicalization. In *ICML*, 2024. (pages 2, 3)
- 658 Bryn Elesedy. Group symmetry in pac learning. In *ICLR workshop on geometrical and topological*  
659 *representation learning*, 2022. (page 7)
- 660
- 661 Bryn Elesedy. *Symmetry and Generalisation in Machine Learning*. PhD thesis, University of Oxford,  
662 2023. (page 7)
- 663 William Falcon. Pytorch lightning. <https://github.com/Lightning-AI/lightning>,  
664 2019.
- 665
- 666 Dario Fasino, Arianna Tonetto, and Francesco Tudisco. Hitting times for non-backtracking random  
667 walks. *arXiv*, 2021. (page 46)
- 668 Bahare Fatemi, Jonathan Halcrow, and Bryan Perozzi. Talk like a graph: Encoding graphs for large  
669 language models. *arXiv*, 2023. (page 47)
- 670
- 671 Uriel Feige. A tight upper bound on the cover time for random walks on graphs. *Random Struct.*  
672 *Algorithms*, 1995. (pages 8, 46)
- 673 Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In  
674 *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- 675
- 676 Robert Fitzner and Remco van der Hofstad. Non-backtracking random walk. *Journal of Statistical*  
677 *Physics*, 2013. (page 18)
- 678 Agelos Georgakopoulos and Peter Winkler. New bounds for edge-cover by random walk. *Comb.*  
679 *Probab. Comput.*, 2014. (page 46)
- 680
- 681 Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural  
682 message passing for quantum chemistry. In *ICML*, 2017.
- 683 Francesco Di Giovanni, T. Konstantin Rusch, Michael M. Bronstein, Andreea Deac, Marc Lackenby,  
684 Siddhartha Mishra, and Petar Velickovic. How does over-squashing affect the power of gnns?  
685 *arXiv*, 2023. (page 47)
- 686
- 687 Jhony H. Giraldo, Konstantinos Skianis, Thierry Bouwmans, and Fragkiskos D. Malliaros. On the  
688 trade-off between over-smoothing and over-squashing in deep graph neural networks. In *CIKM*,  
689 2023. (pages 1, 2, 6, 7, 47)
- 690 Leo J. Grady. Random walks for image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2006.  
691 (pages 10, 47)
- 692
- 693 Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *KDD*, 2016.  
694 (pages 1, 4, 7, 18, 36)
- 695 William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large  
696 graphs. In *NeurIPS*, 2017.
- 697
- 698 Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: decoding-enhanced bert  
699 with disentangled attention. In *ICLR*, 2021. (page 2)
- 700 Xiaoxin He, Xavier Bresson, Thomas Laurent, Adam Perold, Yann LeCun, and Bryan Hooi. Harness-  
701 ing explanations: Llm-to-lm interpreter for enhanced text-attributed graph representation learning.  
In *ICLR*, 2023.

- 702 Jeffery Hein. Wald’s identity. (page 40)  
703
- 704 Kurt Hornik, Maxwell B. Stinchcombe, and Halbert White. Multilayer feedforward networks are  
705 universal approximators. *Neural Networks*, 1989. (page 5)
- 706 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang,  
707 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *ICLR, 2022*. (pages  
708 32, 47)
- 709
- 710 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta,  
711 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. In *NeurIPS*,  
712 2020. (page 10)
- 713 Qian Huang, Horace He, Abhay Singh, Ser-Nam Lim, and Austin R Benson. Combining label  
714 propagation and simple models out-performs graph neural networks. *arXiv*, 2020. (page 10)
- 715
- 716 Satoshi Ikeda, Izumi Kubo, and Masafumi Yamashita. The hitting and cover times of random walks  
717 on finite graphs using local degree information. *Theor. Comput. Sci.*, 2009. (page 46)
- 718 Sergey Ivanov and Evgeny Burnaev. Anonymous walk embeddings. In *ICML*, 2018. (page 1)
- 719
- 720 Svante Janson and Yuval Peres. Hitting times for random walks with restarts. *SIAM J. Discret. Math.*,  
721 2012. (pages 5, 38, 46)
- 722
- 723 Jeff D Kahn, Nathan Linial, Noam Nisan, and Michael E Saks. On the cover time of random walks  
724 on graphs. *Journal of Theoretical Probability*, 1989. (page 46)
- 725
- 726 Mark Kempton. Non-backtracking random walks and a weighted ihara’s theorem. *arXiv*, 2016. (page  
727 46)
- 728 Jinwoo Kim, Dat Nguyen, Seonwoo Min, Sungjun Cho, Moontae Lee, Honglak Lee, and Seunghoon  
729 Hong. Pure transformers are powerful graph learners. In *NeurIPS, 2022*. (page 7)
- 730 Jinwoo Kim, Dat Nguyen, Ayhan Suleymanzade, Hyeokjun An, and Seunghoon Hong. Learning  
731 probabilistic symmetrization for architecture agnostic equivariance. In *NeurIPS, 2023*. (page 7)
- 732
- 733 Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks.  
734 In *ICLR, 2017*.
- 735 Filip Klubicka, Alfredo Maldonado, Abhijit Mahalunkar, and John D. Kelleher. Synthetic, yet  
736 natural: Properties of wordnet random walk corpora and the impact of rare words on embedding  
737 performance. In *GWC, 2019*. (page 47)
- 738 Filip Klubicka, Alfredo Maldonado, Abhijit Mahalunkar, and John D. Kelleher. English wordnet  
739 random walk pseudo-corpora. In *LREC, 2020*. (page 47)
- 740
- 741 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E.  
742 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
743 serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating  
744 Systems Principles, 2023*.
- 745 Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan Zhang.  
746 One for all: Towards training one graph model for all classification tasks. In *ICLR, 2024*. (page 32)
- 747
- 748 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR, 2019*.
- 749
- 750 Andreas Loukas. What graph neural networks cannot learn: depth vs width. In *ICLR, 2020*. (page 7)
- 751 László Lovász. Random walks on graphs. *Combinatorics, Paul erdos is eighty*, 1993. (pages 46, 47)
- 752
- 753 Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Frozen pretrained transformers as  
754 universal computation engines. In *AAAI, 2022*. (page 4)
- 755
- 756 Clare Lyle, Mark van der Wilk, Marta Kwiatkowska, Yarin Gal, and Benjamin Bloem-Reddy. On the  
benefits of invariance in neural networks. *arXiv*, 2020. (page 7)

- 756 Karolis Martinkus, Pál András Papp, Benedikt Schesch, and Roger Wattenhofer. Agent-based graph  
757 neural networks. In *ICLR, 2023*. (pages 1, 7, 9)  
758
- 759 Nathan McNew. Random walks with restarts, 3 examples, 2013. (pages 38, 39, 46)  
760
- 761 meta llama. llama3: The official meta llama 3 github site. [https://github.com/  
762 meta-llama](https://github.com/meta-llama), 2024. (pages 2, 10)
- 763 Silvio Micali and Zeyuan Allen Zhu. Reconstructing markov processes from independent and  
764 anonymous experiments. *Discret. Appl. Math.*, 2016. (pages 1, 2, 4, 46)  
765
- 766 Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav  
767 Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks.  
768 In *AAAI*, 2019.
- 769 Ryan L. Murphy, Balasubramaniam Srinivasan, Vinayak A. Rao, and Bruno Ribeiro. Relational  
770 pooling for graph representations. In *ICML*, 2019. (pages 7, 9)  
771
- 772 Khang Nguyen, Nong Minh Hieu, Vinh Duc Nguyen, Nhat Ho, Stanley J. Osher, and Tan Minh  
773 Nguyen. Revisiting over-smoothing and over-squashing using ollivier-ricci curvature. In *ICML*,  
774 2023. (pages 2, 7, 47)
- 775 Roberto Oliveira. Mixing and hitting times for finite markov chains, 2012. (page 46)  
776
- 777 Yann Ollivier. Ricci curvature of markov chains on metric spaces. *Journal of Functional Analysis*,  
778 2009. (page 47)
- 779 Lawrence Page, Sergey Brin, Rajeev Motwani, Terry Winograd, et al. The pagerank citation ranking:  
780 Bringing order to the web, 1999. (pages 5, 47)  
781
- 782 Athina Panotopoulou. Bounds for edge-cover by random walks, 2013. (page 46)  
783
- 784 Pál András Papp, Karolis Martinkus, Lukas Faber, and Roger Wattenhofer. Dropgnn: Random  
785 dropouts increase the expressiveness of graph neural networks. *NeurIPS*, 2021. (pages 7, 9)  
786
- 787 Seonghyun Park, Narae Ryu, Gahee Kim, Dongyeop Woo, Se-Young Yun, and Sungsoo Ahn. Non-  
788 backtracking graph neural networks. In *NeurIPS Workshop: New Frontiers in Graph Learning*,  
2023. (page 47)
- 789 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor  
790 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward  
791 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,  
792 Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep  
793 learning library. In *NeurIPS*, 2019.
- 794 Yuval Peres and Perla Sousi. Mixing times are hitting times of large sets. *Journal of Theoretical  
795 Probability*, 2015. (page 46)  
796
- 797 Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: online learning of social representations.  
798 In *KDD*, 2014. (pages 1, 3, 4, 7)
- 799 Oleg Platonov, Denis Kuznedelev, Michael Diskin, Artem Babenko, and Liudmila Prokhorenkova. A  
800 critical look at the evaluation of gnns under heterophily: Are we really making progress? *arXiv*,  
801 2023. (page 32)  
802
- 803 Omri Puny, Heli Ben-Hamu, and Yaron Lipman. Global attention improves graph networks general-  
804 ization. *arXiv*, 2020. (page 7)
- 805 Omri Puny, Matan Atzmon, Edward J. Smith, Ishan Misra, Aditya Grover, Heli Ben-Hamu, and  
806 Yaron Lipman. Frame averaging for invariant and equivariant network design. In *ICLR, 2022*.  
807 (pages 2, 3)  
808
- 809 Brian Rappaport, Anuththari Gamage, and Shuchin Aeron. Faster clustering via non-backtracking  
random walks. *arXiv*, 2017. (page 18)

- 810 Danielle Rothermel, Margaret Li, Tim Rocktäschel, and Jakob Foerster. Don't sweep your learning  
811 rate under the rug: A closer look at cross-modal transfer of pretrained transformers. *arXiv*, 2021.  
812 (page 4)  
813
- 814 Clayton Sanford, Bahare Fatemi, Ethan Hall, Anton Tsitsulin, Mehran Kazemi, Jonathan Halcrow,  
815 Bryan Perozzi, and Vahab Mirrokni. Understanding transformer reasoning capabilities via graph  
816 algorithms. *arXiv*, 2024. (page 47)
- 817 Ryoma Sato. Training-free graph neural networks and the power of labels as features. *arXiv*, 2024.  
818 (pages 10, 32)  
819
- 820 Yuval Sieradzki, Nitzan Hodos, Gal Yehuda, and Assaf Schuster. Coin flipping neural networks. In  
821 *ICML*, 2022. (page 7)
- 822 Daniel Spielman. Spectral and algebraic graph theory. *Yale lecture notes, draft of December*, 2019.  
823 (page 47)  
824
- 825 Behrooz Tahmasebi, Derek Lim, and Stefanie Jegelka. The power of recursion in graph neural  
826 networks for counting substructures. In *AISTATS*, 2023. (page 5)  
827
- 828 Yanchao Tan, Zihao Zhou, Hang Lv, Weiming Liu, and Carl Yang. Walklm: A uniform language  
829 model fine-tuning framework for attributed graph embedding. In *NeurIPS*, 2023. (page 7)
- 830 Jan Tönshoff, Martin Ritzert, Hinrikus Wolf, and Martin Grohe. Walking out of the weisfeiler leman  
831 hierarchy: Graph learning beyond message passing. *Transactions on Machine Learning Research*,  
832 2023. (pages 1, 2, 4, 7, 8, 31, 32, 46)  
833
- 834 Jake Topping, Francesco Di Giovanni, Benjamin Paul Chamberlain, Xiaowen Dong, and Michael M.  
835 Bronstein. Understanding over-squashing and bottlenecks on graphs via curvature. In *ICLR*, 2022.  
836 (pages 6, 7, 47)
- 837 Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. Can  
838 language models solve graph problems in natural language? In *NeurIPS*, 2023. (page 47)  
839
- 840 Yuanqing Wang and Kyunghyun Cho. Non-convolutional graph neural networks. *arXiv*, 2024. (page  
841 7)
- 842 Gail Weiss, Yoav Goldberg, and Eran Yahav. Thinking like transformers. In Marina Meila and Tong  
843 Zhang (eds.), *ICML*, 2021. (page 47)  
844
- 845 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,  
846 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. Huggingface's  
847 transformers: State-of-the-art natural language processing. *arXiv*, 2019.
- 848 Xinyi Wu, Amir Ajourlou, Zihui Wu, and Ali Jadbabaie. Demystifying oversmoothing in attention-  
849 based graph neural networks. In *NeurIPS*, 2023. (page 47)  
850
- 851 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural  
852 networks? In *ICLR*, 2019. (page 1)
- 853 Yujun Yan, Milad Hashemi, Kevin Swersky, Yaoqing Yang, and Danai Koutra. Two sides of the  
854 same coin: Heterophily and oversmoothing in graph convolutional neural networks. In *2022 IEEE*  
855 *International Conference on Data Mining (ICDM)*, 2022. (page 33)  
856
- 857 Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. Language is all a graph  
858 needs. In *Findings of EACL*, 2024. (page 47)
- 859 Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank J. Reddi, and Sanjiv Kumar. Are  
860 transformers universal approximators of sequence-to-sequence functions? In *ICLR*, 2020. (pages  
861 4, 5)  
862
- 863 Bohang Zhang, Jingchu Gai, Yiheng Du, Qiwei Ye, Di He, and Liwei Wang. Beyond weisfeiler-  
lehman: A quantitative framework for gnn expressiveness. *arXiv*, 2024. (page 33)



864 Jianan Zhao, Meng Qu, Chaozhuo Li, Hao Yan, Qian Liu, Rui Li, Xing Xie, and Jian Tang. Learning  
865 on large-scale text-attributed graphs via variational inference. *arXiv*, 2022a.  
866

867 Jianan Zhao, Le Zhuo, Yikang Shen, Meng Qu, Kai Liu, Michael M. Bronstein, Zhaocheng Zhu, and  
868 Jian Tang. Graphtext: Graph reasoning in text space. *arXiv*, 2023. (pages 32, 47)

869 Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. *arXiv*, 2019. (page  
870 6)  
871

872 Lingxiao Zhao, Wei Jin, Leman Akoglu, and Neil Shah. From stars to subgraphs: Uplifting any GNN  
873 with local structure awareness. In *ICLR*, 2022b. (pages 9, 33)

874 Xiaojin Zhu and Zoubin Ghahramani. Learning from labeled and unlabeled data with label propaga-  
875 tion. *ProQuest number: information to all users*, 2002. (pages 10, 47)  
876

877 Xiaojin Jerry Zhu. Semi-supervised learning literature survey, 2005. (pages 10, 47)  
878

879 David Zuckerman. On the time to traverse all edges of a graph. *Inf. Process. Lett.*, 1991. (pages 4, 6,  
880 39, 43, 46)  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910  
911  
912  
913  
914  
915  
916  
917

## A APPENDIX

### A.1 SECOND-ORDER RANDOM WALKS (SECTION 2)

In second-order random walks, the probability of choosing  $v_{t+1}$  depends not only on  $v_t$  but also on  $v_{t-1}$ . We consider non-backtracking (Alon et al., 2007; Fitzner & van der Hofstad, 2013) and node2vec (Grover & Leskovec, 2016) random walks as representatives.

A non-backtracking random walk is defined upon a first-order random walk algorithm, e.g. one in Equation 4, by enforcing  $v_{t+1} \neq v_{t-1}$ :

$$\text{Prob}[v_{t+1} = x | v_t = j, v_{t-1} = i] := \begin{cases} \frac{\text{Prob}[v_{t+1} = x | v_t = j]}{\sum_{y \in N(j) \setminus \{i\}} \text{Prob}[v_{t+1} = y | v_t = j]} & \text{for } x \neq i, \\ 0 & \text{for } x = i, \end{cases} \quad (13)$$

where  $\text{Prob}[v_{t+1} = x | v_t = j]$  is the probability of walking  $j \rightarrow x$  given by the underlying first-order random walk. Notice that the probabilities are renormalized over  $N(j) \setminus \{i\}$ . This is ill-defined in the case the walk traverses  $i \rightarrow j$  and reaches a dangling vertex  $j$  which has  $i$  as its only neighbor, since  $N(j) \setminus \{i\} = \emptyset$ . In such cases, we allow the random walk to "begrudgingly" backtrack (Rappaport et al., 2017),  $i \rightarrow j$  and then  $j \rightarrow i$ , given that it is the only possible choice due to the dangling of  $j$ .

In case of node2vec random walk (Grover & Leskovec, 2016), a weighting term  $\alpha(v_{t-1}, v_{t+1})$  with two (hyper)parameters, return  $p$  and in-out  $q$ , is introduced to modify the behavior of a first-order random walk:

$$\text{Prob}[v_{t+1} = x | v_t = j, v_{t-1} = i] := \frac{\alpha(i, x) \text{Prob}[v_{t+1} = x | v_t = j]}{\sum_{y \in N(j)} \alpha(i, y) \text{Prob}[v_{t+1} = y | v_t = j]}, \quad (14)$$

where the probability  $\text{Prob}[v_{t+1} = x | v_t = j]$  is given by the underlying first-order random walk and  $\alpha(v_{t-1}, v_{t+1})$  is defined as follows using the shortest path distance  $d(v_{t-1}, v_{t+1})$ :

$$\alpha(v_{t-1}, v_{t+1}) := \begin{cases} 1/p & \text{for } d(v_{t-1}, v_{t+1}) = 0, \\ 1 & \text{for } d(v_{t-1}, v_{t+1}) = 1, \\ 1/q & \text{for } d(v_{t-1}, v_{t+1}) = 2. \end{cases} \quad (15)$$

Choosing a large return parameter  $p$  reduces backtracking since it decreases the probability of walking from  $v_t$  to  $v_{t-1}$ . Choosing a small in-out parameter  $q$  has a similar effect of avoiding  $v_{t-1}$ , with a slight difference that it avoids the neighbors of  $v_{t-1}$  as well.

We now show an extension of Proposition 2.2 to the above second-order random walks:

**Proposition A.1.** *The non-backtracking random walk in Equation 13 and the node2vec random walk in Equation 14 are invariant if their underlying first-order random walk algorithm is invariant.*

*Proof.* The proof is given in Appendix A.5.3. □

## A.2 MAIN ALGORITHM

We outline the main algorithms for the random walk and the recording function described in Section 2 and used in Section 5. Algorithm 1 shows the random walk algorithm, and Algorithms 2, 3, and 4 show the recording functions. For the random walk algorithm, we use non-backtracking described in Appendix A.1 and the minimum degree local rule conductance  $c_G(\cdot)$  given in Equation 6 by default.

Based on the algorithms, in Figures 4, 5, and 6 we illustrate the task formats used for graph separation experiments in Section 5.2 by providing examples of input graphs, text records of random walks on them, and prediction targets for the language model that processes the records. Likewise, in Figures 7, 8, and 9 we show task formats for transductive classification on arXiv citation network in Section 5.3.

**Algorithm 1:** Random walk algorithm

---

```

983 Algorithm 1: Random walk algorithm
984 Data: Input graph  $G$ , optional input vertex  $v$ , conductance function  $c_G : E(G) \rightarrow \mathbb{R}_+$ , walk
985     length  $l$ , optional restart probability  $\alpha \in (0, 1)$  or period  $k > 1$ 
986 Result: Random walk  $v_0 \rightarrow \dots \rightarrow v_l$ , restart flags  $r_1, \dots, r_l$ 
987 /* transition probabilities  $p : E(G) \rightarrow \mathbb{R}_+$  */
988 1 for  $(u, x) \in E(G)$  do
989 2    $p(u, x) \leftarrow c_G(u, x) / \sum_{y \in N(u)} c_G(u, y)$  // Equation 4
990 /* starting vertex  $v_0$  */
991 3 if  $v$  is given then
992 4   /* query starting vertex */
993 5    $v_0 \leftarrow v$ 
994 else
995 6   /* sample starting vertex */
996 7    $v_0 \sim \text{Uniform}(V(G))$ 
997 /* random walk  $v_1 \rightarrow \dots \rightarrow v_l$ , restart flags  $r_1, \dots, r_l$  */
998 8  $v_1 \sim \text{Categorical}(\{p(v_0, x) : x \in N(v_0)\})$ 
999 9  $r_1 \leftarrow 0$ 
1000 10 for  $t \leftarrow 2$  to  $l$  do
1001 11   /* restart flag  $r_t$  */
1002 12    $r_t \leftarrow 0$ 
1003 13   if  $\alpha$  is given then
1004 14     /* if restarted at  $t-1$ , do not restart */
1005 15     if  $r_{t-1} = 0$  then
1006 16     17      $r_t \sim \text{Bernoulli}(\alpha)$ 
1007 18   else if  $k$  is given then
1008 19     if  $t \equiv 0 \pmod{k}$  then
1009 20     21      $r_t \leftarrow 1$ 
1010 22   /* random walk  $v_t$  */
1011 23   if  $r_t = 0$  then
1012 24      $S \leftarrow N(v_{t-1}) \setminus \{v_{t-2}\}$ 
1013 25     if  $S \neq \emptyset$  then
1014 26     27     /* non-backtracking */
1015 28     29      $v_t \sim \text{Categorical}(\{p(v_{t-1}, x) / \sum_{y \in S} p(v_{t-1}, y) : x \in S\})$  // Equation 13
1016 30     else
1017 31     32     /* begrudgingly backtracking */
1018 33     34      $v_t \leftarrow v_{t-2}$ 
1019 35   else
1020 36     /* restart */
1021 37      $v_t \leftarrow v_0$ 

```

---

1022  
1023  
1024  
1025

---

```

1026
1027 Algorithm 2: Recording function, using anonymization
1028 Data: Random walk  $v_0 \rightarrow \dots \rightarrow v_l$ , restart flags  $r_1, \dots, r_l$ 
1029 Result: Text record  $\mathbf{z}$ 
1030 /* named vertices  $S$ , namespace  $\text{id}(\cdot)$  */
1031 1  $S \leftarrow \{v_0\}$ 
1032 2  $\text{id}(v_0) \leftarrow 1$ 
1033 /* text record  $\mathbf{z}$  */
1034 3  $\mathbf{z} \leftarrow \text{id}(v_0)$ 
1035 4 for  $t \leftarrow 1$  to  $l$  do
1036     /* anonymization  $v_t \mapsto \text{id}(v_t)$  */
1037     5 if  $v_t \notin S$  then
1038         6      $S \leftarrow S \cup \{v_t\}$ 
1039         7      $\text{id}(v_t) \leftarrow |S|$ 
1040     /* text record  $\mathbf{z}$  */
1041     8 if  $r_t = 0$  then
1042         /* record walk  $v_{t-1} \rightarrow v_t$  */
1043         9      $\mathbf{z} \leftarrow \mathbf{z} + "-" + \text{id}(v_t)$ 
1044     else
1045         /* record restart  $v_t = v_0$  */
1046         10      $\mathbf{z} \leftarrow \mathbf{z} + ";" + \text{id}(v_t)$ 

```

---

```

1048 Algorithm 3: Recording function, using anonymization and named neighbors
1049 Data: Random walk  $v_0 \rightarrow \dots \rightarrow v_l$ , restart flags  $r_1, \dots, r_l$ , input graph  $G$ 
1050 Result: Text record  $\mathbf{z}$ 
1051 /* named vertices  $S$ , namespace  $\text{id}(\cdot)$ , recorded edges  $T$  */
1052 1  $S \leftarrow \{v_0\}$ 
1053 2  $\text{id}(v_0) \leftarrow 1$ 
1054 3  $T \leftarrow \emptyset$ 
1055 /* text record  $\mathbf{z}$  */
1056 4  $\mathbf{z} \leftarrow \text{id}(v_0)$ 
1057 5 for  $t \leftarrow 1$  to  $l$  do
1058     /* anonymization  $v_t \mapsto \text{id}(v_t)$  */
1059     6 if  $v_t \notin S$  then
1060         7      $S \leftarrow S \cup \{v_t\}$ 
1061         8      $\text{id}(v_t) \leftarrow |S|$ 
1062     /* text record  $\mathbf{z}$  */
1063     9 if  $r_t = 0$  then
1064         /* record walk  $v_{t-1} \rightarrow v_t$  */
1065         10      $\mathbf{z} \leftarrow \mathbf{z} + "-" + \text{id}(v_t)$ 
1066         11      $T \leftarrow T \cup \{(v_{t-1}, v_t), (v_t, v_{t-1})\}$ 
1067         /* named neighborhood  $U$  */
1068         12      $U \leftarrow N(v_t) \cap S$ 
1069         if  $U \neq \emptyset$  then
1070             /* sort in ascending order  $\text{id}(u_1) < \dots < \text{id}(u_{|U|})$  */
1071             13      $[u_1, \dots, u_{|U|}] \leftarrow \text{SortByKey}(U, \text{id})$ 
1072             for  $u \leftarrow u_1$  to  $u_{|U|}$  do
1073                 /* record named neighbors  $v_t \rightarrow u$  */
1074                 14     if  $(v_t, u) \notin T$  then
1075                     15      $\mathbf{z} \leftarrow \mathbf{z} + "#" + \text{id}(u)$ 
1076                     16      $T \leftarrow T \cup \{(v_t, u), (u, v_t)\}$ 
1077             else
1078                 /* record restart  $v_t = v_0$  */
1079                 17      $\mathbf{z} \leftarrow \mathbf{z} + ";" + \text{id}(v_t)$ 

```

---

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

---

**Algorithm 4:** Recording function for transductive classification on arXiv network (Section 5.3)

---

**Data:** Directed input graph  $G$ , random walk  $v_0 \rightarrow \dots \rightarrow v_l$  and restart flags  $r_1, \dots, r_l$  on the undirected copy of  $G$ , paper titles  $\{\mathbf{t}_v\}_{v \in V(G)}$  and abstracts  $\{\mathbf{a}_v\}_{v \in V(G)}$ , target category labels  $\{\mathbf{y}_v\}_{v \in L}$  for labeled vertices  $L \subset V(G)$

**Result:** Text record  $\mathbf{z}$

```

/* named vertices  $S$ , namespace  $\text{id}(\cdot)$ , recorded edges  $T$  */
1  $S \leftarrow \{v_0\}$ 
2  $\text{id}(v_0) \leftarrow 1$ 
3  $T \leftarrow \emptyset$ 
/* text record  $\mathbf{z}$  */
/* record starting vertex */
4  $\mathbf{z} \leftarrow \text{"Paper 1 - Title: } \{\mathbf{t}_{v_0}\}\text{"}, \text{Abstract: } \{\mathbf{a}_{v_0}\}\text{"}$ 
5 for  $t \leftarrow 1$  to  $l$  do
    /* anonymization  $v_t \mapsto \text{id}(v_t)$  */
    6 if  $v_t \notin S$  then
    7      $S \leftarrow S \cup \{v_t\}$ 
    8      $\text{id}(v_t) \leftarrow |S|$ 
    /* text record  $\mathbf{z}$  */
    9 if  $r_t = 0$  then
    10     /* record walk  $v_{t-1} \rightarrow v_t$  with direction */
    11     if  $(v_{t-1}, v_t) \in E(G)$  then
    12          $\mathbf{z} \leftarrow \mathbf{z} + \text{" Paper } \{\text{id}(v_{t-1})\} \text{ cites Paper } \{\text{id}(v_t)\}\text{"}$ 
    13     else
    14          $\mathbf{z} \leftarrow \mathbf{z} + \text{" Paper } \{\text{id}(v_{t-1})\} \text{ is cited by Paper } \{\text{id}(v_t)\}\text{"}$ 
    15      $T \leftarrow T \cup \{(v_{t-1}, v_t), (v_t, v_{t-1})\}$ 
    16     /* record title  $\mathbf{t}_v$ , abstract  $\mathbf{a}_v$ , and label  $\mathbf{y}_v$  if  $v$  is labeled */
    17     if  $\text{id}(v_t) = |S|$  then
    18          $\mathbf{z} \leftarrow \mathbf{z} + \text{" - } \{\mathbf{t}_v\}\text{"}$ 
    19         if  $v \in L$  then
    20              $\mathbf{z} \leftarrow \mathbf{z} + \text{" , Category: } \{\mathbf{y}_v\}\text{"}$ 
    21          $\mathbf{z} \leftarrow \mathbf{z} + \text{" , Abstract: } \{\mathbf{a}_v\}\text{"}$ 
    22     else
    23          $\mathbf{z} \leftarrow \mathbf{z} + \text{" . }\text{"}$ 
    24     /* named neighborhood  $U$  */
    25      $U \leftarrow N(v_t) \cap S$ 
    26     if  $U \neq \emptyset$  then
    27         /* sort in ascending order  $\text{id}(u_1) < \dots < \text{id}(u_{|U|})$  */
    28          $[u_1, \dots, u_{|U|}] \leftarrow \text{SortByKey}(U, \text{id})$ 
    29         for  $u \leftarrow u_1$  to  $u_{|U|}$  do
    30             /* record named neighbors  $v_t \rightarrow u$  with directions */
    31             if  $(v_t, u) \notin T$  then
    32                 if  $(v_t, u) \in E(G)$  then
    33                      $\mathbf{z} \leftarrow \mathbf{z} + \text{" Paper } \{\text{id}(v_t)\} \text{ cites Paper } \{\text{id}(u)\}\text{"}$ 
    34                 else
    35                      $\mathbf{z} \leftarrow \mathbf{z} + \text{" Paper } \{\text{id}(v_t)\} \text{ is cited by Paper } \{\text{id}(u)\}\text{"}$ 
    36                  $T \leftarrow T \cup \{(v_t, u), (u, v_t)\}$ 
    37     else
    38         /* record restart */
    39          $\mathbf{z} \leftarrow \mathbf{z} + \text{" Restart at Paper 1. }\text{"}$ 

```

---

1134  
1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

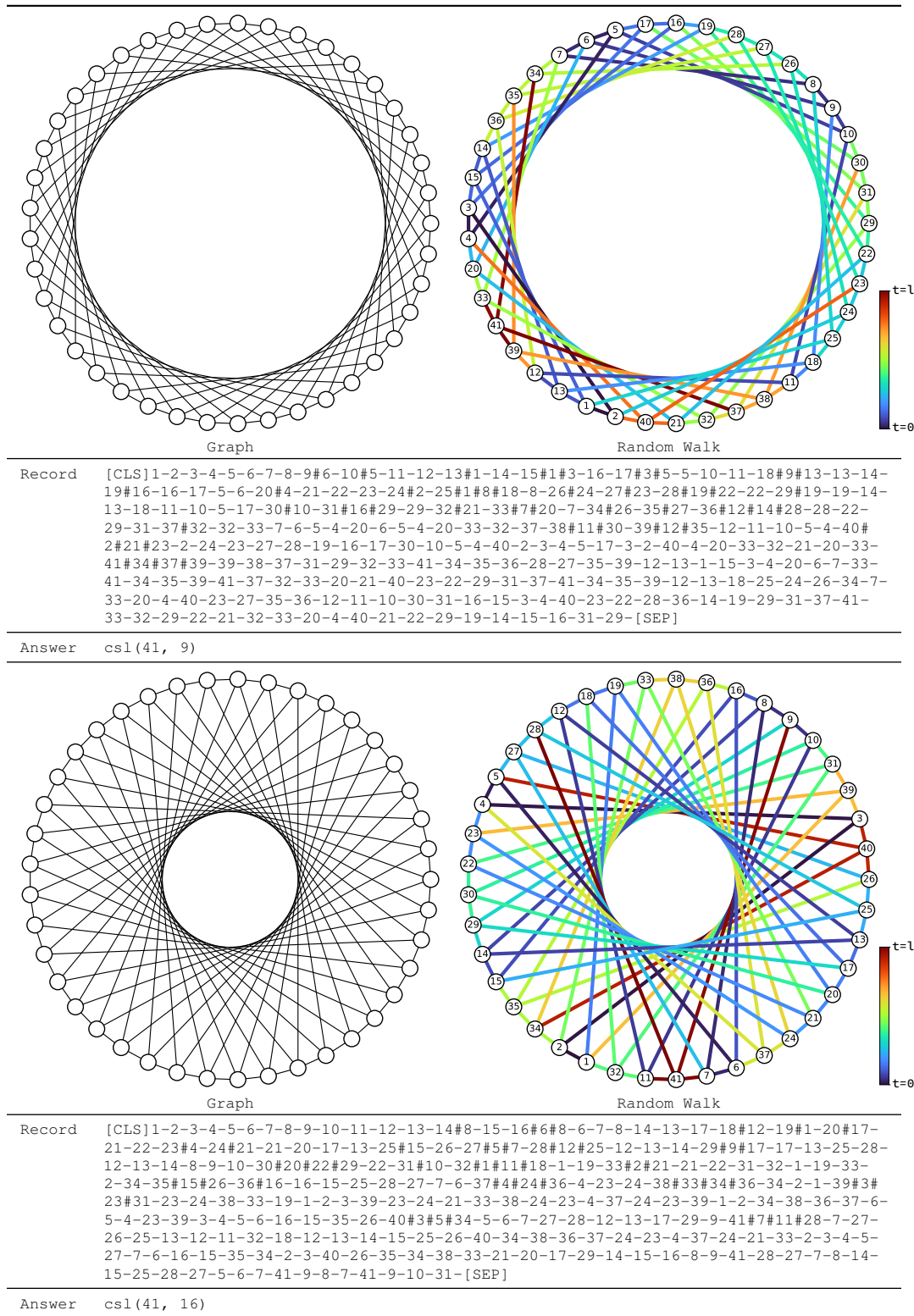


Figure 4: Two CSL graphs and text records of random walks from Algorithms 1 and 3. Task is graph classification into 10 isomorphism types  $csl(41, s)$  for skip length  $s \in \{2, 3, 4, 5, 6, 9, 11, 12, 13, 16\}$ . In random walks, we label vertices by anonymization and color edges by their time of discovery.

1188  
1189  
1190  
1191  
1192  
1193  
1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241

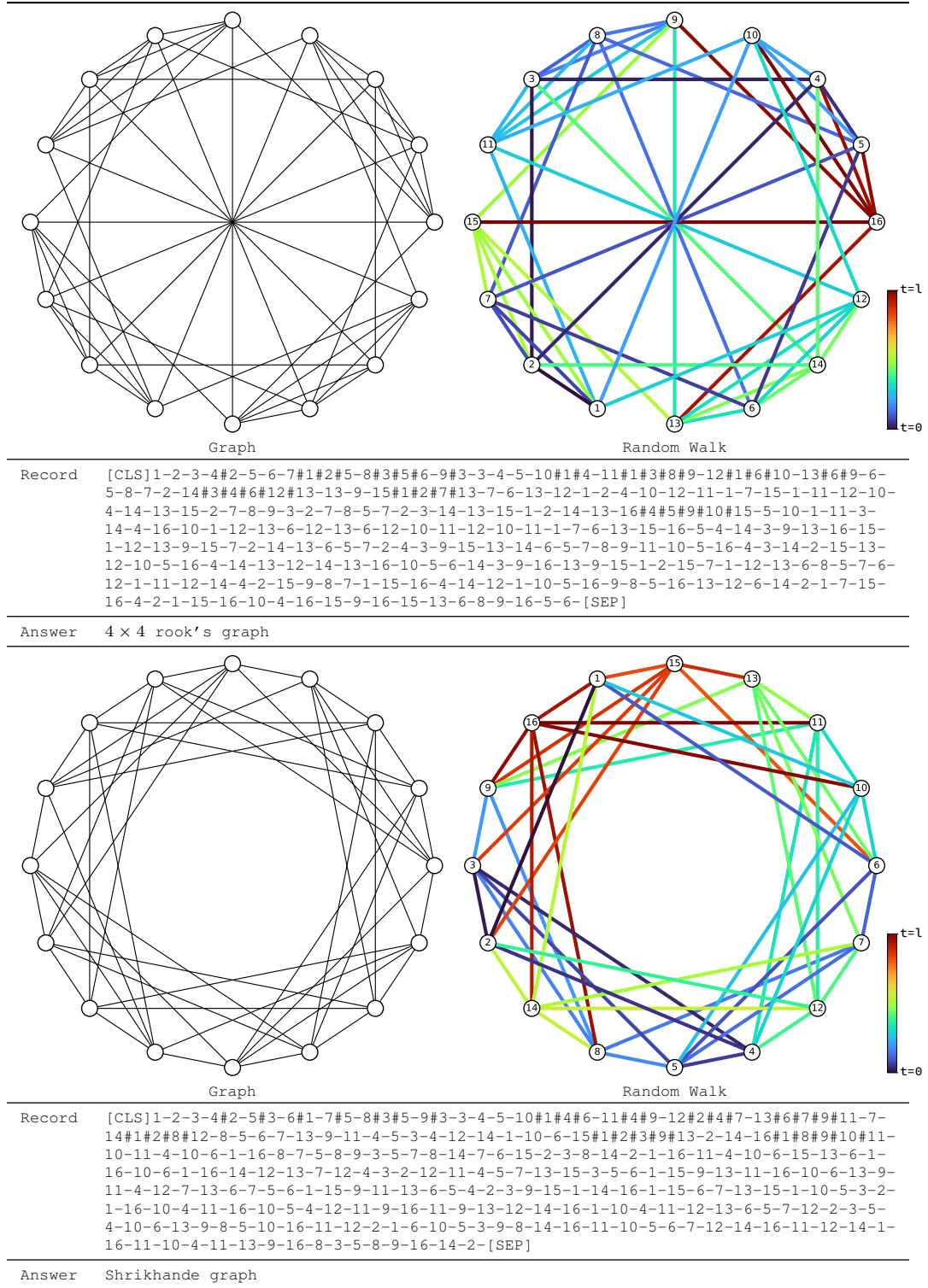


Figure 5: SR16 graphs and text records of random walks from Algorithms 1 and 3. The task is graph classification into two isomorphism types  $\{4 \times 4$  rook's graph, Shrikhande graph $\}$ . In random walks, we label vertices by anonymization and color edges by time of discovery.



1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

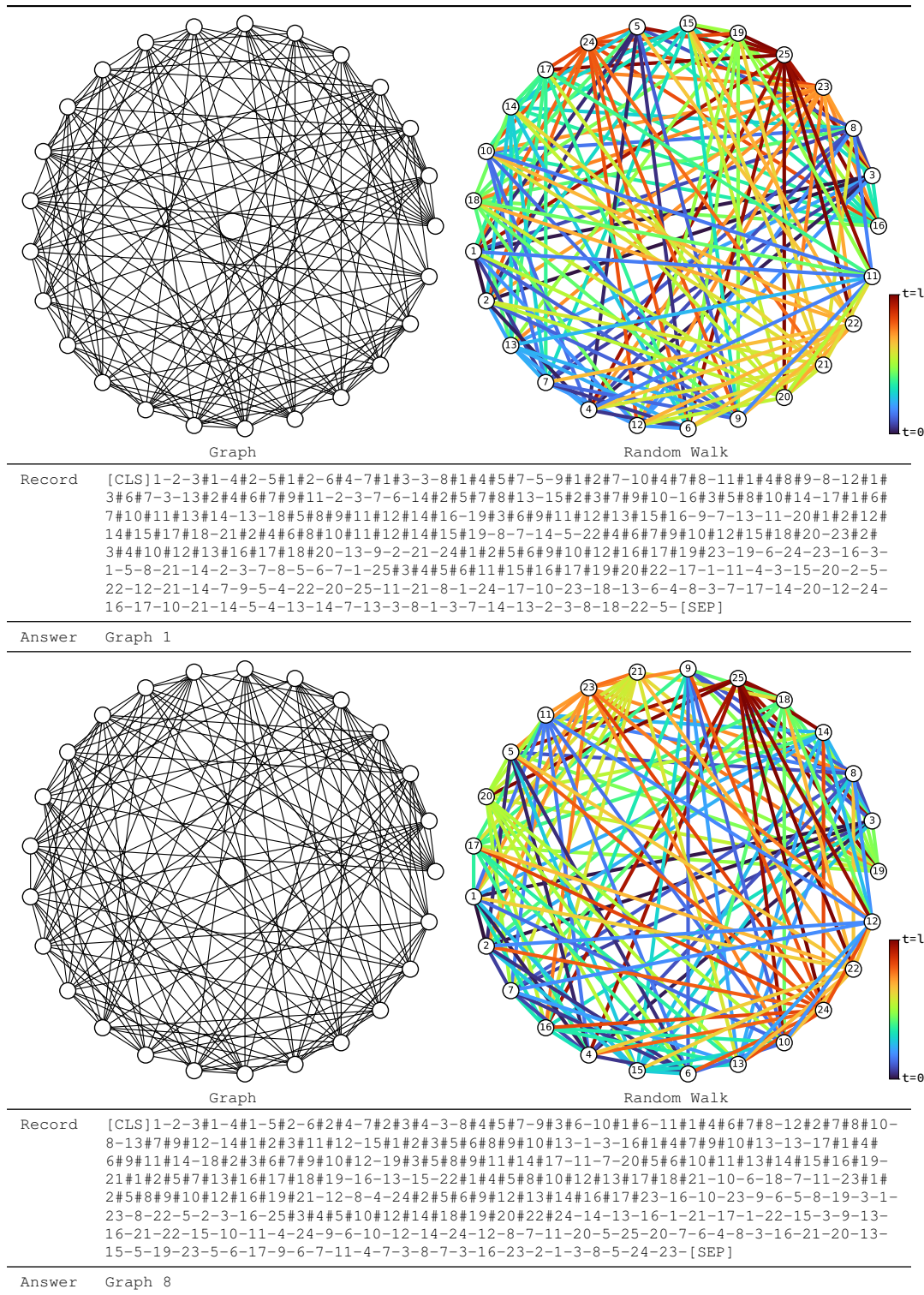


Figure 6: Two SR25 graphs and text records of random walks from Algorithms 1 and 3. The task is graph classification into 15 isomorphism types. In random walks, we label vertices by anonymization and color edges by their time of discovery.



1296	
1297	
1298	
1299	
1300	
1301	
1302	
1303	
1304	System
1305	A walk on the arXiv citation network will be given. Predict the arXiv CS sub-category Paper 1 belongs to. It is one of the following: Numerical Analysis (cs.NA), Multimedia (cs.MM), Logic in Computer Science (cs.LO), ... Discrete Mathematics (cs.DM). Only respond with the answer, do not say any word or explain.
1307	Record
1308	A walk on arXiv citation network is as follows:
1309	Paper 1 - Title: Safety guided deep reinforcement learning via online gaussian process estimation, Abstract: An important facet of reinforcement learning (rl) has to do with how the agent goes about exploring the environment. traditional exploration strategies typically focus on efficiency and ignore safety. however, for practical applications, ensuring safety of the agent during exploration is crucial since performing an unsafe action or reaching an unsafe state could result in irreversible damage to the agent. the main challenge of safe exploration is that characterizing the unsafe states and actions... Restart at 1. Paper 1 is cited by 2 - Learning to walk in the real world with minimal human effort, Abstract: Reliable and stable locomotion has been one of the most fundamental challenges for legged robots. deep reinforcement learning (deep rl) has emerged as a promising method for developing such control policies... Restart at 1. Paper 1 cites 3 - Deep q learning from demonstrations, Category: Artificial Intelligence (cs.AI), Abstract: Deep reinforcement learning (rl) has achieved several high profile successes in difficult decision-making problems. however, these algorithms typically require a huge amount of data before they reach... Restart at 1. Paper 1 cites 4 - Constrained policy optimization, Category: Machine Learning (cs.LG), Abstract: For many applications of reinforcement learning it can be more convenient to specify both a reward function and constraints, rather than trying to design behavior through the reward function. for example,... Paper 4 is cited by 2. Paper 4 is cited by 5 - Artificial intelligence values and alignment, Abstract: This paper looks at philosophical questions that arise in the context of ai alignment. it defends three propositions. first, normative and technical aspects of the ai alignment problem are interrelated,... Paper 5 cites 6 - Agi safety literature review, Category: Artificial Intelligence (cs.AI), Abstract: The development of artificial general intelligence (agi) promises to be a major event. along with its many potential benefits, it also raises serious safety concerns (bostrom, 2014). the intention of... Paper 6 is cited by 7 - Modeling agi safety frameworks with causal influence diagrams, Abstract: Proposals for safe agi systems are typically made at the level of frameworks, specifying how the components of the proposed system should be trained and interact with each other. in this paper, we model... Restart at 1. Paper 1 cites 8 - Safe learning of regions of attraction for uncertain nonlinear systems with gaussian processes, Category: Systems and Control (cs.SY), Abstract: Control theory can provide useful insights into the properties of controlled, dynamic systems. one important property of nonlinear systems is the region of attraction (roa), a safe subset of the state... Paper 8 is cited by 9 - Control theory meets pomdps a hybrid systems approach, Abstract: Partially observable markov decision processes (pomdps) provide a modeling framework for a variety of sequential decision making under uncertainty scenarios in artificial intelligence (ai). since the... Restart at 1.
1310	
1311	
1312	
1313	
1314	
1315	
1316	
1317	
1318	
1319	
1320	
1321	
1322	
1323	
1324	
1325	
1326	
1327	
1328	
1329	
1330	
1331	
1332	
1333	
1334	
1335	
1336	...
1337	Which arXiv CS sub-category does Paper 1 belong to? Only respond with the answer, do not say any word or explain.
1338	
1339	Answer
1340	Machine Learning (cs.LG)

Figure 7: Transductive classification on arXiv citation network using text record of random walk from Algorithms 1 and 4. Colors indicate task instruction, walk information, and label information.

1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

1350	System	Title and abstract of an arXiv paper will be given. Predict the arXiv CS sub-category the paper belongs to. It is one of the following: Numerical Analysis (cs.NA), Multimedia (cs.MM), Logic in Computer Science (cs.LO), ... Discrete Mathematics (cs.DM). Only respond with the answer, do not say any word or explain.
1351		
1352		
1353	Context	Title: Safety guided deep reinforcement learning via online gaussian process estimation
1354		Abstract: An important facet of reinforcement learning (rl) has to do with how the agent goes about exploring the environment. traditional exploration strategies typically focus on efficiency and ignore safety. however, for practical applications, ensuring safety of the agent during exploration is crucial since performing an unsafe action or reaching an unsafe state could result in irreversible damage to the agent. the main challenge of safe exploration is that characterizing the unsafe states and actions...
1355		
1356		
1357		
1358		
1359		Which arXiv CS sub-category does this paper belong to?
1360		
1361	Answer	Machine Learning (cs.LG)

Figure 8: Zero-shot format for vertex classification on arXiv citation network. Task instruction and label information are colored.

1366	System	Title and abstract of an arXiv paper will be given. Predict the arXiv CS sub-category the paper belongs to. It is one of the following: Numerical Analysis (cs.NA), Multimedia (cs.MM), Logic in Computer Science (cs.LO), ... Discrete Mathematics (cs.DM).
1367		
1368		
1369	Context	Title: Deeptrack learning discriminative feature representations online for robust visual tracking
1370		Abstract: Deep neural networks, albeit their great success on feature learning in various computer vision tasks, are usually considered as impractical for online visual tracking, because they require very long...
1371		Category: cs.CV
1372		
1373		
1374		Title: Perceived audiovisual quality modelling based on decision trees genetic programming and neural networks...
1375		Abstract: Our objective is to build machine learning based models that predict audiovisual quality directly from a set of correlated parameters that are extracted from a target quality dataset. we have used the...
1376		Category: cs.MM
1377		
1378		...
1379		
1380		Title: Safety guided deep reinforcement learning via online gaussian process estimation
1381		Abstract: An important facet of reinforcement learning (rl) has to do with how the agent goes about exploring the environment. traditional exploration strategies typically focus on efficiency and ignore safety. however, for practical applications, ensuring safety of the agent during exploration is crucial since performing an unsafe action or reaching an unsafe state could result in irreversible damage to the agent. the main challenge of safe exploration is that characterizing the unsafe states and actions...
1382		
1383		
1384		
1385		
1386		Which arXiv CS sub-category does this paper belong to?
1387	Answer	cs.LG

Figure 9: One-shot format for transductive classification on arXiv citation network. 40 labeled examples are given in form of multi-turn dialogue. Task instruction and label information are colored.

1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403

1404 [...] A walk on arXiv citation network is as follows:  
 1405 Paper 1 - Title: Unsupervised and interpretable scene discovery with discrete attend infer  
 1406 repeat, Abstract: In this work we present discrete attend infer repeat (discrete-air), a  
 1407 recurrent auto-encoder with structured latent distributions containing discrete categori-  
 1408 cal distributions, continuous attribute distributions, and factorised spatial attention.  
 1409 while inspired by the original air model and retaining air model’s capability in identify-  
 1410 ing objects in an image, discrete-air provides direct interpretability of the latent codes.  
 1411 we show that for multi-mnist and a multiple-objects version of dsprites... Restart at 1.  
 1412 [...] Restart at 1. Paper 1 cites 12 - Attend infer repeat fast scene understanding with  
 1413 generative models, Category: Computer Vision and Pattern Recognition (cs.CV), Abstract: We  
 1414 present a framework for efficient inference in structured image models that explicitly rea-  
 1415 son about objects. we achieve this by performing probabilistic inference using a recurrent  
 1416 neural network that... Paper 12 cites 3. Paper 12 cites 7. Paper 12 is cited by 9. Paper  
 1417 12 is cited by 10. Paper 12 cites 11. Restart at 1. Paper 1 cites 11. Restart at 1. Paper  
 1418 1 cites 12. Paper 12 is cited by 13 - An architecture for deep hierarchical generative mod-  
 1419 els, Category: Machine Learning (cs.LG), Abstract: We present an architecture which lets us  
 1420 train deep, directed generative models with many layers of latent variables. we include de-  
 1421 terministic paths between all latent variables and the generated output,... Paper 13 cites  
 1422 3. Restart at 1. [...]

---

Attention (Layer 13/30)

---

Answer	cs.CV
Prediction	cs.CV

---

1423 Figure 10: Example of attention in a frozen Llama 3 8B applied on arXiv transductive classification.  
 1424 Attention weights are averaged over all 30 heads.

### 1425 A.3 ATTENTION VISUALIZATIONS

1426 In Figure 10, we show an example of self-attention in the frozen Llama 3 8B model applied to arXiv  
 1427 transductive classification (Section 5.3). We show attention weights on text record of random walk  
 1428 from the generated `cs` token as query, which is just before finishing the prediction e.g. `cs . CV`. We  
 1429 color the strongest activation with orange, and do not color values below 1% of it. The model invests  
 1430 a nontrivial amount of attention on walk information such as `Paper 1 cites 12`, while also  
 1431 making use of titles and labels of labeled vertices. This indicates the model is utilizing the graph  
 1432 structure recorded by the random walk in conjunction with other information to make predictions.  
 1433

1434 In Figures 11, 12, and 13, we show examples of self-attention in DeBERTa models trained for graph  
 1435 separation (Section 5.2). We first show attention weights on text record of random walk from the  
 1436 [CLS] query token. Then, we show an alternative visualization where the attention weights are  
 1437 mapped onto the original input graph. For example, for each attention on  $v_{t+1}$  where the walk has  
 1438 traversed  $v_t \rightarrow v_{t+1}$ , we regard it as an attention on the edge  $(v_t, v_{t+1})$  in the input graph. We ignore  
 1439 [CLS] and [SEP] tokens since they do not appear on the input graph. We observe that the models  
 1440 often focus on sparse, connected substructures, which presumably provide discriminative information  
 1441 on the isomorphism types. In particular, for CSL graphs (Figure 11) we observe an interpretable  
 1442 pattern of approximate cycles composed of skip links. This is presumably related to measuring the  
 1443 lengths of skip links, which provides sufficient information of isomorphism types of CSL graphs.  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

1458  
1459  
1460  
1461  
1462  
1463  
1464  
1465  
1466  
1467  
1468  
1469  
1470  
1471  
1472  
1473  
1474  
1475  
1476  
1477  
1478  
1479  
1480  
1481  
1482  
1483  
1484  
1485  
1486  
1487  
1488  
1489  
1490  
1491  
1492  
1493  
1494  
1495  
1496  
1497  
1498  
1499  
1500  
1501  
1502  
1503  
1504  
1505  
1506  
1507  
1508  
1509  
1510  
1511

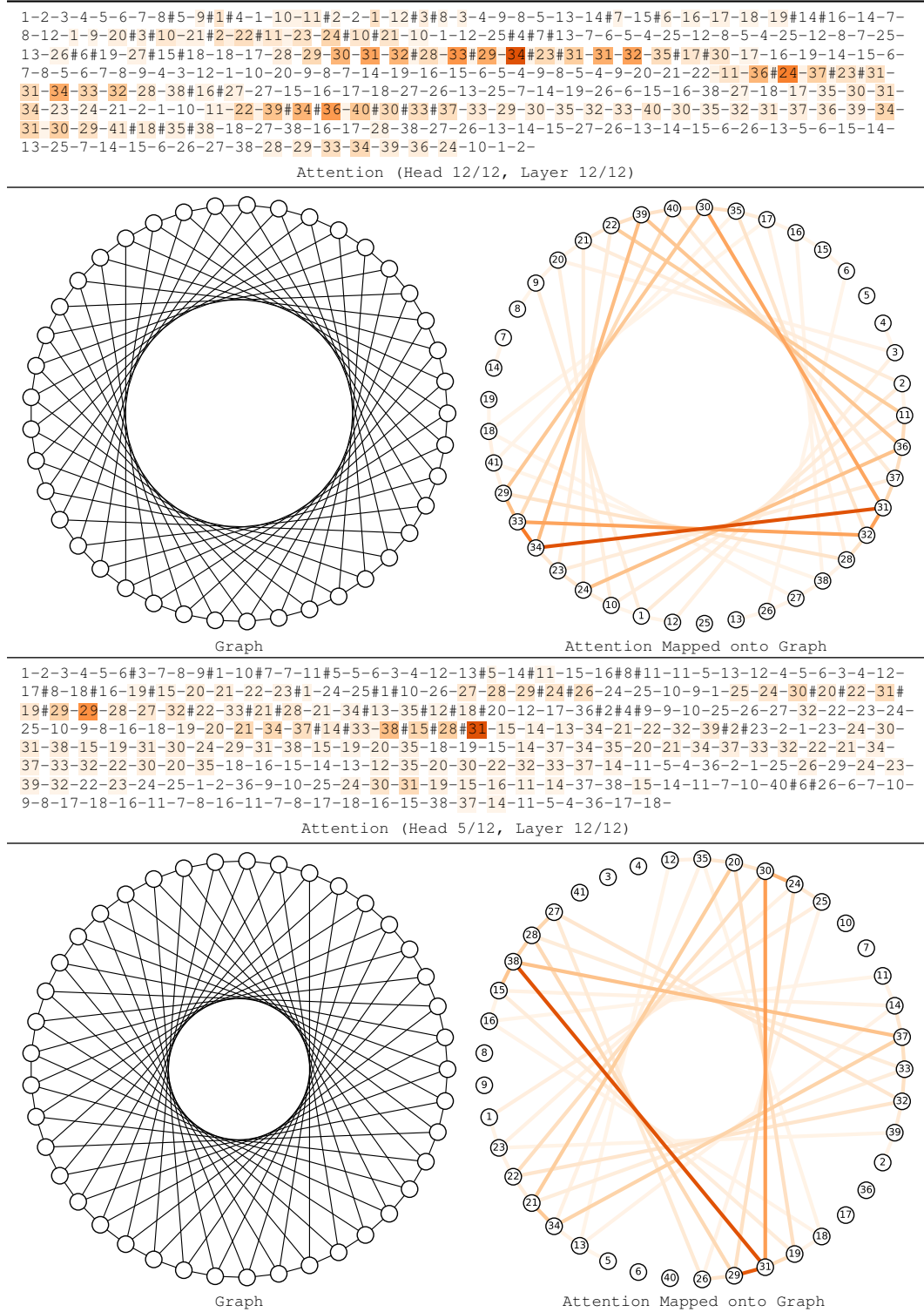
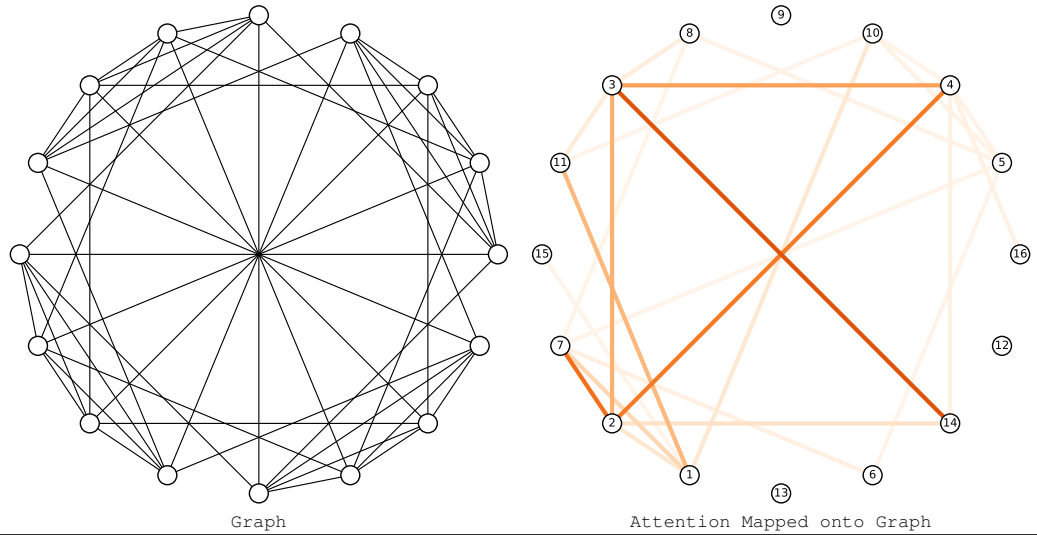


Figure 11: Examples of attention from [CLS] token in a DeBERTa trained on CSL graph separation, forming approximate cycles using skip links. This is presumably related to measuring the lengths of skip links, which provides sufficient information to infer isomorphism types of CSL graphs.

1512  
 1513  
 1514  
 1515  
 1516  
 1517  
 1518  
 1519  
 1520  
 1521  
 1522  
 1523  
 1524  
 1525  
 1526  
 1527  
 1528  
 1529  
 1530  
 1531  
 1532  
 1533  
 1534  
 1535  
 1536  
 1537  
 1538  
 1539  
 1540  
 1541  
 1542  
 1543  
 1544  
 1545  
 1546  
 1547  
 1548  
 1549  
 1550  
 1551  
 1552  
 1553  
 1554  
 1555  
 1556  
 1557  
 1558  
 1559  
 1560  
 1561  
 1562  
 1563  
 1564  
 1565

1-2-3-4#2-5-6-7#1#2#5-8#3#5#6-9#3-3-4-5-10#1#4-11#1#3#8#9-12#1#6#10-13#6#9-6-5-8-7-2-14#3#4#6#12#13-13-9-15#1#2#7#13-7-6-13-12-1-2-4-10-12-11-1-7-15-1-11-12-10-4-14-13-15-2-7-8-9-3-2-7-8-5-7-2-3-14-13-15-1-2-14-13-16#4#5#9#10#15-5-10-1-11-3-14-4-16-10-1-12-13-6-12-13-6-12-10-11-12-10-11-1-7-6-13-15-16-5-4-14-3-9-13-16-15-1-12-13-9-15-7-2-14-13-6-5-7-2-4-3-9-15-13-14-6-5-7-8-9-11-10-5-16-4-3-14-2-15-13-12-10-5-16-4-14-13-12-14-13-16-10-5-6-14-3-9-16-13-9-15-1-2-15-7-1-12-13-6-8-5-7-6-12-1-11-12-14-4-2-15-9-8-7-1-15-16-4-14-12-1-10-5-16-9-8-5-16-13-12-6-14-2-1-7-15-16-4-2-1-15-16-10-4-16-15-9-16-15-13-6-8-9-16-5-6-

Attention (Head 11/12, Layer 10/12)



1-2-3-4#2#5#3-6#1-7#5-8#3#5-9#3-3-4-5-10#1#4#6-11#4#9-12#2#4#7-13#6#7#9#11-7-14#1#2#8#12-8-5-6-7-13-9-11-4-5-3-4-12-14-1-10-6-15#1#2#3#9#13-2-14-16#1#8#9#10#11-10-11-4-10-6-1-16-8-7-5-8-9-3-5-7-8-14-7-6-15-2-3-8-14-2-1-16-11-4-10-6-15-13-6-1-16-10-6-1-16-14-12-13-7-12-4-3-2-12-11-4-5-7-13-15-3-5-6-1-15-9-13-11-16-10-6-13-9-11-4-12-7-13-6-7-5-6-1-15-9-11-13-6-5-4-2-3-9-15-1-14-16-1-15-6-7-13-15-1-10-5-3-2-1-16-10-4-11-16-10-5-4-12-11-9-16-11-9-13-12-14-16-1-10-4-11-12-13-6-5-7-12-2-3-5-4-10-6-13-9-8-5-10-16-11-12-2-1-6-10-5-3-9-8-14-16-11-10-5-6-7-12-14-16-11-12-14-1-16-11-10-4-11-13-9-16-8-3-5-8-9-16-14-2-

Attention (Head 11/12, Layer 7/12)

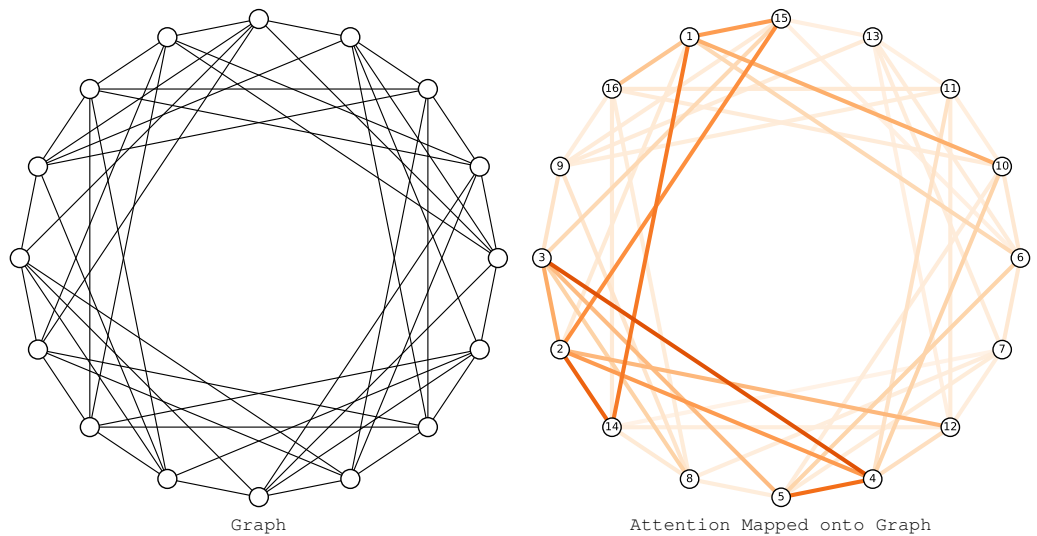


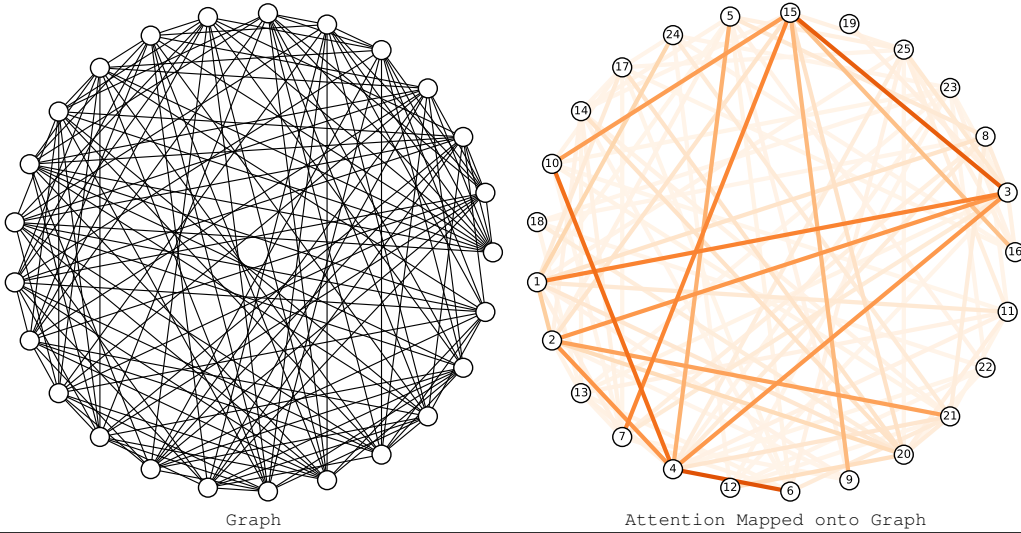
Figure 12: Examples of attention from [CLS] token in a DeBERTa trained on SR16 graph separation. The model tend to focus on neighborhood records that form a sparse connected substructure, which we conjecture to provide discriminative information on the isomorphism type of SR16 graphs.



1566  
1567  
1568  
1569  
1570  
1571  
1572  
1573  
1574  
1575  
1576  
1577  
1578  
1579  
1580  
1581  
1582  
1583  
1584  
1585  
1586  
1587  
1588  
1589  
1590  
1591  
1592  
1593  
1594  
1595  
1596  
1597  
1598  
1599  
1600  
1601  
1602  
1603  
1604  
1605  
1606  
1607  
1608  
1609  
1610  
1611  
1612  
1613  
1614  
1615  
1616  
1617  
1618  
1619

1-2-3#1-4#2-5#1#2-6#4-7#1#3-3-8#1#4#5#7-5-9#1#2#7-10#4#7#8-11#1#4#8#9-8-12#1#3#6#7-3-13#2#4#6#7#9#11-2-3-7-6-14#2#5#7#8#13-15#2#3#7#9#10-16#3#5#8#10#14-17#1#6#7#10#11#13#14-13-18#5#8#9#11#12#14#16-19#3#6#9#11#12#13#15#16-9-7-13-11-20#1#2#12#14#15#17#18-21#2#4#6#8#10#11#12#14#15#19-8-7-14-5-22#4#6#7#9#10#12#15#18#20-23#2#3#4#10#12#13#16#17#18#20-13-9-2-21-24#1#2#5#6#9#10#12#16#17#19#23-19-6-24-23-16-3-1-5-8-21-14-2-3-7-8-5-6-7-1-25#3#4#5#6#11#15#16#17#19#20#22-17-1-11-4-3-15-20-2-5-22-12-21-14-7-9-5-4-22-20-25-11-21-8-1-24-17-10-23-18-13-6-4-8-3-7-17-14-20-12-24-16-17-10-21-14-5-4-13-14-7-13-3-8-1-3-7-14-13-2-3-8-18-22-5-

Attention (Head 4/12, Layer 9/12)



1-2-3#1-4#1-5#2-6#2#4-7#2#3#4-8#4#5#7-9#3#6-10#1#6-11#1#4#6#7#8-12#2#7#8#10-8-13#7#9#12-14#1#2#3#11#12-15#1#2#3#5#6#8#9#10#13-1-3-16#1#4#7#9#10#13-13-17#1#4#6#9#11#14-18#2#3#6#7#9#10#12-19#3#5#8#9#11#14#17-11-7-20#5#6#10#11#13#14#15#16#19-21#1#2#5#7#13#16#17#18#19-16-13-15-22#1#4#5#8#10#12#13#17#18#21-10-6-18-7-11-23#1#2#5#8#9#10#12#16#19#21-12-8-4-24#2#5#6#9#12#13#14#16#17#23-16-10-23-9-6-5-8-19-3-1-23-8-22-5-2-3-16-25#3#4#5#10#12#14#18#19#20#22#24-14-13-16-1-21-17-1-22-15-3-9-13-16-21-22-15-10-11-4-24-9-6-10-12-14-24-12-8-7-11-20-5-25-20-7-6-4-8-3-16-21-20-13-15-5-19-23-5-6-17-9-6-7-11-4-7-3-8-7-3-16-23-2-1-3-8-5-24-23-

Attention (Head 10/12, Layer 8/12)

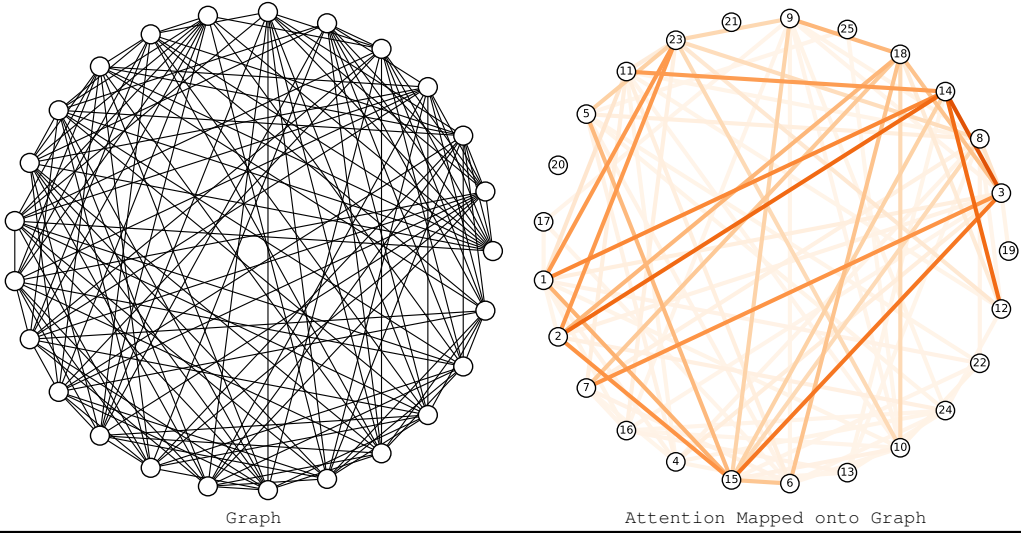


Figure 13: Examples of attention from [CLS] token in a DeBERTa trained on SR25 graph separation. The model tend to focus on neighborhood records that form a sparse connected substructure, which we conjecture to provide discriminative information on the isomorphism type of SR25 graphs.

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

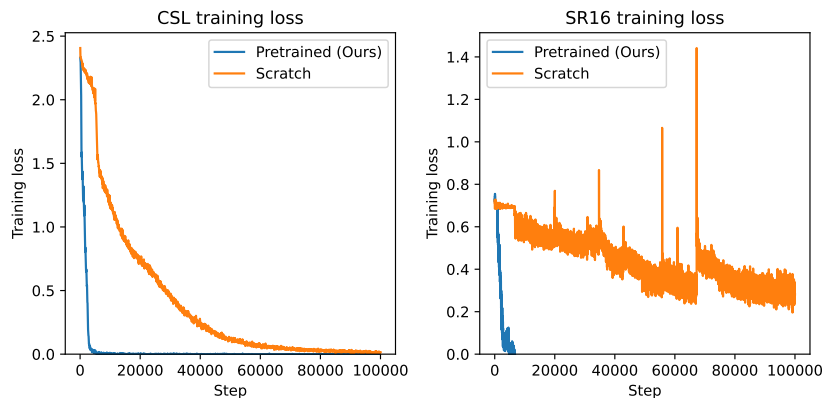


Figure 14: The impact of language pre-training for graph isomorphism learning (Section 5.2).

---

```
Record 1 [O, sp2]=π2 [C, d:3, sp2, r]-3 [C, cw, d:4, H:1, sp3, r]-4 [N, d:3, H:1, sp2, r]-
π5 [C, d:3, sp2, r]=π6 [O, sp2]=π5-π4-3-2=π1=π2-3-7 [C, d:4, H:2, sp3, r]-8 [S, d:2, sp3, r]-
9 [S, d:2, sp3, r]-10 [C, d:4, H:2, sp3, r]-11 [C, ccw, d:4, H:1, sp3, r]-12 [N, d:3, H:1, sp2, r]-
π13 [C, d:3, sp2, r]=π14 [O, sp2]=π13-π12-11-15 [C, d:3, sp2, r]-π16 [N, d:3, H:1, sp2, r]-
17 [C, cw, d:4, H:1, sp3, r]-18 [C, d:3, sp2, r]=π19 [O, sp2]=π18-17-20 [C, d:4, H:2, sp3]-
21 [C, d:4, H:2, sp3]-22 [C, d:3, sp2]=π23 [O, sp2]=π22-π24 [O, d:2, H:1, sp2]-π22=π23=
π22-π24-π22=π23=π22-21-20-17-16-π15=π25 [O, sp2]=π15-π16-17-20-21-22-π24-π22=
π23=π22-21-20-17-16-π15=π25=π15-π16-17-20-21-22-π24-π22=π23=π22-π24-π22-21-
20-17-18=π19=π18-π26 [N, d:3, H:1, sp2, r]-27 [C, ccw, d:4, H:1, sp3, r]-28 [C, d:3, sp2]=
π29 [O, sp2]=π28-π30 [N, d:3, H:1, sp2]-31 [C, ccw, d:4, H:1, sp3]-32 [C, d:4, H:1, sp3]-
33 [C, d:4, H:3, sp3]-32-34 [C, d:4, H:3, sp3]-32-33-32-34-32-33-32-34-32-33-31-30-
π28=π29=π28-π30-31-32-34-32-31-30-π28=π29=π28-π30-31-32-33-32-34-32-33-32-34-
32-31-30-π28=π29=π28-π30-31-35 [C, d:3, sp2]=π36 [O, sp2]=π35-π37 [N, d:3, H:1, sp2]-
38 [C, cw, d:4, H:1, sp3]-39 [C, d:3, sp2]-π40 [N, d:3, H:1, sp2]-41 [C, cw, d:4, H:1, sp3]-
42 [C, d:4, H:2, sp3]-43 [C, d:3, sp2, *, r]*π44 [C, d:3, H:1, sp2, *, r]*π45 [N, d:3, H:1, sp2, *,
r]*π46 [C, d:3, H:1, sp2, *, r]*π47 [N, d:2, sp2, *, r]**π43*π43*π44*π45*π46*π47*π43-42-41-
40-π39=π48 [O, sp2]=π39-38-37-π35=π36=π35-π37-38-49 [C, d:4, H:2, sp3]-50 [C
```

---

Figure 15: An example text record for Peptides-func graph classification.

Table 5: Peptides-func graph classification. The baseline scores are from Tönshoff et al. (2023).

Method	Test AP
GCN	0.5930
GINE	0.5498
GatedGCN	0.6069
Transformer	0.6326
SAN	0.6439
CRaWI	0.7074
<b>RWNN-DeBERTa (Ours)</b>	<b>0.7124 ± 0.0016</b>

#### A.4 SUPPLEMENTARY EXPERIMENTS

We include supplementary experiments we could not include in the main text due to space constraints.

##### A.4.1 THE BENEFIT OF LANGUAGE PRE-TRAINING (SECTION 5.2)

In Section 5.2, we have used language pre-trained DeBERTa for graph isomorphism learning, even though the records of random walks do not resemble natural language (Appendix A.2 and A.3). To test if language pre-training is useful, we additionally trained RWNN-DeBERTa on CSL and SR16 datasets using the same configurations to our reported models but from random initialization. The training curves are given in Figure 14. We find that pre-training has significant benefits over random initialization for isomorphism learning tasks, as randomly initialized models converge very slowly for CSL, and fails to converge to near-zero training loss for SR16.

Table 6: Fine-tuning for arXiv, extending Table 4. † denotes using validation labels as in Table 4.

Method	Training-free?	Accuracy
Llama3-8b zero-shot	○	50.20%
Llama3-8b one-shot	○	52.49%
Llama3-8b one-shot†	○	52.54%
Llama3-70b zero-shot	○	65.33%
Llama3-70b one-shot	○	67.57%
RWNN-Llama3-8b (Ours)	○	71.10%
RWNN-Llama3-8b (Ours)†	○	73.11%
RWNN-Llama3-70b (Ours)†	○	74.75%
RWNN-Llama3-8b (Ours), LoRA	×	<b>74.95%</b>

Table 7: Real-world transductive classification test accuracy on several homophilic (Cora, Citeseer) and heterophilic (Amazon Ratings) datasets. The baseline results were collected from (Sato, 2024; Zhao et al., 2023; Chen et al., 2024a; Liu et al., 2024; Chen et al., 2024b; Platonov et al., 2023). † denotes using validation labels as in Table 4.

Method	Training-free?	Cora 20-shot	Cora	Citeseer	Amazon Ratings
Training-free GNN	○	60.00%	-	-	-
GCN	×	81.40%	89.13%	74.92%	48.70%
GAT	×	80.80%	89.68%	75.39%	49.09%
GraphText	○	68.30%	67.77%	68.98%	-
LLaGA	○	-	59.59%	-	-
GraphText	×	-	87.11%	74.77%	-
LLaGA	×	-	88.86%	-	28.20%
OFA	×	75.90%	-	-	51.44%
RWNN-Llama3-8b (Ours)	○	72.05%	86.72%	78.37%	44.29%
RWNN-Llama3-8b (Ours)†	○	79.40%	88.01%	78.53%	48.83%
RWNN-Llama3-8b (Ours), LoRA	×	79.45%	86.72%	81.03%	53.90%

#### A.4.2 REAL-WORLD GRAPH CLASSIFICATION

We conduct a preliminary experiment on the Peptides-func dataset (Dwivedi et al., 2022) used in Tönshoff et al. (2023). Our model is a pre-trained DeBERTa-base (identical to Section 5.2) fine-tuned on text records of non-backtracking MDLR random walks with anonymization and neighborhood recording. The recording function is designed to properly incorporate the vertex and edge attributes provided in the dataset, including the atom and bond types. An example text is provided in Figure 15.

In Table 5, we report the test average precision (AP) at best validation accuracy, with 40 random predicted logits averaged at test time. We report the mean and standard deviation of the test performances for five repeated tests. The baseline scores, including CRaW1, are from Tönshoff et al. (2023). We see that RWNN-DeBERTa, which has been successful in graph isomorphism learning (Section 5.2), also shows a promising result in real-world protein graph classification, despite its simplicity of recording random walks in text and processing them with a fine-tuned language model.

#### A.4.3 FINE-TUNING FOR ARXIV TRANSDUCTIVE CLASSIFICATION (SECTION 5.3)

In Section 5.3, our RWNNs are training-free. We provide a preliminary result on fine-tuning on training labels using low-rank adaptation (LoRA) (Hu et al., 2022). The result is in Table 6, showing a promising performance.

#### A.4.4 ADDITIONAL REAL-WORLD TRANSDUCTIVE CLASSIFICATION

We provide additional experiments on real-world transductive classification datasets Cora, Citeseer, and Amazon Ratings. Our models are similar to ones in Section 5.3. We compile the baseline scores from (Zhao et al., 2023; Chen et al., 2024a; Liu et al., 2024; Chen et al., 2024b; Platonov et al., 2023), including MPNNs and language models on graphs. The results are in Table 7. Ours is competitive as a training-free method, and is reasonably good when fine-tuned with LoRA. Especially, on Amazon Ratings which is heterophilic (Platonov et al., 2023), our fine-tuned model outperforms all baselines.



Table 8: Substructure (8-cycle) counting. We report normalized test mean absolute error (MAE) at best validation. We trained MP, Subgraph GNN, Local 2-GNN, and Local 2-FGNN using configurations in Zhang et al. (2024), and trained MP-RNI, MP-Dropout, AgentNet, and CRaWI as in Table 3.

Method	Graph-level	Vertex-level
MP	.0917	.1156
Subgraph GNN	.0515	.0715
Local 2-GNN	.0433	.0478
Local 2-FGNN	.0422	.0438
Local 2-FGNN (large)	.0932	.1121
MP-RNI	.3610	.3650
MP-Dropout	.1359	.1393
AgentNet	.1107	N/A
CRaWI	.0526	.0725
RWNN-DeBERTa (Ours)	.0459	.0465

Table 9: Real-world link prediction. The baseline scores are from Table 5 of Cai et al. (2021), and we reproduced LGLP ourselves. We report the aggregated test AP for 3 randomized runs.

Method	YST	KHN	ADV
Katz	81.63 $\pm$ 0.41	83.04 $\pm$ 0.38	91.76 $\pm$ 0.15
PR	82.08 $\pm$ 0.46	87.18 $\pm$ 0.26	92.43 $\pm$ 0.17
SR	76.02 $\pm$ 0.49	75.87 $\pm$ 0.66	83.22 $\pm$ 0.20
node2vec	76.61 $\pm$ 0.94	80.60 $\pm$ 0.74	76.70 $\pm$ 0.82
SEAL	86.45 $\pm$ 0.25	90.37 $\pm$ 0.16	93.52 $\pm$ 0.13
LGLP	88.54 $\pm$ 0.77	90.80 $\pm$ 0.33	93.51 $\pm$ 0.32
RWNN-transformer (Ours)	88.18 $\pm$ 0.28	90.63 $\pm$ 0.44	93.52 $\pm$ 0.20

This supports our results in Section 3.2 that RWNNs avoid over-smoothing, as avoiding it is known to be important for handling heterophily (Yan et al., 2022).

#### A.4.5 SUBSTRUCTURE COUNTING

We test RWNN on additional challenging tasks that require both expressive power and generalization to unseen graphs. We use a synthetic dataset of 5,000 random regular graphs from Chen et al. (2020), and consider the task of counting substructures. We choose 8-cycles since they are known to require a particularly high expressive power (Zhang et al., 2024). We experiment with both graph- and vertex-level counting i.e. counting 8-cycles containing the queried vertex, following Zhang et al. (2024). We use the data splits from previous work (Chen et al., 2020; Zhao et al., 2022b; Zhang et al., 2024), while excluding disconnected graphs following our modeling assumption in Section 2.1. Our RWNN uses the same design to Section 5.2, and we train them with L1 loss for at most 250k steps using batch size of 128 graphs for graph-level counting and 8 graphs for vertex-level. At test-time, we ensemble 64 and 32 predictions by averaging for graph- and vertex-level tasks, respectively.

The results are in Table 8, and overall in line with Section 5.2. While MPNNs show a low performance, variants with higher expressive powers such as local 2-GNN and local 2-FGNN achieve high performances. This is consistent with the observation of Zhang et al. (2024). On the other hand, MPNNs with stochastic symmetry-breaking often show learning difficulties and does not achieve good performance. AgentNet was not directly applicable for vertex-level counting since its vertex encoding depends on which vertices the agent visits, and a learning difficulty was observed for graph-level counting. Our approach based on DeBERTa demonstrates competitive performance, outperforming random walk baselines and approaching performances of local 2-GNN and 2-FGNN.

#### A.4.6 LINK PREDICTION

We test an extension to real-world link prediction, based on Cai et al. (2021) where link prediction is cast as vertex classification on a line graph. We follow Cai et al. (2021) and only change their 3-layer 32-dim DGCNN with an RWNN that uses a 3-layer 32-dim transformer encoder. As the task is vertex-level, we use periodic restarts with  $k = 4$ . The results in Table 9 shows that RWNN performs reasonably well.

1782 A.5 PROOFS  
1783

1784 We recall that an isomorphism between two graphs  $G$  and  $H$  is a bijection  $\pi : V(G) \rightarrow V(H)$  such  
1785 that any two vertices  $u$  and  $v$  are adjacent in  $G$  if and only if  $\pi(u)$  and  $\pi(v)$  are adjacent in  $H$ . If an  
1786 isomorphism  $\pi$  exists between  $G$  and  $H$ , the graphs are isomorphic and written  $G \simeq H$  or  $G \stackrel{\pi}{\simeq} H$ .  
1787 By  $X \stackrel{d}{=} Y$  we denote the equality of two random variables  $X$  and  $Y$  in distributions.  
1788

1789 In the proofs, we write by  $\text{id}(\cdot)$  the namespace obtained by anonymization of a walk  $v_0 \rightarrow \dots \rightarrow v_l$   
1790 (Section 2) that maps each vertex  $v_t$  to its unique integer name  $\text{id}(v_t)$  starting from  $\text{id}(v_0) = 1$ .  
1791

1792 Let us define some notations related to cover times. These will be used from Appendix A.5.5. We  
1793 denote by  $H(u, v)$  the expected number of steps a random walk starting from  $u$  takes until reaching  $v$ .  
1794 This is called the hitting time between  $u$  and  $v$ . For a graph  $G$ , let  $C_V(G, u)$  be the expected number  
1795 of steps a random walk starting from  $u$  takes until visiting every vertices of  $G$ . The vertex cover time  
1796  $C_V(G)$ , or the cover time, is this quantity given by the worst possible  $u$ :

$$1796 \quad C_V(G) := \max_{u \in V(G)} C_V(G, u). \quad (16)$$

1798 Likewise, let  $C_E(G, u)$  be the expected number of steps a random walk starting from  $u$  takes until  
1799 traversing every edge of  $G$ . The edge cover time  $C_E(G)$  is given by the worst possible  $u$ :

$$1801 \quad C_E(G) := \max_{u \in V(G)} C_E(G, u). \quad (17)$$

1803 For local balls  $B_r(u)$ , we always set the starting vertex of the random walk to  $u$  (Section 2.2). Thus,  
1804 we define their cover times as follows:

$$1805 \quad C_V(B_r(u)) := C_V(B_r(u), u), \quad (18)$$

$$1806 \quad C_E(B_r(u)) := C_E(B_r(u), u). \quad (19)$$

1808 A.5.1 PROOF OF PROPOSITION 2.1 (SECTION 2)  
1809

1810 **Proposition 2.1.**  $X_\theta(\cdot)$  is invariant in probability, if its random walk algorithm is invariant in  
1811 probability and its recording function is invariant.

1812 *Proof.* We recall the probabilistic invariance of random walk algorithm in Equation 3:

$$1814 \quad \pi(v_0) \rightarrow \dots \rightarrow \pi(v_l) \stackrel{d}{=} u_0 \rightarrow \dots \rightarrow u_l, \quad \forall G \stackrel{\pi}{\simeq} H, \quad (20)$$

1816 where  $v_{[\cdot]}$  is a random walk on  $G$  and  $u_{[\cdot]}$  is a random walk on  $H$ . We further recall the invariance of  
1817 recording function  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$  in Equation 7:

$$1819 \quad q(v_0 \rightarrow \dots \rightarrow v_l, G) = q(\pi(v_0) \rightarrow \dots \rightarrow \pi(v_l), H), \quad \forall G \stackrel{\pi}{\simeq} H, \quad (21)$$

1820 for any given random walk  $v_{[\cdot]}$  on  $G$ . Combining Equation 20 and equation 21, we have the following:  
1821

$$1822 \quad q(v_0 \rightarrow \dots \rightarrow v_l, G) \stackrel{d}{=} q(u_0 \rightarrow \dots \rightarrow u_l, H), \quad \forall G \simeq H. \quad (22)$$

1824 Then, since the reader neural network  $f_\theta$  is a deterministic map, we have:

$$1825 \quad f_\theta(q(v_0 \rightarrow \dots \rightarrow v_l, G)) \stackrel{d}{=} f_\theta(q(u_0 \rightarrow \dots \rightarrow u_l, H)), \quad \forall G \simeq H, \quad (23)$$

1827 which leads to:

$$1828 \quad X_\theta(G) \stackrel{d}{=} X_\theta(H), \quad \forall G \simeq H. \quad (24)$$

1830 This shows Equation 2 and completes the proof.  $\square$

1832 A.5.2 PROOF OF PROPOSITION 2.2 (SECTION 2)  
1833

1834 We first show a useful lemma.

1835 **Lemma A.2.** The random walk in Equation 4 is invariant in probability, if the probability distributions  
of the starting vertex  $v_0$  and each transition  $v_{t-1} \rightarrow v_t$  are invariant.

1836 *Proof.* We prove by induction. Let us write probabilistic invariance of the starting vertex  $v_0$  as:  
1837

$$1838 \pi(v_0) \stackrel{d}{=} u_0, \quad \forall G \stackrel{\pi}{\simeq} H, \quad (25)$$

1839 and probabilistic invariance of each transition  $v_{t-1} \rightarrow v_t$  as follows, by fixing  $v_{t-1}$  and  $u_{t-1}$ :  
1840

$$1841 \pi(v_{t-1}) \rightarrow \pi(v_t) \stackrel{d}{=} u_{t-1} \rightarrow u_t, \quad \forall G \stackrel{\pi}{\simeq} H, v_{t-1} \in V(G), u_{t-1} := \pi(v_{t-1}). \quad (26)$$

1842 Let us assume the following for some  $t \geq 1$ :  
1843

$$1844 \pi(v_0) \rightarrow \cdots \rightarrow \pi(v_{t-1}) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_{t-1}, \quad G \stackrel{\pi}{\simeq} H. \quad (27)$$

1845 Then, from the chain rule of joint distributions, the first-order Markov property of random walk in  
1846 Equation 4, and probabilistic invariance in Equation 26, we obtain the following:  
1847

$$1848 \pi(v_0) \rightarrow \cdots \rightarrow \pi(v_t) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_t, \quad G \stackrel{\pi}{\simeq} H. \quad (28)$$

1849 By induction from the initial condition  $\pi(v_0) \stackrel{d}{=} u_0$  in Equation 25, the following holds  $\forall l > 0$ :  
1850

$$1851 \pi(v_0) \rightarrow \cdots \rightarrow \pi(v_l) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_l, \quad \forall G \stackrel{\pi}{\simeq} H. \quad (29)$$

1852 This shows Equation 3 and completes the proof.  $\square$   
1853

1854 We now prove Proposition 2.2.  
1855

1856 **Proposition 2.2.** *The random walk in Equation 4 is invariant in probability if its conductance  $c_{[\cdot]}(\cdot)$   
1857 is invariant:*

$$1858 c_G(u, v) = c_H(\pi(u), \pi(v)), \quad \forall G \stackrel{\pi}{\simeq} H. \quad (30)$$

1859 It includes constant conductance, and any choice that only uses degrees of endpoints  $\deg(u)$ ,  $\deg(v)$ .  
1860

1861 *Proof.* We recall Equation 4 for a given conductance function  $c_G : E(G) \rightarrow \mathbb{R}_+$ :  
1862

$$1863 \text{Prob}[v_t = x | v_{t-1} = u] := \frac{c_G(u, x)}{\sum_{y \in N(u)} c_G(u, y)}. \quad (31)$$

1864 From Lemma A.2, it is sufficient to show the probabilistic invariance of each transition  $v_{t-1} \rightarrow v_t$  as  
1865 we sample  $v_0$  from invariant distribution  $\text{Uniform}(V(G))$  (Algorithm 1). We rewrite Equation 26 as:  
1866

$$1867 \text{Prob}[v_t = x | v_{t-1} = u] = \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(u)], \quad \forall G \stackrel{\pi}{\simeq} H. \quad (32)$$

1868 If the conductance function  $c_{[\cdot]}(\cdot)$  is invariant (Equation 30), we can show Equation 32 by:  
1869

$$1870 \begin{aligned} 1871 \text{Prob}[v_t = x | v_{t-1} = u] &:= \frac{c_G(u, x)}{\sum_{y \in N(u)} c_G(u, y)}, \\ 1872 &= \frac{c_H(\pi(u), \pi(x))}{\sum_{y \in N(u)} c_H(\pi(u), \pi(y))}, \\ 1873 &= \frac{c_H(\pi(u), \pi(x))}{\sum_{\pi(y) \in N(\pi(u))} c_H(\pi(u), \pi(y))}, \\ 1874 &= \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(u)], \quad \forall G \stackrel{\pi}{\simeq} H. \end{aligned} \quad (33)$$

1875 In the third equality, we used the fact that isomorphism  $\pi(\cdot)$  preserves adjacency. It is clear that any  
1876 conductance function  $c_{[\cdot]}(\cdot)$  with a constant conductance such as  $c_G(\cdot) = 1$  is invariant. Further-  
1877 more, any conductance function that only uses degrees of endpoints  $\deg(u)$ ,  $\deg(v)$  is invariant as  
1878 isomorphism  $\pi(\cdot)$  preserves the degree of each vertex. This completes the proof.  $\square$   
1879

## A.5.3 PROOF OF PROPOSITION A.1 (APPENDIX A.1)

We extend the proof of Proposition 2.2 given in Appendix A.5.2 to second-order random walks discussed in Appendix A.1. We first show a useful lemma:

**Lemma A.3.** *A second-order random walk algorithm is invariant in probability, if the probability distributions of the starting vertex  $v_0$ , the first transition  $v_0 \rightarrow v_1$ , and each transition  $(v_{t-2}, v_{t-1}) \rightarrow v_t$  for  $t > 1$  are invariant.*

*Proof.* We prove by induction. The probabilistic invariance of starting vertex  $v_0$  and first transition  $v_0 \rightarrow v_1$  are:

$$\pi(v_0) \stackrel{d}{=} u_0, \quad \forall G \stackrel{\pi}{\simeq} H, \quad (34)$$

$$\pi(v_0) \rightarrow \pi(v_1) \stackrel{d}{=} u_0 \rightarrow u_1, \quad \forall G \stackrel{\pi}{\simeq} H, v_0 \in V(G), u_0 := \pi(v_0), \quad (35)$$

and probabilistic invariance of each transition  $(v_{t-2}, v_{t-1}) \rightarrow v_t$  for  $t > 1$  can be written as follows by fixing  $(v_{t-2}, v_{t-1})$  and  $(u_{t-2}, u_{t-1})$ :

$$\begin{aligned} \pi(v_{t-2}) \rightarrow \pi(v_{t-1}) \rightarrow \pi(v_t) \stackrel{d}{=} u_{t-2} \rightarrow u_{t-1} \rightarrow u_t, \quad \forall G \stackrel{\pi}{\simeq} H, \\ (v_{t-2}, v_{t-1}) \in E(G), \\ (u_{t-2}, u_{t-1}) := (\pi(v_{t-2}), \pi(v_{t-1})). \end{aligned} \quad (36)$$

Let us assume the following for some  $t \geq 2$ :

$$\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_{t-2}) \rightarrow \pi(v_{t-1}) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_{t-2} \rightarrow u_{t-1}, \quad G \stackrel{\pi}{\simeq} H. \quad (37)$$

Then, from the chain rule of joint distributions, the second-order Markov property of second-order random walks, and probabilistic invariance in Equation 36, we obtain the following:

$$\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_{t-1}) \rightarrow \pi(v_t) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_{t-1} \rightarrow u_t, \quad G \stackrel{\pi}{\simeq} H. \quad (38)$$

By induction from the initial condition  $\pi(v_0) \rightarrow \pi(v_1) \stackrel{d}{=} u_0 \rightarrow u_1$  given from Equations 34 and 35, the following holds  $\forall l > 0$ :

$$\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_l) \stackrel{d}{=} u_0 \rightarrow \cdots \rightarrow u_l, \quad \forall G \stackrel{\pi}{\simeq} H. \quad (39)$$

This shows Equation 3 and completes the proof.  $\square$

We now prove Proposition A.1.

**Proposition A.1.** *The non-backtracking random walk in Equation 13 and the node2vec random walk in Equation 14 are invariant in probability, if their underlying first-order random walk algorithm is invariant in probability.*

*Proof.* We assume that the first transition  $v_0 \rightarrow v_1$  is given by the underlying first-order random walk algorithm. This is true for non-backtracking since there is no  $v_{t-2}$  to avoid, and true for the official implementation of node2vec (Grover & Leskovec, 2016). Then from Lemma A.3, it is sufficient to show probabilistic invariance of each transition  $(v_{t-2}, v_{t-1}) \rightarrow v_t$  for  $t > 1$  as we sample  $v_0$  from the invariant distribution  $\text{Uniform}(V(G))$  (Algorithm 1) and the first transition  $v_0 \rightarrow v_1$  is given by the first-order random walk which is assumed to be invariant in probability. Rewrite Equation 36 as:

$$\text{Prob}[v_t = x | v_{t-1} = j, v_{t-2} = i] = \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(j), u_{t-2} = \pi(i)], \quad \forall G \stackrel{\pi}{\simeq} H. \quad (40)$$

For non-backtracking, we first handle the case where  $i \rightarrow j$  reaches a dangling vertex  $j$  which has  $i$  as its only neighbor. In this case the walk begrudgingly backtracks  $i \rightarrow j \rightarrow i$  (Appendix A.1). Since isomorphism  $\pi(\cdot)$  preserves adjacency,  $\pi(i) \rightarrow \pi(j)$  also reaches a dangling vertex  $\pi(j)$  and must begrudgingly backtrack  $\pi(i) \rightarrow \pi(j) \rightarrow \pi(i)$ . By interpreting the distributions on  $v_t$  and  $u_t$  as one-hot at  $i$  and  $\pi(i)$ , respectively, we can see that Equation 40 holds. If  $i \rightarrow j$  does not reach a

dangling vertex, the walk  $j \rightarrow x$  follows the invariant probability of the underlying first-order random walk  $\text{Prob}[v_{t+1} = x | v_t = j]$  renormalized over  $N(j) \setminus \{i\}$ . Then we can show Equation 40 by:

$$\begin{aligned}
& \text{Prob}[v_t = x | v_{t-1} = j, v_{t-2} = i] \\
& := \begin{cases} \frac{\text{Prob}[v_t = x | v_{t-1} = j]}{\sum_{y \in N(j) \setminus \{i\}} \text{Prob}[v_t = y | v_{t-1} = j]} & \text{for } x \neq i, \\ 0 & \text{for } x = i, \end{cases} \\
& = \begin{cases} \frac{\text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(j)]}{\sum_{\pi(y) \in N(\pi(j)) \setminus \{\pi(i)\}} \text{Prob}[u_t = \pi(y) | u_{t-1} = \pi(j)]} & \text{for } \pi(x) \neq \pi(i), \\ 0 & \text{for } \pi(x) = \pi(i), \end{cases} \\
& = \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(j), u_{t-2} = \pi(i)], \quad \forall G \stackrel{\pi}{\simeq} H. \tag{41}
\end{aligned}$$

In the second equality, we have used the fact that isomorphism  $\pi(\cdot)$  is a bijection and preserves adjacency, and the assumption that the first-order random walk algorithm is invariant in probability. For node2vec walk, we first show the invariance of the weighting term  $\alpha(i, x)$  in Equation 15 using the fact that isomorphism preserves shortest path distances between vertices  $d(i, x) = d(\pi(i), \pi(x))$ :

$$\begin{aligned}
\alpha(i, x) & := \begin{cases} 1/p & \text{for } d(i, x) = 0, \\ 1 & \text{for } d(i, x) = 1, \\ 1/q & \text{for } d(i, x) = 2. \end{cases} \\
& = \begin{cases} 1/p & \text{for } d(\pi(i), \pi(x)) = 0, \\ 1 & \text{for } d(\pi(i), \pi(x)) = 1, \\ 1/q & \text{for } d(\pi(i), \pi(x)) = 2. \end{cases} \\
& = \alpha(\pi(i), \pi(x)). \tag{42}
\end{aligned}$$

Then we can show Equation 40 by:

$$\begin{aligned}
\text{Prob}[v_t = x | v_{t-1} = j, v_{t-2} = i] & := \frac{\alpha(i, x) \text{Prob}[v_t = x | v_{t-1} = j]}{\sum_{y \in N(j)} \alpha(i, y) \text{Prob}[v_t = y | v_{t-1} = j]}, \\
& = \frac{\alpha(\pi(i), \pi(x)) \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(j)]}{\sum_{\pi(y) \in N(\pi(j))} \alpha(\pi(i), \pi(y)) \text{Prob}[u_t = \pi(y) | u_{t-1} = \pi(j)]}, \\
& = \text{Prob}[u_t = \pi(x) | u_{t-1} = \pi(j), u_{t-2} = \pi(i)], \quad \forall G \stackrel{\pi}{\simeq} H. \tag{43}
\end{aligned}$$

In the second equality, we have used the fact that isomorphism  $\pi(\cdot)$  preserves adjacency, and the assumption that the first-order walk algorithm is invariant in probability. This completes the proof.  $\square$

#### A.5.4 PROOF OF PROPOSITION 2.3 (SECTION 2)

**Proposition 2.3.** *A recording function  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$  that uses anonymization, optionally with named neighbors, is invariant.*

*Proof.* Given a walk  $v_0 \rightarrow \dots \rightarrow v_l$  on  $G$ , we can write the anonymized namespace  $\text{id}(\cdot)$  as follows:

$$\text{id}(v_t) = 1 + \arg \min_i [v_i = v_t]. \tag{44}$$

Consider any walk  $\pi(v_0) \rightarrow \dots \rightarrow \pi(v_l)$  on some  $H$  which is isomorphic to  $G$  via  $\pi$ . Let us denote the anonymization of this walk as  $\text{id}'(\cdot)$ . Then we have:

$$\begin{aligned}
\text{id}'(\pi(v_t)) & = 1 + \arg \min_i [\pi(v_i) = \pi(v_t)], \\
& = 1 + \arg \min_i [v_i = v_t], \\
& = \text{id}(v_t), \quad \forall G \stackrel{\pi}{\simeq} H. \tag{45}
\end{aligned}$$

1998 The second equality is due to the fact that  $\pi$  is a bijection. This completes the proof for anonymization.  
 1999 For the named neighbors, we prove by induction. Let us fix some  $t \geq 1$  and assume:

$$2000 \quad q(v_0 \rightarrow \cdots \rightarrow v_{t-1}, G) = q(\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_{t-1}), H), \quad \forall G \stackrel{\pi}{\simeq} H. \quad (46)$$

2002 Let  $T \subseteq E(G)$  be the set of edges recorded by  $\mathbf{z} := q(v_0 \rightarrow \cdots \rightarrow v_{t-1}, G)$ , and  $T' \subseteq E(H)$  be the  
 2003 set of edges recorded by  $\mathbf{z}' := q(\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_{t-1}), H)$ . Then,  $T$  and  $T'$  are isomorphic via  $\pi$ :

$$2004 \quad T' = \{(\pi(u), \pi(v)) \mid (u, v) \in T\}. \quad (47)$$

2006 The equality is shown as follows. ( $\supseteq$ ) Any  $(u, v) \in T$  is recorded in  $\mathbf{z}$  as  $(\text{id}(u), \text{id}(v))$ . From  
 2007 Equations 46 and 45, we have  $\mathbf{z} = \mathbf{z}'$  and  $(\text{id}(u), \text{id}(v)) = (\text{id}(\pi(u)), \text{id}(\pi(v)))$ . Thus,  $(\pi(u), \pi(v))$   
 2008 is recorded in  $\mathbf{z}'$  as  $(\text{id}(\pi(u)), \text{id}(\pi(v)))$ , i.e.  $(\pi(u), \pi(v)) \in T'$ . ( $\subseteq$ ) This follows from symmetry.  
 2009 Now, we choose some  $v_t \in N(v_{t-1})$  and consider the set  $F \subseteq E(G)$  of all unrecorded edges from  $v_t$ :

$$2010 \quad F := \{(v_t, u) : u \in N(v_t)\} \setminus T. \quad (48)$$

2012 Then, the set  $D$  of unrecorded neighbors of  $v_t$  is given as:

$$2013 \quad D := \{u : (v_t, u) \in F\}. \quad (49)$$

2015 Since  $\pi(\cdot)$  preserves adjacency, the set  $F' \subseteq E(H)$  of all unrecorded edges from  $\pi(v_t)$  is given by:

$$\begin{aligned} 2016 \quad F' &:= \{(\pi(v_t), \pi(u)) : \pi(u) \in N(\pi(v_t))\} \setminus T', \\ 2017 \quad &= \{(\pi(v_t), \pi(u)) : u \in N(v_t)\} \setminus \{(\pi(u'), \pi(v')) \mid (u', v') \in T\}, \\ 2018 \quad &= \{(\pi(v_t), \pi(u)) : u \in N(v_t), (v_t, u) \notin T\}, \\ 2019 \quad &= \{(\pi(v_t), \pi(u)) : (v_t, u) \in F\}, \end{aligned} \quad (50)$$

2021 and consequently:

$$\begin{aligned} 2022 \quad D' &:= \{\pi(u) : (\pi(v_t), \pi(u)) \in F'\}, \\ 2023 \quad &= \{\pi(u) : u \in D\}. \end{aligned} \quad (51)$$

2025 While  $D$  and  $D'$  may contain vertices not yet visited by the walks and hence not named by  $\text{id}(\cdot)$ ,  
 2026 what we record are named neighbors. Let  $S$  be the set of all named vertices in  $G$ . It is clear that the  
 2027 set  $S'$  of named vertices in  $H$  is given by  $S' = \{\pi(v) : v \in S\}$ . Then, the named neighbors in  $G$   
 2028 to be newly recorded for  $v_{t-1} \rightarrow v_t$  is given by  $U = D \cap S$ , and for  $\pi(v_{t-1}) \rightarrow \pi(v_t)$  in  $H$  the set  
 2029 is given by  $U' = D' \cap S'$ . Since  $\pi(\cdot)$  is a bijection, we can see that  $U' = \{\pi(u) : u \in U\}$ . From  
 2030 the invariance of anonymization in Equation 45,  $U$  and  $U'$  are named identically  $\{\text{id}(u) : u \in U\} =$   
 2031  $\{\text{id}(\pi(u)) : \pi(u) \in U'\}$ , and therefore will be recorded identically. As a result, the information to be  
 2032 added to the records at time  $t$  are identical, and we have:

$$2033 \quad q(v_0 \rightarrow \cdots \rightarrow v_t, G) = q(\pi(v_0) \rightarrow \cdots \rightarrow \pi(v_t), H), \quad \forall G \stackrel{\pi}{\simeq} H. \quad (52)$$

2035 Then, by induction from the initial condition:

$$2036 \quad q(v_0, G) = q(\pi(v_0), H) = \text{id}(v_0) = \text{id}(\pi(v_0)) = 1, \quad (53)$$

2038 the recording function using anonymization and named neighbors is invariant.  $\square$

## 2040 A.5.5 PROOF OF THEOREM 2.4 (SECTION 2.2)

2041 **Theorem 2.4.** *For a uniform random walk on an infinite graph  $G$  starting at  $v$ , the vertex and edge*  
 2042 *cover times of the finite local ball  $B_r(v)$  are not always finitely bounded.*

2044 *Proof.* We revisit the known fact that hitting times on the infinite line  $\mathbb{Z}$  is infinite (Dumitriu et al.,  
 2045 2003; Janson & Peres, 2012; carnage, 2012; McNew, 2013). Consider a local ball  $B_1(0)$  of radius  
 2046 1 from the origin 0 such that  $V(B_1(0)) = \{-1, 0, 1\}$ . Let us assume that the expected number of  
 2047 steps for a random walk starting at 0 to visit 1 for the first time is finite, and denote it by  $T$ . In  
 2048 the first step of the walk, we have to go either to  $-1$  or 1. If we go to  $-1$ , we have to get back  
 2049 to 0 and then to 1 which would take  $2T$  in expectation (by translation symmetry). This leads to  
 2050  $T = 1/2 \cdot 1 + 1/2 \cdot (1 + 2T) = 1 + T$ , which is a contradiction. Since visiting all vertices or  
 2051 traversing all edges in  $B_1(0)$  clearly involves visiting 1, the vertex and edge cover times cannot be  
 finitely bounded.  $\square$

## A.5.6 PROOF OF THEOREM 2.5 (SECTION 2.2)

Let us recall that our graph  $G$  is locally finite (Section 2.2), and denote its maximum degree as  $\Delta$ . We further denote by  $d(u, v)$  the shortest path distance between  $u$  and  $v$ . We show some useful lemmas. We first show that  $H(u, x)$ , the expected number of steps of a random walk starting at  $u$  takes until reaching  $x$ , is finite for any  $u, x$  in  $B_r(v)$  for any nonzero restart probability  $\alpha$  or restart period  $k \geq r$ .

**Lemma A.4.** *For any  $u$  and  $x$  in  $B_r(v)$ ,  $H(u, x)$  is bounded by the following:*

$$H(u, x) \leq \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^r + \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r}, \quad (54)$$

if the random walk restarts at  $v$  with any probability  $\alpha \in (0, 1)$ , and:

$$H(u, x) \leq k + k\Delta^r, \quad (55)$$

if the random walk restarts at  $v$  with any period  $k \geq r$ .

*Proof.* We first prove for random restarts. The proof is inspired by Theorem 1.1 of McNew (2013) and Lemma 2 of Zuckerman (1991). We clarify some measure-theoretic details. Let  $W$  be the set of all (infinite) random walks on  $G$ . We denote by  $P$  the probability measure defined on the space  $W$ , equipped with its cylinder  $\sigma$ -algebra, by the random walk restarting at  $v$  with probability  $\alpha \in (0, 1)$ . For a vertex  $x$  in  $G$ , we define the hitting time function  $|\cdot|_x : W \rightarrow \mathbb{N} \cup \{\infty\}$  as follows:

$$|w|_x := \arg \min_i [(w)_i = x], \quad \forall w \in W. \quad (56)$$

Let  $W(a)$  be the set of all random walks that start at  $a$ ,  $W(a, b)$  be the set of random walks that start at  $a$  and visit  $b$ , and  $W(a, b^-)$  be the set of random walks that start at  $a$  and do not visit  $b$ . Then, the measurability of  $|\cdot|_x$  follows from the measurability of the sets  $W(a)$ ,  $W(a, b)$ , and  $W(a, b^-)$ , which we show as follows.  $\forall a, b \in V(G)$ ,  $W(a)$  is clearly measurable,  $W(a, b)$  is measurable as it equals  $\bigcup_{n \in \mathbb{N}} \{w \in W \mid (w)_1 = a, (w)_n = b\}$ , and  $W(a, b^-)$  is measurable as it is  $W(a) \setminus W(a, b)$ .

Consider  $W(u, x^-)$ , the set of random walks that start at  $u$  and never visit  $x$ . We start by showing that this set is of measure zero:

$$P(W(u, x^-)) = 0. \quad (57)$$

For this, we use the fact that  $W(u, x^-) = W(u, v^-, x^-) \cup W(u, v, x^-)$ , where  $W(u, v^-, x^-)$  is the set of random walks starting at  $u$  that do not visit  $v$  nor  $x$ , and  $W(u, v, x^-)$  is the set of walks starting at  $u$  that visit  $v$  and do not visit  $x$ . We have the following:

$$\begin{aligned} P(W(u, x^-)) &= P(W(u, v^-, x^-)) + P(W(u, v, x^-)), \\ &\leq P(W(u, v^-)) + P(W(u, v, x^-)). \end{aligned} \quad (58)$$

We show Equation 57 by showing  $P(W(u, v^-)) = 0$  and  $P(W(u, v, x^-)) = 0$  in Equation 58.

1. ( $P(W(u, v^-)) = 0$ ) Consider  $W(u, v^-)$ , the set of all random walks that start at  $u$  and never visit  $v$ . Since restart sends a walk to  $v$ , every walk in this set never restarts. The probability of a walk not restarting until step  $t$  is  $(1 - \alpha)^t$ . Denote this probability by  $p_t$ , and let  $p$  be the probability of a random walk to never restart. Then  $p_t \downarrow p = 0$  and  $P(W(u, v^-)) \leq p = 0$ .
2. ( $P(W(u, v, x^-)) = 0$ ) Assume  $P(W(u, v, x^-)) > 0$ . Then  $P(W(v, x^-)) > 0$  since each walk step is independent. Let  $W_N(v, x^-)$  be the set of walks that start at  $v$  and do not reach  $x$  within  $N$  restarts. Then we have  $W_N(v, x^-) \downarrow W(v, x^-)$ . If a walk restarts at  $v$ , the probability of reaching  $x$  before the next restart is at least the probability of exactly walking the shortest path from  $v$  to  $x$ , which is  $\geq (\frac{1-\alpha}{\Delta})^{d(v, x)} \in (0, 1)$ . Then  $P(W_N(v, x^-)) \leq (1 - (\frac{1-\alpha}{\Delta})^{d(v, x)})^N \downarrow 0$ , leading to  $P(W(v, x^-)) = 0$ . This is a contradiction, so we have  $P(W(u, v, x^-)) = 0$ .

We are now ready to bound  $H(u, x)$ . Using the hitting time function  $|\cdot|_x$  in Equation 56, we have:

$$H(u, x) = \mathbb{E}[|w|_x \mid w \in W(u)]. \quad (59)$$

Since  $W(u) = W(u, x) \cup W(u, x^-)$ , and  $P(W(u, x^-)) = 0$  from Equation 57, we have:

$$\begin{aligned} H(u, x) &= \mathbb{E}[|w|_x | w \in W(u, x)], \\ &= \int_{W(u, x)} |w|_x dP(w | w \in W(u, x)). \end{aligned} \quad (60)$$

Each walk in  $W(u, x)$  starts at  $u$  and, when it reaches  $x$ , it either has or has not visited  $v$ . We treat the two cases separately. Let  $W(a, b, c)$  be the set of walks that start at  $a$  and reach  $c$  after  $b$ , and let  $W(a, b^-, c)$  be the set of walks that start at  $a$  and reach  $c$  before  $b$ . Then we have:

$$\begin{aligned} H(u, x) &= \int_{W(u, v, x) \cup W(u, v^-, x)} |w|_x dP(w | W(u, x)), \\ &= \int_{W(u, v, x)} |w|_x dP(w | W(u, x)) + \int_{W(u, v^-, x)} |w|_x dP(w | W(u, x)). \end{aligned} \quad (61)$$

Let  $\hat{H}(a, b, c)$  (or  $\hat{H}(a, b^-, c)$ ) be the expected number of steps for a walk starting at  $a$  to reach  $c$  after reaching  $b$  (or before  $b$ ), given that it is in  $W(a, b, c)$  (or in  $W(a, b^-, c)$ ). If a walk from  $u$  reaches  $v$  before  $x$ , the expected steps is given by  $\hat{H}(u, x^-, v) + H(v, x)$ . If the walk reaches  $x$  before  $v$ , the expected steps is  $\hat{H}(u, v^-, x)$ . Then, if  $W(u, v^-, x) \neq \emptyset$ , we can write  $H(u, x)$  as follows:

$$\begin{aligned} H(u, x) &= \left[ \hat{H}(u, x^-, v) + H(v, x) \right] P(W(u, v, x) | W(u, x)) \\ &\quad + \hat{H}(u, v^-, x) P(W(u, v^-, x) | W(u, x)). \end{aligned} \quad (62)$$

On the other hand, if  $W(u, v^-, x) = \emptyset$ , we simply have:

$$H(u, x) = \hat{H}(u, x^-, v) + H(v, x). \quad (63)$$

We first consider  $\hat{H}(u, x^-, v)$ , the expected number of steps from  $u$  to reach  $v$  before  $x$ . We show:

$$\hat{H}(u, x^-, v) \leq \frac{1}{\alpha}. \quad (64)$$

To see this, note that  $\hat{H}(u, x^-, v)$  is equivalent to the expectation  $\mathbb{E}[T]$  of the number of steps  $T$  to reach  $v$  from  $u$ , on the graph  $G$  with the vertex  $x$  deleted and transition probabilities renormalized. If  $u$  is isolated on this modified graph, the only walk in  $W(u, x^-, v)$  is the one that immediately restarts at  $v$ , giving  $\mathbb{E}[T] = 1$ . Otherwise, we have  $\mathbb{E}[T] \leq 1/\alpha$  due to the following. Let  $T'$  be the number of steps until the first restart at  $v$ . Since restart can be treated as a Bernoulli trial with probability of success  $\alpha$ ,  $T'$  follows geometric distribution with expectation  $1/\alpha$ . Since it is possible for the walk to reach  $v$  before the first restart, we have  $\mathbb{E}[T] \leq \mathbb{E}[T'] = 1/\alpha$ , which gives Equation 64.

We now consider  $H(v, x)$ , the expectation  $\mathbb{E}[T]$  of the number of steps  $T$  to reach  $x$  from  $v$ . Let  $T'$  be the steps until the walk restarts at  $v$ , then exactly walks the shortest path from  $v$  to  $x$  for the first time, and then restarts at  $v$ . Since  $T'$  walks until restart after walking the shortest path to  $x$ , and it is possible for a walk to reach  $x$  before walking the shortest path to it, we have  $\mathbb{E}[T] \leq \mathbb{E}[T']$ . Then, we split the walk of length  $T'$  into  $N$  trials, where each trial consists of restarting at  $v$  and walking until the next restart. A trial is successful if it immediately walks the shortest path from  $v$  to  $x$ . Then  $N$  is the number of trials until we succeed, and it follows geometric distribution with probability of success at least  $(\frac{1-\alpha}{\Delta})^{d(v, x)}$  due to bounded degrees. Its expectation is then bounded as:

$$\mathbb{E}[N] \leq \left( \frac{\Delta}{1-\alpha} \right)^{d(v, x)}. \quad (65)$$

Let  $S_i$  be the length of the  $i$ -th trial. Since each  $S_i$  is i.i.d. with finite mean  $\mathbb{E}[S_i] = 1/\alpha$ , and  $N$  is stopping time, we can apply Wald's identity (Hein) to compute the expectation of  $T'$ :

$$\begin{aligned} \mathbb{E}[T'] &= \mathbb{E}[S_1 + \dots + S_N], \\ &= \mathbb{E}[N] \mathbb{E}[S_1], \\ &\leq \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{d(v, x)}. \end{aligned} \quad (66)$$



2160 We remark that  $H(v, x) \leq \mathbb{E}[T']$ . Combining this result with Equation 64, we have:

$$2161 \hat{H}(u, x^-, v) + H(v, x) \leq \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{d(v,x)}. \quad (67)$$

2164 If  $W(u, v^-, x) = \emptyset$ , we have  $H(u, x) = \hat{H}(u, x^-, v) + H(v, x)$  from Equation 63 and it is finitely  
2165 bounded for any  $\alpha \in (0, 1)$  by the above. If  $W(u, v^-, x) \neq \emptyset$ , Equations 62 and 67 lead to:

$$2166 H(u, x) \leq \hat{H}(u, x^-, v) + H(v, x) + \hat{H}(u, v^-, x) P(W(u, v^-, x)|W(u, x)), \\ 2167 \leq \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{d(v,x)} + \hat{H}(u, v^-, x) P(W(u, v^-, x)|W(u, x)), \quad (68)$$

2171 and it suffices to bound  $\hat{H}(u, v^-, x) P(W(u, v^-, x)|W(u, x))$ . We show the following:

$$2172 \hat{H}(u, v^-, x) = \int_{W(u, v^-, x)} |w|_x dP(w|W(u, v^-, x)), \\ 2173 = \sum_{k=1}^{\infty} \int_{\{w \in W(u, v^-, x) \wedge |w|_x = k\}} |w|_x dP(w|W(u, v^-, x)), \\ 2174 = \sum_{k=1}^{\infty} k \int_{\{w \in W(u, v^-, x) \wedge |w|_x = k\}} dP(w|W(u, v^-, x)), \\ 2175 = \sum_{k=1}^{\infty} k P(\{w \in W(u, v^-, x) \wedge |w|_x = k\} | W(u, v^-, x)), \\ 2176 \leq \sum_{k=1}^{\infty} k P(\{w : |w|_x = k\} | W(u, v^-, x)), \\ 2177 = \sum_{k=1}^{\infty} k \frac{P(\{w : |w|_x = k \wedge w \in W(u, v^-, x)\})}{P(W(u, v^-, x))}, \\ 2178 \leq \sum_{k=1}^{\infty} k (1-\alpha)^k \frac{1}{P(W(u, v^-, x))}, \\ 2179 = \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \frac{1}{P(W(u, v^-, x))}. \quad (69)$$

2183 We have used Fubini's theorem for the second equality. Then we have:

$$2184 \hat{H}(u, v^-, x) P(W(u, v^-, x)|W(u, x)) \leq \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \frac{P(W(u, v^-, x)|W(u, x))}{P(W(u, v^-, x))}, \\ 2185 = \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \frac{1}{P(W(u, x))}. \quad (70)$$

2186  $P(W(u, x))$  is at least the probability of precisely walking the shortest path from  $u$  to  $x$ , which has  
2187 length  $d(u, x) \leq 2r$  since  $u$  and  $x$  are both in  $B_r(v)$ . This gives us the following:

$$2188 P(W(u, x)) \geq \left( \frac{1-\alpha}{\Delta} \right)^{d(u,x)}, \\ 2189 \geq \left( \frac{1-\alpha}{\Delta} \right)^{2r}. \quad (71)$$

2190 Combining this with Equation 70, we have:

$$2191 \hat{H}(u, v^-, x) P(W(u, v^-, x)|W(u, x)) \leq \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r}. \quad (72)$$

2192 Combining with Equations 62 and 67, we have:

$$2193 H(u, x) \leq \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{d(v,x)} + \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r}, \quad (73)$$

for any  $\alpha \in (0, 1)$ . Notice that, while this bound is for the case  $W(u, v^-, x) \neq \emptyset$ , it subsumes the bound for the case  $W(u, v^-, x) = \emptyset$  (Equation 67). Then, using  $d(v, x) \leq r$ , we get Equation 54. This completes the proof for random restarts.

We now prove for periodic restarts. If the walk starting at  $u$  reaches  $x$  before restarting at  $v$ , the steps taken is clearly less than  $k$ . If the walk starting at  $u$  restarts at  $v$  at step  $k$  before reaching  $x$ , it now needs to reach  $x$  from  $v$ , while restarting at  $v$  at every  $k$  steps. Let  $T$  be the steps taken to reach  $x$  from  $v$ . Let  $T'$  be the number of steps until the walk restarts at  $v$ , then exactly follows the shortest path from  $v$  to  $x$  for the first time, and then restarts at  $v$ . It is clear that  $\mathbb{E}[T] \leq \mathbb{E}[T']$ . Then, we split the walk of length  $T'$  into  $N$  trials, where each trial consists of restarting at  $v$  and walking  $k$  steps until the next restart. A trial is successful if it immediately walks the shortest path from  $v$  to  $x$ . Then  $N$  is the number of trials until we get a success, and it follows geometric distribution with probability of success at least  $(1/\Delta)^{d(v,x)}$  only for  $k \geq d(v, x)$ , and zero for  $k < d(v, x)$  since the walk cannot reach  $x$  before restart. Hence, its expectation is at most  $\Delta^{d(v,x)}$  for  $k \geq d(v, x)$ , and we have  $\mathbb{E}[T] \leq \mathbb{E}[T'] = k\mathbb{E}[N] \leq k\Delta^{d(v,x)}$ . Adding the  $k$  steps until the first restart at  $v$ , we have:

$$H(u, x) \leq k + k\Delta^{d(v,x)}, \quad (74)$$

for any  $k \geq d(v, x)$ . Using  $d(v, x) \leq r$ , we get Equation 55. This completes the proof.  $\square$

We now extend Lemma A.4 to edges. Let  $H(u, (x, y))$  be the expected number of steps of a random walk starting at  $u$  takes until traversing an edge  $(x, y)$  by  $x \rightarrow y$ . We show that  $H(u, (x, y))$  is finite for any  $u$  and adjacent  $x, y$  in  $B_r(v)$  for any nonzero restart probability  $\alpha$  or restart period  $k \geq r + 1$ .

**Lemma A.5.** *For any  $u$  and adjacent  $x, y$  in  $B_r(v)$ ,  $H(u, (x, y))$  is bounded by the following:*

$$H(u, (x, y)) \leq \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{r+1} + \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r+1}, \quad (75)$$

if the random walk restarts at  $v$  with any probability  $\alpha \in (0, 1)$ , and:

$$H(u, (x, y)) \leq k + k\Delta^{r+1}, \quad (76)$$

if the random walk restarts at  $v$  with any period  $k \geq r + 1$ .

*Proof.* The proof is almost identical to Lemma A.4, except the target  $x$  of reaching is substituted by  $(x, y)$  in the direction of  $x \rightarrow y$ , and all arguments that use the shortest path from  $u$  or  $v$  to  $x$  instead use the shortest path to  $x$  postfixed by  $x \rightarrow y$ , which adds  $+1$  to several terms in the bounds.  $\square$

We are now ready to prove Theorem 2.5.

**Theorem 2.5.** *In Theorem 2.4, if the random walk restarts at  $v$  with any nonzero probability  $\alpha$  or any period  $k \geq r + 1$ , the vertex and edge cover times of  $B_r(v)$  are always finite.*

*Proof.* The proof is inspired by the spanning tree argument of Aleliunas et al. (1979). Let us consider a depth first search of  $B_r(v)$  starting from  $v$ . We denote by  $T$  the resulting spanning tree with vertices  $V(T) = V(B_r(v))$ . We consider the expected time for a random walk starting at  $v$  to visit every vertex in the precise order visited by the depth first search by traversing each edge twice. It is clear that this upper-bounds the vertex cover time of  $B_r(v)$  starting at  $v$  (Equation 18):

$$C_V(B_r(v)) \leq \sum_{(x,y) \in E(T)} [H(x, y) + H(y, x)]. \quad (77)$$

Then, using the bounds from Lemma A.4, the property of spanning trees  $|E(T)| = |V(T)| - 1$ , and the fact that  $|V(T)| = |V(B_r(v))| \leq \Delta^r$  from bounded degree, we obtain:

$$C_V(B_r(v)) \leq 2(\Delta^r - 1) \left( \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^r + \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r} \right), \quad (78)$$

if the random walk restarts at  $v$  with any probability  $\alpha \in (0, 1)$ , and:

$$C_V(B_r(v)) \leq 2(\Delta^r - 1)(k + k\Delta^r), \quad (79)$$

if the random walk restarts at  $v$  with any period  $k \geq r$ . This completes the proof for the vertex cover time. For the edge cover time, we consider the expected time for a random walk starting at  $v$  to visit every edge in the precise order discovered<sup>6</sup> by the depth first search by traversing each edge twice. It is clear that this upper-bounds the edge cover time of  $B_r(v)$  starting at  $v$  (Equation 19):

$$C_E(B_r(v)) \leq \sum_{(x,y) \in E(B_r(v))} [H(x, (x, y)) + H(y, (y, x))]. \quad (80)$$

Then, using Lemma A.4 and the fact that  $|E(B_r(v))| \leq \Delta^{2r} - 1$  from bounded degree, we obtain:

$$C_E(B_r(v)) \leq 2(\Delta^{2r} - 1) \left( \frac{1}{\alpha} + \frac{1}{\alpha} \left( \frac{\Delta}{1-\alpha} \right)^{r+1} + \frac{1}{\alpha} \left( \frac{1}{\alpha} - 1 \right) \left( \frac{\Delta}{1-\alpha} \right)^{2r+1} \right), \quad (81)$$

if the random walk restarts at  $v$  with any probability  $\alpha \in (0, 1)$ , and:

$$C_E(B_r(v)) \leq 2(\Delta^{2r} - 1) \cdot (k + k\Delta^{r+1}), \quad (82)$$

if the random walk restarts at  $v$  with any period  $k \geq r + 1$ . This completes the proof.  $\square$

While our proof shows finite bounds for the cover times, it is possible that they can be made tighter, for instance based on Zuckerman (1991). We leave improving the bounds as a future work.

#### A.5.7 PROOF OF THEOREM 3.2 (SECTION 3.1)

We recall universal approximation of graph-level functions in probability (Definition 3.1):

**Definition 3.1.** We say  $X_\theta(\cdot)$  is a universal approximator of graph-level functions in probability if, for all invariant functions  $\phi : \mathbb{G}_n \rightarrow \mathbb{R}$  for a given  $n \geq 1$ , and  $\forall \epsilon, \delta > 0$ , there exist choices of length  $l$  of the random walk and network parameters  $\theta$  such that the following holds:

$$\text{Prob}[|\phi(G) - X_\theta(G)| < \epsilon] > 1 - \delta, \quad \forall G \in \mathbb{G}_n. \quad (83)$$

We remark that an RWNN  $X_\theta(\cdot)$  is composed of a random walk algorithm, a recording function  $q : (v_0 \rightarrow \dots \rightarrow v_l, G) \mapsto \mathbf{z}$ , and a reader neural network  $f_\theta : \mathbf{z} \mapsto \hat{\mathbf{y}} \in \mathbb{R}$ .

Intuitively, if the record  $\mathbf{z}$  of the random walk always provides complete information of the input graph  $G$ , we may invoke universal approximation of  $f_\theta$  to always obtain  $|\phi(G) - f_\theta(\mathbf{z})| < \epsilon$ , and thus  $|\phi(G) - X_\theta(G)| < \epsilon$ . However, this is not always true as the random walk may e.g. fail to visit some vertices of  $G$ , in which case the record  $\mathbf{z}$  would be incomplete. As we show below, this uncertainty leads to the probabilistic bound  $> 1 - \delta$  of the approximation.

Let us denote the collection of all possible random walk records as  $\{\mathbf{z}\} := \text{Range}(q)$ , and consider a decoding function  $\psi : \{\mathbf{z}\} \rightarrow \mathbb{G}_n$  that takes the record  $\mathbf{z} := q(v_0 \rightarrow \dots \rightarrow v_l, G)$  of a given random walk  $v_{[\cdot]}$  and outputs the graph  $\psi(\mathbf{z}) \in \mathbb{G}_n$  composed of all recorded vertices  $V(H) := \{\text{id}(v_t) : v_t \in \{v_0, \dots, v_l\}\}$  and all recorded edges  $E(H) \subset V(H) \times V(H)$ . We show the following lemma:

**Lemma A.6.** Let  $G_{\mathbf{z}}$  be the subgraph of  $G$  whose vertices and edges are recorded by  $\mathbf{z}$ . Then the graph  $\psi(\mathbf{z})$  decoded from the record  $\mathbf{z}$  is isomorphic to  $G_{\mathbf{z}}$  through the namespace  $\text{id}(\cdot)$ :

$$G_{\mathbf{z}} \stackrel{\text{id}}{\simeq} \psi(\mathbf{z}). \quad (84)$$

Furthermore, the decoded graph  $\psi(\mathbf{z})$  reconstructs  $G$  up to isomorphism, that is,

$$G \stackrel{\text{id}}{\simeq} \psi(\mathbf{z}), \quad (85)$$

if the recording function  $q(\cdot)$  and the random walk  $v_{[\cdot]}$  satisfies either of the following:

- $q(\cdot)$  uses anonymization, and  $v_{[\cdot]}$  has traversed all edges of  $G$ .
- $q(\cdot)$  uses anonymization and named neighbors, and  $v_{[\cdot]}$  has visited all vertices of  $G$ .

<sup>6</sup>We remark that depth first search discovers all edges of a graph, while not necessarily visiting all of them.

2322 *Proof.* Equation 84 is straightforward from the fact that the namespace  $\text{id}(\cdot)$  defines a bijection  
 2323 from  $V(G_{\mathbf{z}})$  to  $[|V(G_{\mathbf{z}})|]$ , and the recording function uses names  $\text{id}(v_t)$  to record vertices and edges.  
 2324 Equation 85 is satisfied when all vertices and edges of  $G$  have been recorded, i.e.  $G_{\mathbf{z}} = G$ , which  
 2325 is possible either when the random walk has traversed all edges of  $G$ , or when it has traversed all  
 2326 vertices of  $G$  and named neighbors are used to record the induced subgraph  $G[V(G)] = G$ .  $\square$

2327

2328 We further remark Markov’s inequality for any nonnegative random variable  $T$  and  $a > 0$ :

2329

$$2330 \text{Prob}[T \geq a] \leq \frac{\mathbb{E}[T]}{a}. \quad (86)$$

2331

2332 We are now ready to prove Theorem 3.2.

2333

2334 **Theorem 3.2.** *An RWNN  $X_{\theta}(\cdot)$  with a sufficiently powerful  $f_{\theta}$  is a universal approximator of*  
 2335 *graph-level functions in probability (Definition 3.1) if it satisfies either of the below:*

2336

- *It uses anonymization to record random walks of lengths  $l > C_E(G)/\delta$ .*

2337

- *It uses anonymization and named neighbors to record walks of lengths  $l > C_V(G)/\delta$ .*

2338

2339

2340

2341

*Proof.* Instead of directly approximating the target function  $\phi : \mathbb{G}_n \rightarrow \mathbb{R}$ , it is convenient to define a  
 proxy target function on random walk records  $\phi' : \{\mathbf{z}\} \rightarrow \mathbb{R}$  where  $\{\mathbf{z}\} := \text{Range}(q)$  as follows:

2342

$$2343 \phi' := \phi \circ \psi, \quad (87)$$

2344

where  $\psi : \{\mathbf{z}\} \rightarrow \mathbb{G}_n$  is the decoding function of walk records. Then, for a given  $\mathbf{z}$ , we have:

2345

$$2346 G \simeq \psi(\mathbf{z}) \implies \phi(G) = \phi'(\mathbf{z}), \quad (88)$$

2347

2348

which is because  $\phi$  is an invariant function, so  $\phi(G) = \phi(\psi(\mathbf{z})) = \phi \circ \psi(\mathbf{z}) = \phi'(\mathbf{z})$ . Then we have:

2349

$$2350 \text{Prob}[G \simeq \psi(\mathbf{z})] \leq \text{Prob}[|\phi(G) - \phi'(\mathbf{z})| = 0]. \quad (89)$$

2351

2352

We now invoke universality of  $f_{\theta}$  to approximate  $\phi'$ . If  $f_{\theta}$  is a universal approximator of functions on  
 its domain  $\{\mathbf{z}\} := \text{Range}(q)$ , for any  $\epsilon > 0$  there exists a choice of  $\theta$  such that the below holds:

2353

$$2354 |\phi'(\mathbf{z}) - f_{\theta}(\mathbf{z})| < \epsilon, \quad \forall \mathbf{z} \in \text{Range}(q). \quad (90)$$

2355

2356

Combining Equations 89 and 90, we have:

2357

$$2358 \text{Prob}[G \simeq \psi(\mathbf{z})] \leq \text{Prob}[|\phi(G) - \phi'(\mathbf{z})| + |\phi'(\mathbf{z}) - f_{\theta}(\mathbf{z})| < \epsilon]. \quad (91)$$

2359

2360

We remark triangle inequality of distances on  $\mathbb{R}$ , for a given  $\mathbf{z}$ :

2361

$$2362 |\phi(G) - f_{\theta}(\mathbf{z})| \leq |\phi(G) - \phi'(\mathbf{z})| + |\phi'(\mathbf{z}) - f_{\theta}(\mathbf{z})|, \quad (92)$$

2363

2364

which implies, for a given  $\mathbf{z}$ :

2365

$$2366 |\phi(G) - \phi'(\mathbf{z})| + |\phi'(\mathbf{z}) - f_{\theta}(\mathbf{z})| < \epsilon \implies |\phi(G) - f_{\theta}(\mathbf{z})| < \epsilon, \quad (93)$$

2367

2368

and hence:

2369

$$2370 \text{Prob}[|\phi(G) - \phi'(\mathbf{z})| + |\phi'(\mathbf{z}) - f_{\theta}(\mathbf{z})| < \epsilon] \leq \text{Prob}[|\phi(G) - f_{\theta}(\mathbf{z})| < \epsilon]. \quad (94)$$

2371

2372

Combining Equations 91 and 94, we have:

2373

$$2374 \text{Prob}[|\phi(G) - f_{\theta}(\mathbf{z})| < \epsilon] \geq \text{Prob}[G \simeq \psi(\mathbf{z})], \quad (95)$$

2375

2376

which can be written as follows:

2377

$$2378 \text{Prob}[|\phi(G) - X_{\theta}(G)| < \epsilon] \geq \text{Prob}[G \simeq \psi(\mathbf{z})]. \quad (96)$$

2379

2380

2381

2382

2383

2384

2385

We now consider the probability of the event  $G \simeq \psi(\mathbf{z})$  based on Lemma A.6. We first consider the  
 case where the recording function  $q(\cdot)$  uses anonymization. In this case,  $G \simeq \psi(\mathbf{z})$  is achieved if the  
 random walk of length  $l$  has traversed all edges of  $G$ . Let  $T_E(G, v_0)$  be the number of steps that a  
 random walk starting at  $v_0$  takes until traversing all edges of  $G$ . Since the edge cover time  $C_E(G)$  is  
 its expectation taken at the worst possible starting vertex (Equation 17), we have the following:

$$2386 \mathbb{E}[T_E(G, v_0)] \leq C_E(G), \quad \forall v_0 \in V(G), \quad (97)$$

2376 which leads to the following from Markov’s inequality (Equation 86):  
 2377

$$2378 \text{Prob}[T_E(G, v_0) < l] \geq 1 - \frac{\mathbb{E}[T_E(G)]}{l} \geq 1 - \frac{C_E(G)}{l}. \quad (98)$$

2380 For a given random walk  $v_0 \rightarrow \dots \rightarrow v_l$  and its record  $\mathbf{z}$ , the following holds:  
 2381

$$2382 T_E(G, v_0) < l \implies G \simeq \psi(\mathbf{z}), \quad (99)$$

2383 which implies the following:  
 2384

$$2385 \text{Prob}[T_E(G, v_0) < l] \leq \text{Prob}[G \simeq \psi(\mathbf{z})]. \quad (100)$$

2386 Combining Equations 96, 98, and 100, we have:  
 2387

$$2388 \text{Prob}[|\phi(G) - X_\theta(G)| < \epsilon] \geq 1 - \frac{C_E(G)}{l}. \quad (101)$$

2390 Therefore, for any  $\delta > 0$ , if we choose  $l > C_E(G)/\delta$  we would have the following:  
 2391

$$2392 \text{Prob}[|\phi(G) - X_\theta(G)| < \epsilon] > 1 - \delta. \quad (102)$$

2393 This completes the proof for anonymization. The proof is identical for the recording function  
 2394 that uses anonymization and named neighbors, except that the edge cover time is changed to the  
 2395 vertex cover time  $C_V(G)$  (Equation 16). This is because neighborhood recording automatically  
 2396 records the induced subgraph of visited vertices, thus visiting all vertices implies recording all edges,  
 2397  $G[V(G)] = G$ .  $\square$   
 2398

#### 2399 A.5.8 PROOF OF THEOREM 3.4 (SECTION 3.1) 2400

2401 **Theorem 3.4.** *An RWNN  $X_\theta(\cdot)$  with a sufficiently powerful  $f_\theta$  and any nonzero restart probability  
 2402  $\alpha$  or restart period  $k \geq r + 1$  is a universal approximator of vertex-level functions in probability  
 2403 (Definition 3.3) if it satisfies either of the below for all  $B_r(v) \in \mathbb{B}_r$ :*  
 2404

- 2405 • *It uses anonymization to record random walks of lengths  $l > C_E(B_r(v))/\delta$ .*
- 2406 • *It uses anonymization and named neighbors to record walks of lengths  $l > C_V(B_r(v))/\delta$ .*  
 2407

2408 *Proof.* The proof is almost identical to Theorem 3.2, except  $G \in \mathbb{G}_n$  are substituted by  $B_r(v) \in \mathbb{B}_r$ ,  
 2409 and the decoding function  $\psi : \{\mathbf{z}\} \rightarrow \mathbb{B}_r$  is defined to ignore all recorded vertices  $\text{id}(x)$  whose  
 2410 shortest path distance from the starting vertex  $\text{id}(v) = \text{id}(v_0) = 1$  exceeds  $r$ . The latter is necessary  
 2411 to restrict the range of the decoding function  $\psi$  to  $\mathbb{B}_r$ . In addition, any nonzero restart probability  $\alpha$   
 2412 or restart period  $k \geq r + 1$  is sufficient to make the cover times  $C_E(B_r(v))$  and  $C_V(B_r(v))$  finite  
 2413 (Theorem 2.5), thereby guaranteeing the existence of a finite choice of  $l$ .  $\square$   
 2414

#### 2415 A.5.9 PROOF OF THEOREM 3.5 (SECTION 3.2) 2416

2417 **Theorem 3.5.** *The simple RWNN outputs  $\mathbf{h}^{(l)} \rightarrow \mathbf{x}^\top \boldsymbol{\pi}$  as  $l \rightarrow \infty$ .*  
 2418

2419 *Proof.* Since  $G$  is connected and non-bipartite, the uniform random walk on it defines an ergodic  
 2420 Markov chain with a unique stationary distribution  $\boldsymbol{\pi}$ . The limiting frequency of visits on each vertex  
 2421  $v$  is precisely the stationary probability  $\pi_v$ . Since the model reads  $\mathbf{x}_{v_0} \rightarrow \dots \rightarrow \mathbf{x}_{v_l}$  by average  
 2422 pooling, the output is given by weighted mean  $\sum_v \pi_v \mathbf{x}_v$  which is  $\mathbf{x}^\top \boldsymbol{\pi}$ .  $\square$   
 2423

#### 2424 A.5.10 PROOF OF THEOREM 3.6 (SECTION 3.2) 2425

2426 **Theorem 3.6.** *Let  $\mathbf{h}_u^{(l)}$  be output of the simple RWNN queried with  $u$ . Then:*  
 2427

$$2428 \mathbb{E} \left[ \left\| \frac{\partial \mathbf{h}_u^{(l)}}{\partial \mathbf{x}_v} \right\| \right] = \frac{1}{l+1} \left[ \sum_{t=0}^l P^t \right]_{uv} \rightarrow \boldsymbol{\pi}_v \quad \text{as } l \rightarrow \infty. \quad (103)$$

*Proof.* Since the model reads  $\mathbf{x}_{v_0} \rightarrow \dots \rightarrow \mathbf{x}_{v_l}$  by average pooling, the feature Jacobian  $|\partial \mathbf{h}_u^{(l)} / \partial \mathbf{x}_v|$  is given as number of visits to the vertex  $v$  in the random walk  $v_0 \rightarrow \dots \rightarrow v_l$  starting at  $v_0 = u$ , divided by length  $l + 1$ . Let us denote the expected number of these visits by  $J(u, v, l)$ . Let  $\mathbb{1}_{v_t=v}$  be the indicator function that equals 1 if  $v_t = v$  and 0 otherwise. Then we can write  $J(u, v, l)$  as:

$$\begin{aligned} J(u, v, l) &= \mathbb{E} \left[ \sum_{t=0}^l \mathbb{1}_{v_t=v} | v_0 = u \right], \\ &= \sum_{t=0}^l \mathbb{E}[\mathbb{1}_{v_t=v} | v_0 = u]. \end{aligned} \quad (104)$$

We have used linearity of expectations for the second equality.  $\mathbb{E}[\mathbb{1}_{v_t=v} | v_0 = u]$  is the probability of being at  $v$  at step  $t$  given that the walk started at  $u$ . This probability is precisely  $[P^t]_{uv}$ . Therefore:

$$J(u, v, l) = \sum_{t=0}^l [P^t]_{uv} = \left[ \sum_{t=0}^l P^t \right]_{uv}, \quad (105)$$

which gives the equality in Equation 103. Furthermore, since  $G$  is connected and non-bipartite, the uniform random walk on it defines an ergodic Markov chain with a unique stationary distribution  $\pi$ . The limiting frequency of visits on vertex  $v$  is precisely the stationary probability  $\pi_v$ , which gives the convergence in Equation 103. This completes the proof.  $\square$

## A.6 EXTENDED RELATED WORK

**Anonymized random walks (Micali & Zhu, 2016)** The initial work on anonymization by Micali & Zhu (2016) has stated an important result, that a sufficiently long anonymized walk starting from a vertex  $v$  encodes sufficient information to reconstruct the local subgraph  $B_r(v)$  up to isomorphism. While the focus of Micali & Zhu (2016) was a probabilistic graph reconstruction algorithm that uses a set of independent anonymized walks and accesses the oracle set of all possible anonymized walks, we adopt the idea in a neural processing context to acquire universality in probability (Section 3.1). In addition, the invariance property of anonymization has rarely been noticed formally in the literature, which is our key motivation for using it, on top of being able to recover the whole graph (Section 2).

**In-depth comparison with CRaWI (Tönshoff et al., 2023)** Our approach is related to CRaWI in two key aspects: **(1)** identity and connectivity encodings of CRaWI contains analogous information to anonymization and neighborhood recording, respectively, and **(2)** CRaWI uses 1D CNNs as the reader NN. A key technical difference lies in the first part. The identity and connectivity encodings of CRaWI are defined within a fixed window size (denoted  $s$  in (Tönshoff et al., 2023)), which puts a locality constraint on the recorded information. Precisely, the encodings at step  $t$  can encode the information of the walk from step  $t - s$  to  $t$  (precisely, its induced subgraph), referred to as a walklet. The window size  $s$  is a hyperparameter that controls the expressive power of CRaWI, and this dependency makes CRaWI non-universal (Section 3.2, (Tönshoff et al., 2023)). Our choice of anonymization and neighborhood recording are not under such a local constraint, and they encode the full information of a given walk globally. This property underlies our universality results in Section 3.1, which also naturally motivates our choice of universal reader NNs, e.g. a transformer language model, which were not explicitly considered in CRaWI.

**Random walks on graphs** Our work builds upon theory of random walks on graphs, i.e. Markov chains on discrete spaces. Their statistical properties such as hitting and mixing times have been well-studied (Aleliunas et al., 1979; Lovász, 1993; Coppersmith et al., 1996; Feige, 1995; Oliveira, 2012; Peres & Sousi, 2015), and our method (Section 2) is related to vertex cover time (Aleliunas et al., 1979; Kahn et al., 1989; Ding et al., 2011; Abdullah, 2012), edge cover time (Zuckerman, 1991; Bussian, 1996; Panotopoulou, 2013; Georgakopoulos & Winkler, 2014), and improving them, using local degree information (Ikeda et al., 2009; Abdullah et al., 2015; David & Feige, 2018), non-backtracking (Alon et al., 2007; Kempton, 2016; Arrigo et al., 2019; Fasino et al., 2021), or restarts in case of infinite graphs (Dumitriu et al., 2003; McNew, 2013; Janson & Peres, 2012). Our work is also inspired by graph algorithms based on random walks, such as anonymous observation (Micali & Zhu, 2016), sublinear algorithms (Dasgupta et al., 2014; Chiericetti et al., 2016; Ben-Hamou et al., 2018;

2484 Bera & Seshadhri, 2020), and personalized PageRank for search (Page et al., 1999). While we adopt  
2485 their techniques to make our walks and their records well-behaved, the difference is that we use a  
2486 deep neural network to process the records and directly make predictions. Our method is also related  
2487 to label propagation algorithms (Zhu & Ghahramani, 2002; Zhu, 2005; Grady, 2006) that perform  
2488 transductive learning on graphs based on random walks, which we further discuss in Section 5.3.

2489  
2490 **Over-smoothing and over-squashing** Prior work on over-smoothing and over-squashing of  
2491 MPNNs (Barceló et al., 2020; Topping et al., 2022; Giovanni et al., 2023; Giraldo et al., 2023;  
2492 Black et al., 2023; Nguyen et al., 2023; Park et al., 2023; Wu et al., 2023) often make use of structural  
2493 properties of graphs, such as effective resistance (Doyle & Snell, 1984; Chandra et al., 1997), Lapla-  
2494 cian eigen-spectrum (Lovász, 1993; Spielman, 2019), and discrete Ricci curvatures (Ollivier, 2009;  
2495 Devriendt & Lambiotte, 2022). Interestingly, random walks and their statistical properties are often  
2496 closely related to these properties, indicating some form of parallelism between MPNNs and RWNNs.  
2497 This has motivated our analysis in Section 3.2, where we transfer the prior results on over-smoothing  
2498 and over-squashing of MPNNs based on these properties into the results on our approach.

2499 **Language models for learning on graphs** While not based on random walks, there have been  
2500 prior attempts on applying language models for problems on graphs (Wang et al., 2023; Chen et al.,  
2501 2023b; Zhao et al., 2023; Fatemi et al., 2023; Ye et al., 2024), often focusing on prompting methods  
2502 on problems involving simulations of graph algorithms. We take a more principle-oriented approach  
2503 based on invariance and expressive power, and thus demonstrate our approach mainly on the related  
2504 tasks, e.g. graph separation. We believe extending our work to simulating graph algorithms requires  
2505 a careful treatment (Weiss et al., 2021; Delétang et al., 2023; de Luca & Fountoulakis, 2024; Sanford  
2506 et al., 2024) and plan to investigate it as future work.

## 2507 2508 A.7 LIMITATIONS AND FUTURE WORK

2509 Our current main limitation is the cost of performing training and inference if using language models,  
2510 especially on long walks. We believe efficient fine-tuning with e.g. low-rank adaptation (Hu et al.,  
2511 2022; Chen et al., 2023a) may help overcome the issue. A question for future work is whether we  
2512 can train a language model on a large pseudocorpus of random walks (Klubicka et al., 2019; 2020)  
2513 to build a foundation model for graphs which is language-compatible. We plan to investigate this  
2514 direction in the future.

2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529  
2530  
2531  
2532  
2533  
2534  
2535  
2536  
2537