

Disobeying Directions: Switching Random Walk Filters for Unsupervised Node Embedding Learning on Directed Graphs

Ciwan Ceylan

*Division of Robotics, Perception and Learning
KTH Royal Institute of Technology*

ciwan@kth.se

Kambiz Ghoorchian

SEB Group

kambiz.ghoorchian@seb.com

Danica Kragic

*Division of Robotics, Perception and Learning
KTH Royal Institute of Technology*

dani@kth.se

Reviewed on OpenReview: <https://openreview.net/forum?id=yngjRgVA5A>

Abstract

Unsupervised learning of node embeddings for directed graphs (digraphs) requires careful handling to ensure unbiased modelling. This paper addresses two key challenges: (1) the obstruction of information propagation in random walk and message-passing methods due to local sinks, and (2) the representation of multiple multi-step directed neighbourhoods, arising from the distinction between in- and out-neighbours. These challenges are interconnected—local sinks can be mitigated by treating the graph as undirected, but this comes at the cost of discarding all directional information. We make two main contributions to unsupervised embedding learning for digraphs. First, we introduce ReachNEs (Reachability Node Embeddings), a general framework for analysing embedding models and diagnosing local sink behaviour on digraphs. ReachNEs defines the reachability filter, a matrix polynomial over normalized adjacency matrices that captures multi-step, direction-sensitive proximity. It unifies the analysis of message-passing and random walk models, making its insights applicable across a wide range of embedding methods. Second, we propose DirSwitch, a novel embedding model that resolves both local sink bias and neighbourhood multiplicity via switching random walks. These walks use directed edges for local steps, preserving directional structure, then switch to undirected edges for long-range transitions, enabling escape from local sinks and improving information dispersal. Empirical results on node classification benchmarks demonstrate that DirSwitch consistently outperforms state-of-the-art unsupervised digraph proximity embedding methods, and also serves as a flexible digraph extension for self-supervised graph neural networks.

1 Introduction

Directed graphs (digraphs) and node embeddings are both essential modelling tools for graph data. Digraphs are ubiquitous for describing networks with asymmetric relationships, such as citation networks (McCallum et al., 2000), online dating platforms (Takac & Zabojsky, 2012), financial transactions (Fire & Guestrin, 2020), brain connectomes (Winding et al., 2023), and the Internet (Faloutsos et al., 1999). Similarly, node embeddings are key to solving machine learning and data mining tasks on graphs, including node classification (Hou et al., 2023), node clustering (Wang et al., 2017b), graph alignment (Heimann et al., 2018), and link prediction (Virinchi & Saladi, 2023). Despite their combined prevalence, the field of unsupervised embedding learning for digraphs remains under-explored. Most research has focused on undirected graphs and (semi-)supervised learning, even though task-specific labelled data is often scarce (Hu et al., 2020b).

Embedding modelling for digraphs presents several challenges. A primary requirement is the ability to capture the asymmetric nature of digraphs. Conventional spectral-based approaches fall short, as asymmetric matrices do not admit eigendecomposition (Zhang et al., 2021b). As a result, it is common practice to treat digraphs as undirected, thereby discarding important relational information (Rossi et al., 2023).

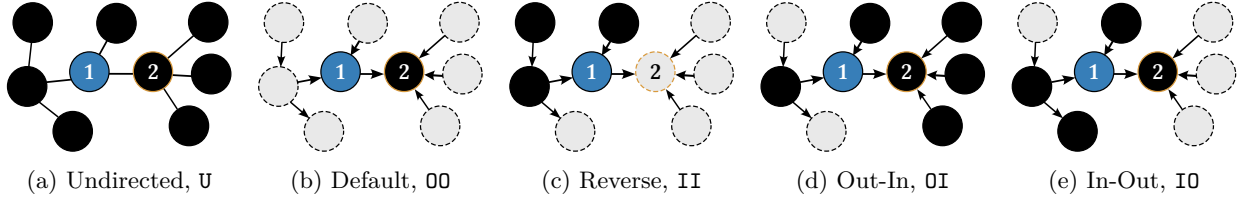


Figure 1: Visualization of local sinks and neighbourhood multiplicity. Each digraph shows a valid definition of the 2-step neighbourhood of node 1, highlighted in black. Note that each directed neighbourhood is unique, carrying independent information. (a) uses undirected edges, resulting in the maximum number of reachable nodes at the expense of the digraph structure. (b) shows the neighbourhood using the default, outgoing edge directions. The local sink, node 2, obstructs the information flow, so the neighbourhood consists of a single node. In (c), the reverse, incoming edges are followed, and (d) and (e) use mixed edge directions. (d) uses outgoing edges for the first step and incoming edges for the second, and vice versa for (e).

Recent research on digraph node embeddings has focused on adapting random walk and message-passing approaches to capture digraph asymmetry (Zhou et al., 2017; Khosla et al., 2020; Tong et al., 2020; Zhang et al., 2021b; Virinchi & Saladi, 2023; Rossi et al., 2023). These approaches rely on information propagation between nodes to generate embeddings that represent the connectivity of multi-step node neighbourhoods. To account for digraph asymmetry, these models separate propagation along outgoing (default) edges from reverse (incoming) edges, creating two sets of paired node embeddings.

While these models address digraph asymmetry, other challenges have been overlooked. This work tackles two such challenges: *local sinks* and *neighbourhood multiplicity*, visualized in Figure 1. The unidirectionality of digraph edges means information flows in only one direction. As a result, digraphs contain *local sinks*, i.e., nodes without outgoing edges. As shown in Figure 1b, these sinks obstruct information propagation and hinder exploration of multi-step neighbourhoods. This leads to a disproportionate influence of sink nodes, resulting in biased embedding models. For instance, in Figure 1b, node 1 only sees the sink node (node 2).

As highlighted in Figure 1, the multi-step neighbourhood in digraphs is not uniquely defined. This multiplicity arises from the fact that each node has two distinct sets of neighbours: out-neighbours and in-neighbours. When multi-step neighbourhoods are considered, this multiplicity is exponentially amplified, as each node can have up to 2^r distinct r -step neighbourhoods. The challenge for unsupervised node embedding models is to represent this multitude of distinct neighbourhoods as accurately as possible, since it is generally unknown a priori which neighbourhoods are key for solving downstream tasks.

In this work, we present two key contributions to address the challenges of unsupervised embedding learning for digraphs: (1) ReachNEs (Reachability Node Embeddings), a unifying framework for analysing node embedding models, and (2) DirSwitch, a flexible approach for mitigating the issues of local sinks and neighbourhood multiplicity, applicable to both random walk and message-passing embedding models.¹

ReachNEs is a mathematical framework designed to study and analyse the challenges of digraph node embedding modelling. Flexibility and analytical tractability are prioritized in its design. To achieve this, ReachNEs is built around the *reachability graph filter*, a matrix computed as a polynomial of normalized adjacency matrices, capturing asymmetric node proximity through multi-step random walk transition probabilities between node pairs in the graph.

The reachability matrix unifies random walk and message-passing embedding models. To generate random walk-based proximity embeddings (Zhou et al., 2017; Zhu et al., 2021a), matrix factorization is used to decompose the reachability matrix into low-dimensional factors. Conversely, the reachability matrix can serve as a neighbourhood smoothing filter to produce linear message-passing embeddings (Wu et al., 2019). This versatility makes the insights drawn from ReachNEs analysis applicable to a wide range of embedding techniques.

The ReachNEs framework is also flexible with respect to modelling different neighbourhood scales, supporting multi-scale embeddings, which are crucial in the unsupervised setting (Rozemberczki et al., 2021). This control is provided by the coefficients of the reachability matrix’s defining polynomial, which are interpreted as random walk length probabilities.

¹Our source code: <https://github.com/ciwanceylan/dirswitch-experiments-tmlr2025.git>

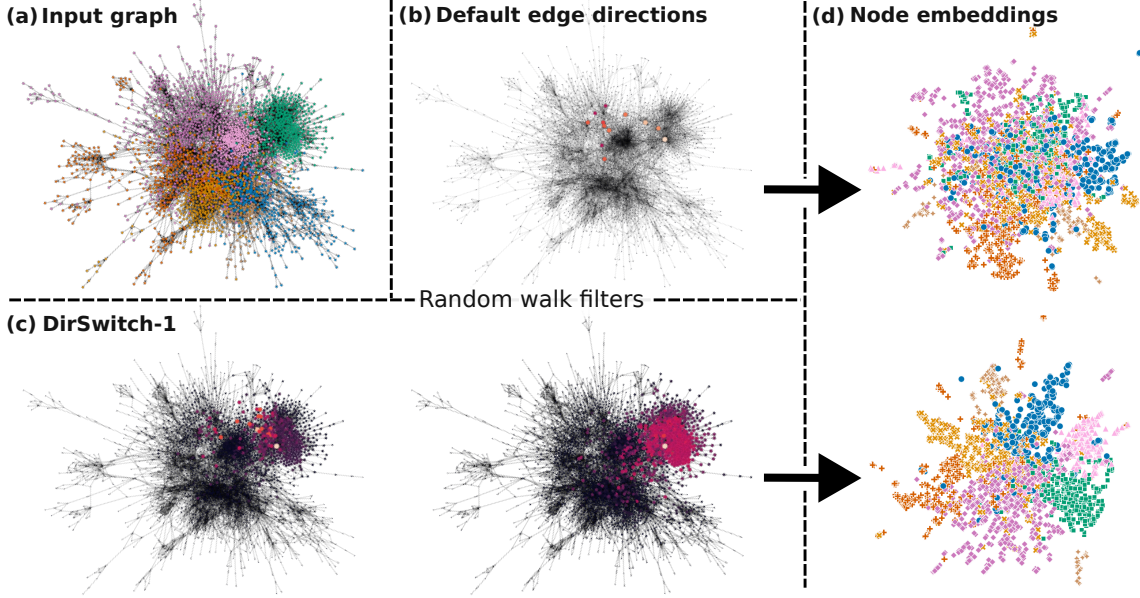


Figure 2: A qualitative illustration of DirSwitch. (a): The Cora-ML citation graph, with nodes coloured by class label (for visualization purposes only; not used in the model). (b) and (c): Random walk transition probabilities from a specific initialization node. The graph in (b) shows the default edge directions, while the graphs in (c) employ our DirSwitch-1 approach. DirSwitch-1 performs one directed step using either default (*left*) or reverse (*right*) edge directions, followed by undirected transitions. (d): UMAP (Sainburg et al., 2021) visualizations of the resulting node embeddings. The DirSwitch-1 embedding space better reflects the spatial ground truth class separation, demonstrating its improved neighbourhood smoothing on digraphs.

We use ReachNEs to quantify the effects of local sinks on information flow and to mathematically characterize the multiplicity of distinct directed neighbourhoods in digraphs. This analysis reveals a tension between these two issues. The local sink problem can be easily resolved by treating digraphs as undirected, allowing information to flow freely within each weakly connected component. However, this comes at the cost of losing the distinction between directed neighbourhoods, reducing the expressivity of the embedding model. Therefore, we seek a model that can address local sinks while preserving the capture of directed neighbourhoods.

Our second contribution, DirSwitch, achieves this balance. The core idea of DirSwitch is to decouple local, short-range random walk behaviour from global, long-range behaviour. Specifically, DirSwitch preserves local edge direction expressivity by performing directed random walks in the initial steps. It then switches to treating the graph as undirected, allowing information to propagate freely and escape local sinks. The effectiveness of this approach is illustrated in Figure 2.

We empirically evaluate DirSwitch both within and beyond the ReachNEs framework. First, we verify that DirSwitch both improves information propagation on digraphs compared to fully directed approaches, without compromising the representation of local neighbourhood multiplicity.

We then demonstrate that this results in higher-quality embeddings and increased effectiveness by benchmarking DirSwitch on 14 standard node classification benchmark datasets. Compared to ReachNEs models that do not distinguish between local and global random walk behaviours, DirSwitch either achieves the highest accuracy or performs within one standard deviation of the best result across all datasets. This holds true for both the random walk-based proximity embedding setting and the message-passing setting. Importantly, DirSwitch performs well on both homophilic and heterophilic datasets (Zhu et al., 2020). This is desirable as it implies a broad applicability of the DirSwitch in real-world unsupervised settings where the validity of a homophily assumption is in question (Wang et al., 2023).

To further assess DirSwitch’s practical effectiveness, we evaluate it against recent state-of-the-art digraph methods. On six proximity embedding benchmarks without node attributes, DirSwitch achieves the highest accuracy on all datasets except one, where it performs within one standard deviation of the best.

For message-passing embeddings, we apply DirSwitch as a digraph extension of the GraphSAGE graph neural network (GNN) model (Hamilton et al., 2017) and integrate it with self-supervised loss functions, which have previously been limited to undirected graphs (Zhang et al., 2021a; Hou et al., 2023). This demonstrates how insights from ReachNEs can be leveraged to enhance embedding models more broadly. As a baseline, we use the recent digraph GNN extension by Rossi et al. (2023), designed for supervised learning. Again, DirSwitch achieves higher accuracy on all but one dataset, underscoring the need for additional considerations when modelling digraphs in the unsupervised setting.

2 ReachNEs: Reachability Node Embeddings

This section introduces the ReachNEs framework for learning and analysing node embeddings in digraphs. It consists of two components: the reachability graph filter, which represents multi-step node relations, and reduction methods, which extract embeddings. Section 2.1 derives reachability filters from random walks, while Section 2.2 describes the walk length distributions that determine the filter’s smoothing properties. Finally, Section 2.3 presents reduction methods that unify message-passing and proximity-based embeddings.

2.1 The random walk reachability filter

Let $\mathcal{G} = (\mathbb{N}, \mathbb{M})$ be a directed, potentially weighted graph without self-loops, where \mathbb{N} is the set of nodes and \mathbb{M} is the set of edges. The graph has $n = |\mathbb{N}|$ nodes and $m = |\mathbb{M}|$ edges. The adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is defined as follows, along with the out-degree and in-degree of a node j :

$$A_{i,j} = \begin{cases} W_{j \rightarrow i} & \text{if } (j, i) \in \mathbb{M}, \\ 0 & \text{otherwise,} \end{cases} \quad \deg_0(j) = \sum_{k=1}^n A_{k,j}, \quad \deg_I(j) = \sum_{k=1}^n A_{j,k}. \quad (1)$$

Here, $W_{j \rightarrow i} \in \mathbb{R}_{\geq 0}$ is a nonnegative weight associated with the edge (j, i) . Note that the columns of \mathbf{A} correspond to out-neighbourhoods, while the rows correspond to in-neighbourhoods.

We treat undirected graphs as a special case of the above definitions, where $A_{i,j} = A_{j,i}$ for all edges in \mathbb{M} . Consequently, a digraph can be transformed into its corresponding undirected graph by adding the reverse edge (j, i) to \mathbb{M} for every edge $(i, j) \in \mathbb{M}$, forming the edge set $\mathbb{M}_{\text{undir}}$. In terms of the adjacency matrix, this symmetrization is achieved by summing \mathbf{A} and its transpose \mathbf{A}^\top , $\mathbf{A}_{\text{undir}} = \mathbf{A} + \mathbf{A}^\top$.

Next, we define three *random walk normalized adjacency matrices*, denoted as \mathbf{A}_0 , \mathbf{A}_I , and \mathbf{A}_U , which we collectively refer to as $\mathbf{A}_* \in [0, 1]^{n \times n}$. Each of these matrices is a column-stochastic matrix (Horn & Johnson, 2012, Ch. 8.7) and serves as a random walk state transition matrix. Specifically, given a probability state vector $\mathbf{p}^{(k)} \in [0, 1]^n$, where $p_i^{(k)}$ represents the probability of a random walker being at node i after k steps, the state vector after $k + 1$ steps is computed as $\mathbf{p}^{(k+1)} = \mathbf{A}_* \mathbf{p}^{(k)}$.

Each normalized adjacency matrix defines a different random walk behaviour with respect to the edge directions. For \mathbf{A}_0 , a random walker follows the default edge directions, transitioning from a node to its neighbours along outgoing edges. Conversely, transitions using \mathbf{A}_I follow the *reverse* edge directions, meaning the walker traverses incoming edges. For \mathbf{A}_U , the edges are treated as undirected, allowing the random walker to move across both outgoing and incoming edges. These behaviours are formally defined as follows:

$$A_{0,i,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg_0(j) = 0, \\ \frac{A_{i,j}}{\deg_0(j)} & \text{otherwise,} \end{cases} \quad A_{I,i,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg_I(j) = 0, \\ \frac{A_{j,i}}{\deg_I(j)} & \text{otherwise,} \end{cases}$$

$$A_{U,i,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg(j) = 0, \\ \frac{A_{\text{undir},i,j}}{\deg(j)} & \text{otherwise,} \end{cases} \quad \deg(i) = \sum_{k=1}^n A_{\text{undir},k,i}.$$

To ensure column stochasticity, diagonal elements are set to 1 for nodes without legal edges to follow. Such nodes are referred to as *sink nodes*, as a random walker gets stuck on these nodes. For digraphs without sink nodes, the above definitions can be expressed compactly in matrix form:

$$\mathbf{A}_0 = \mathbf{A} \mathbf{D}_0^{-1}, \quad \mathbf{A}_I = \mathbf{A}^\top \mathbf{D}_I^{-1}, \quad \mathbf{A}_U = \mathbf{A}_{\text{undir}} \mathbf{D}^{-1},$$

where \mathbf{D}_0 , \mathbf{D}_I , and \mathbf{D} are diagonal matrices of the respective node degrees.

Each \mathbf{A}_* describe immediate neighbourhoods, which can be extended to capture multi-step connectivity through matrix powers. Specifically, the j th column of the power matrix \mathbf{A}_*^k encodes the transition probabilities of a random walk of length k initialized at node j . This follows directly from the state transition formula $\mathbf{p}^{(k+1)} = \mathbf{A}_* \mathbf{p}^{(k)}$. Accordingly, we denote element (i, j) of \mathbf{A}_*^k as $P(j \rightsquigarrow i|k)$, representing the probability of transitioning from node j to node i in exactly k steps.

Each matrix power \mathbf{A}_*^k provides a snapshot of transition probabilities at a specific walk length k . However, any fixed k inherently restricts the description of node neighbourhoods and introduces potential biases. For example, setting $k = 2$ might fail to capture immediate neighbours if no 2-step paths exist between them, resulting in zero transition probabilities for such nodes.

To mitigate the bias and sensitivity introduced by single walk lengths, we instead impose a probability distribution over the walk length: $P_w(k)$. Using this distribution, the node transition probabilities are computed by marginalizing over the walk length. This leads to the definition of the *reachability matrix* $\mathbf{R} \in \mathbb{R}^{n \times n}$, where each element $R_{i,j}$ represents the probability of transitioning from node j to node i , with the walk length sampled from $P_w(k)$:

$$\mathbf{R}(\mathbf{A}_*; P_w) = \sum_{k=0}^{\infty} P_w(k) \mathbf{A}_*^k, \quad R_{i,j}(\mathbf{A}_*; P_w) = P(j \rightsquigarrow i | P_w) = \sum_{k=0}^{\infty} P_w(k) P(j \rightsquigarrow i | k). \quad (2)$$

Figure 3a shows an example of the reachability matrix $\mathbf{R}(\mathbf{A}_u; P_w)$ using the U.S. political blogs graph (Adamic & Glance, 2005). Each column represents the transition probabilities from a random walk starting at the corresponding node. The two prominent blocks reflect the political alignments in the dataset. Figure 3b visualizes one such column vector, overlaid on a ridiculogram of the graph. The separation between the two political communities appears clearly as two clusters.

2.2 Walk length distributions

As defined in Equation 2, the walk length distribution $P_w(k)$ controls the weighting of each adjacency power \mathbf{A}_*^k , thereby shaping the filter’s smoothing behaviour. While ReachNEs is compatible with any such distribution, in this work we focus on four commonly used ones: geometric, binomial, Poisson, and uniform.

Table 1 lists the corresponding probability mass functions (pmfs), each parameterized by the expected walk length $\tau = \mathbb{E}[P_w]$. Figure 3c visualizes the pmfs for $\tau = 2$. We briefly describe each distribution below; further analysis appears in Appendix C.

The geometric distribution is central to the PageRank algorithm (Page et al., 1999), where $\alpha = \frac{\tau}{\tau+1}$ denotes the probability of continuing a walk, and $1 - \alpha$ the restart probability. It is monotonically decreasing for all τ , and remains mode-centred at $k = 0$, even as it flattens with increasing τ . Owing to PageRank’s influence, the geometric distribution is widely used in embedding methods (Zhou et al., 2017; Yan et al., 2024).

The binomial distribution has a finite maximum walk length K , where $\alpha = \frac{\tau}{K}$ denotes the probability of stepping, and $1 - \alpha$ the probability of remaining in place. As discussed in Appendix C.5, this distribution naturally arises when self-loops are added to graphs, a common practice in graph convolutional networks (Kipf & Welling, 2017; Wu et al., 2019).

Table 1: Overview of the studied random walk length distributions $P_w(k; \tau)$, where k represents the walk length variable. Each distribution is parameterized by the average walk length τ . For readability, we use α as a substitution variable in the geometric and binomial distributions. The binomial distribution includes an additional parameter, K , representing the maximum walk length. Distribution plots are shown in Figure 3c.

NAME	GEOMETRIC	BINOMIAL	POISSON	UNIFORM
ABBREVIATION	Geom	Binom	Pois	\mathcal{U}
$P_w(k; \tau)$	$(1 - \alpha)\alpha^k, \quad \alpha = \frac{\tau}{\tau + 1}$	$\binom{K}{k}(1 - \alpha)^{K-k}\alpha^k, \quad \alpha = \frac{\tau}{K}$	$e^{-\tau} \frac{\tau^k}{k!}$	$\begin{cases} \frac{1}{2\tau + 1} & \text{if } k \in [0, 2\tau], \\ 0 & \text{otherwise} \end{cases}$

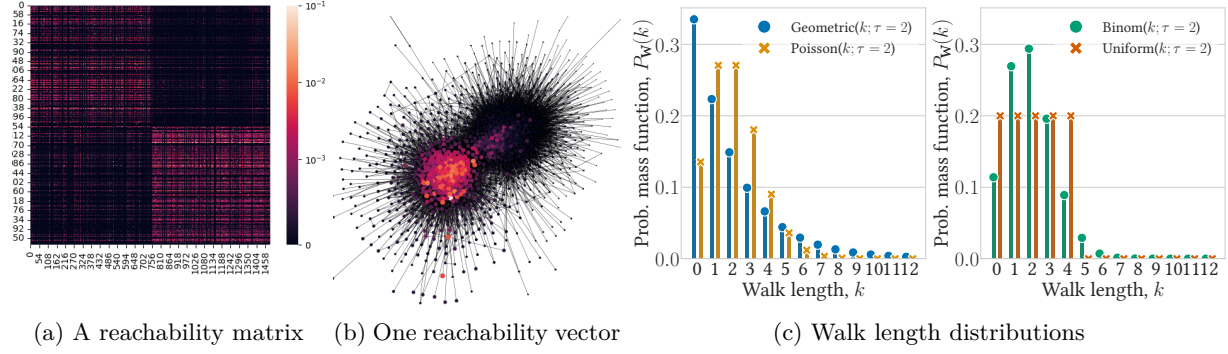


Figure 3: (a) The reachability matrix $R(\mathbf{A}_U; \text{Pois}(\tau = 2))$ for the Polblogs graph. Each column encodes the transition probabilities of a random walk starting at the corresponding node. (b) A single column of R visualized on a ridculogram of the graph, illustrating the locality of the reachability distribution. (c) The four walk length distributions used in this work, evaluated at $\tau = 2$: Geometric, Poisson, Binomial, and Uniform. These distributions weight the contribution of different walk lengths in the computation of R .

In the limit $K \rightarrow \infty$, the binomial distribution converges to the Poisson distribution. The Poisson distribution has been used in structural embedding methods (Donnat et al., 2018; Zhu et al., 2021a; Ceylan et al., 2022) and connects reachability to heat diffusion on graphs.

Unlike the geometric distribution, both the binomial and Poisson distributions tend toward non-informative filters as $\tau \rightarrow \infty$. Their modes shift with τ , emphasizing larger k values, which leads to over-smoothing: columns of R become nearly uniform, reducing embedding discriminability. This effect, analysed in detail in Appendix C, plays a key role in interpreting some of our experimental results.

Finally, the *uniform distribution* is widely used in node embedding models, as it naturally arises from common negative sampling schemes (Perozzi et al., 2014; Qiu et al., 2018; Chanpuriya & Musco, 2020).

2.3 Reachability reduction into node embeddings

To complete the ReachNEs framework, the pairwise multi-step relationships encoded in the reachability matrix R must be transformed into node embeddings. We refer to these transformations as *reduction methods*, as they reduce the $n \times n$ matrix R into a lower-dimensional embedding matrix $Z \in \mathbb{R}^{n \times p}$.

The reduction method serves two main purposes. First, it projects nodes into a lower-dimensional space ($p < n$), improving computational efficiency for downstream tasks and alleviating the curse of dimensionality (Hastie et al., 2009, Ch. 2.5). Second, the choice of reduction method encodes specific equivalence relationships into the embeddings (Zhu et al., 2021a). We explore two reduction approaches: proximity embeddings and message-passing embeddings, which capture structural and automorphic equivalence respectively.

2.3.1 Proximity embedding reduction and structural equivalence

Proximity embeddings capture mutual node connectivity and closeness, and are often referred to as *positional embeddings* in the context of undirected graphs (Zhu et al., 2021a). However, the asymmetric nature of digraphs makes the concept of “position” in a symmetric, Euclidean sense inappropriate. Instead, we define node proximity asymmetrically using the reachability matrix $R(\mathbf{A}_*; \mathbf{P}_w)$, where the element $R_{i,j}$ quantifies the closeness of node j to node i . The walk length distribution \mathbf{P}_w controls the resolution of this proximity.

To derive proximity embeddings from R , we employ Singular Value Decomposition (SVD) (Golub & Van Loan, 2013, Ch. 2.4), a widely used technique in proximity embedding methods (Qiu et al., 2018; Zhu et al., 2021a). Let $\mathbf{U}, \Sigma, \mathbf{V}^\top = \text{SVD}(R)$ represent the SVD of R , and let $\mathbf{U}_{:,q}$, $\Sigma_{q,q}$, and $\mathbf{V}_{:,q}$ denote the q -truncation of \mathbf{U} , Σ , and \mathbf{V} , respectively, where $q = p/2$. The proximity embeddings are then defined as:

$$\mathbf{Z} = [\mathbf{Z}_U \quad \mathbf{Z}_V] \in \mathbb{R}^{n \times 2q}, \quad \mathbf{Z}_U = \mathbf{U}_{:,q} \sqrt{\Sigma_{q,q}} \in \mathbb{R}^{n \times q}, \quad \mathbf{Z}_V = \mathbf{V}_{:,q} \sqrt{\Sigma_{q,q}} \in \mathbb{R}^{n \times q}. \quad (3)$$

The proximity embeddings capture node proximity in two key ways. First, the inner product of the left and right embeddings approximates the reachability matrix: $\mathbf{R} \approx \mathbf{Z}_U \mathbf{Z}_V^\top = \mathbf{U}_{:,q} \boldsymbol{\Sigma}_{:,q,q} \mathbf{V}_{:,q}^\top$. Second, the Euclidean distance between embeddings in \mathbf{Z} reflects the overlap of multi-step neighbourhoods. The latter is particularly important for downstream tasks such as clustering, where embedding distances are often utilized.

The following equality (derived in Appendix A.3) formalizes the relationship between reachability similarity and embedding distance:

$$\|\mathbf{R}_{i,:} - \mathbf{R}_{j,:}\|_2^2 + \|\mathbf{R}_{:,i} - \mathbf{R}_{:,j}\|_2^2 = \left\| (\mathbf{Z}_{i,:} - \mathbf{Z}_{j,:}) \sqrt{\hat{\boldsymbol{\Sigma}}} \right\|_2^2, \quad \hat{\boldsymbol{\Sigma}} = \begin{bmatrix} \boldsymbol{\Sigma}_{:,q,q} & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{:,q,q} \end{bmatrix}. \quad (4)$$

The left-hand side measures the difference between rows and columns of \mathbf{R} for nodes i and j , quantifying the overlap of their multi-step neighbourhoods. The right-hand side represents the distance between their respective proximity embeddings in \mathbf{Z} . This equivalence implies that, for $q \leq \text{rank}(\mathbf{R})$, two nodes have identical embeddings in \mathbf{Z} if and only if their multi-step neighbourhoods perfectly overlap.

This property is particularly valuable for distance-based downstream tasks. However, its validity relies on the first term in the definition of \mathbf{R} in Equation 2, i.e., $\mathbf{P}_w(k=0)\mathbf{I}_n$, being zero. This term introduces a constant, non-neighbourhood-dependent contribution that prevents structurally equivalent nodes from having identical embeddings. To address this, we compute \mathbf{R} for proximity embeddings using a shifted walk-length distribution; meaning that the coefficients in Equation 2 are replaced with $\mathbf{P}_w(k-1)$, using $\mathbf{P}_w(-1) = 0$.

Another modification to the definition of \mathbf{R} is the application of the elementwise thresholding function $f(R_{i,j}) = \log(\max(nR_{i,j}, 1))$ before performing the SVD. This contrast-enhancing function amplifies the difference between elements where $R_{i,j} < \frac{1}{n}$ and those where $R_{i,j} \in \left[\frac{1}{n}, \frac{1}{\log n}\right]$. This type of contrast enhancement has been shown to yield practical benefits for undirected graphs (Qiu et al., 2018; Chanpuriya & Musco, 2020; Zhu et al., 2021a), and we find that these advantages extend to digraphs as well.

Structural equivalence. Two nodes i and j are *structurally equivalent* if they have identical in- and out-neighbours, i.e., $\mathbf{A}_{i,:} = \mathbf{A}_{j,:}$ and $\mathbf{A}_{:,i} = \mathbf{A}_{:,j}$ (Borgatti & Everett, 1992). In undirected graphs, this implies that structurally equivalent nodes are at most two steps apart, highlighting the close connection between structural equivalence and node proximity.

Proximity ReachNEs embeddings preserve this equivalence under mild conditions. If nodes i and j are structurally equivalent and have nonzero in- and out-degrees, and if the walk length distribution satisfies $\mathbf{P}_w(k=0) = 0$, then $\mathbf{Z}_{i,:} = \mathbf{Z}_{j,:}$ holds as long as $q \leq \text{rank}(\mathbf{R})$. The reason is straightforward: structural equivalence implies that the i -th and j -th rows and columns of \mathbf{A}_*^k are identical for all $k \geq 1$, and thus also in $\mathbf{R} = \sum_{k \geq 1} \mathbf{P}_w(k) \mathbf{A}_*^k$. It follows from Equation 4 that the embeddings must be equal, since the diagonal elements of $\hat{\boldsymbol{\Sigma}}$ are positive under the assumption $q \leq \text{rank}(\mathbf{R})$.

An exception occurs for sink nodes. When a node lacks outgoing edges, the normalization of \mathbf{A}_0 sets the corresponding diagonal entry to 1, violating the condition that $\mathbf{A}_{*,i,:} = \mathbf{A}_{*,j,:}$. As a result, structurally equivalent sink nodes (with identical incoming edges) can receive different embeddings. The same issue applies under $\mathbf{A}_\mathbf{I}$ for nodes without incoming edges.

This limitation is shared by many random-walk-based proximity embeddings (Perozzi et al., 2014; Grover & Leskovec, 2016; Zhou et al., 2017; Qiu et al., 2018; Khosla et al., 2020). A simple workaround is to assign the zero vector $\mathbf{0}$ to all zero-degree sink nodes in post-processing. Whether this improves downstream performance remains an open research question.

2.3.2 Message-passing embedding reduction and automorphic equivalence

Message-passing embedding reduction requires an initial set of node representations, i.e., a matrix of node attribute vectors $\mathbf{X} \in \mathbb{R}^{n \times d}$. These attributes typically supplement the graph structure, such as text embeddings of paper abstracts in citation graphs (Hu et al., 2020a). Alternatively, \mathbf{X} may consist of graph-derived properties, such as node degrees or local clustering coefficients.

The ReachNEs message-passing embeddings are smoothed versions of the node attribute vectors, with the columns of the reachability matrix \mathbf{R} defining the multi-step smoothing filter over each node’s local neighbourhood. Thus, initial d -dimensional message-passing embeddings are straightforwardly obtained via the matrix multiplication $\mathbf{Z} = \mathbf{R}^\top \mathbf{X}$. To further reduce the dimensionality to p dimensions, techniques like PCA (Murphy, 2012, Ch. 12.2) can be applied, yielding the embeddings $\mathbf{Z} = \text{PCA}(\mathbf{R}^\top \mathbf{X})$.

Connection to message-passing. Message-passing is commonly used in the context of graph neural networks (Gilmer et al., 2017). In typical descriptions, each node receives messages from its neighbours in the form of embedding vectors, which are aggregated and then used to update its own embedding representation. Given this, it may not be immediately clear how the linear model $\mathbf{R}^\top \mathbf{X}$ relates to message-passing.

To establish this connection, we expand $\mathbf{R}^\top \mathbf{X}$ by bringing the transpose inside the sum of Equation 2:

$$\mathbf{R}^\top \mathbf{X} = \left(\sum_{k=0}^{\infty} \mathbf{P}_w(k) \mathbf{A}_*^k \right)^\top \mathbf{X} = \sum_{k=0}^{\infty} \mathbf{P}_w(k) (\mathbf{A}_*^k)^\top \mathbf{X} = \sum_{k=0}^{\infty} \mathbf{P}_w(k) \mathbf{A}_*^\top{}^k \mathbf{X}. \quad (5)$$

Next, we let $\mathbf{H}^{(0)} = \mathbf{X}$ and $\mathbf{Z}^{(0)} = \mathbf{P}_w(0) \mathbf{X}$, and express the sum in Equation 5 as an iterative algorithm:

$$\text{For } k \in \{1, 2, \dots\}, \quad \begin{cases} \mathbf{H}^{(k)} &= \mathbf{A}_*^\top \mathbf{H}^{(k-1)}, \\ \mathbf{Z}^{(k)} &= \mathbf{Z}^{(k-1)} + \mathbf{P}_w(k) \mathbf{H}^{(k)}. \end{cases} \quad (6)$$

This iterative formulation connects directly to message-passing. To clarify the connection, consider the update $\mathbf{H}^{(k)} = \mathbf{A}_*^\top \mathbf{H}^{(k-1)}$ for the undirected random walk matrix. Assuming no sink nodes and uniform edge weights for simplicity, the process can be expressed as:

$$\mathbf{A}_*^\top \mathbf{H} = \mathbf{D}^{-1} \mathbf{A}_{\text{undir}}^\top \mathbf{H}, \quad [\mathbf{A}_*^\top \mathbf{H}]_{i,:} = \frac{1}{\deg(i)} \sum_{j:(j,i) \in \mathbb{M}_{\text{undir}}} \mathbf{H}_{j,:}. \quad (7)$$

The right-hand side describes *mean neighbourhood aggregation*, where messages (represented by \mathbf{H}) are averaged over a node’s immediate neighbours. This mirrors the aggregation function used in GraphSAGE (Hamilton et al., 2017, Alg. 1). Similarly, \mathbf{A}_0^\top and \mathbf{A}_1^\top aggregate messages over outgoing and incoming neighbourhoods, respectively, as in the directed GraphSAGE extension proposed by Rossi et al. (2023). Full mathematical expressions for the directed case are provided in Appendix A.2.

Automorphic equivalence. Two nodes i and j are *automorphically equivalent* if there exists a permutation matrix $\mathbf{P}_\pi \in \{0,1\}^{n \times n}$ such that $\mathbf{A} = \mathbf{P}_\pi \mathbf{A} \mathbf{P}_\pi^\top$ and $\mathbf{P}_{\pi j,i} = 1$. Intuitively, this means that i and j are indistinguishable in terms of all graph-theoretic properties, such as degree, centrality, and clustering coefficients, and differ only in their labels. Automorphic equivalence is implied by structural equivalence, but not vice versa. Notably, automorphically equivalent nodes can be arbitrarily far apart in the graph, or even belong to disconnected components.

Message-passing ReachNEs preserves automorphic equivalence, provided that the node attributes are also invariant under the automorphism. Specifically, if $\mathbf{X} = \mathbf{P}_\pi \mathbf{X}$, which holds whenever \mathbf{X} consists of graph-derived features, then the embeddings satisfy $\mathbf{P}_\pi \mathbf{Z} = \mathbf{Z}$. In this case, automorphically equivalent nodes receive identical message-passing embeddings. A proof and further discussion are provided in Appendix B.

3 Local sinks and directed neighbourhood multiplicity

In this section, we use the ReachNEs framework to study local sinks and neighbourhood multiplicity. We examine how local sinks obstruct information propagation and quantify their effects using entropy. We formalize the multiplicity of directed neighbourhoods within ReachNEs and discuss its implications for embedding model expressivity. Finally, we explore how expressivity can be improved through embedding concatenation.



(a) Cora-ML. *Left: Default (D), Right: Undirected (U)* (b) Fly Larva. *Left: Default (D), Right: Undirected (U)*

Figure 4: Visualization of transition probability accumulation in local sinks. Node size and colour indicate the magnitude of reachability, with black denoting zero. (a) shows the Cora-ML citation graph using default directed edges (D, left) and undirected edges (U, right). In the directed case, reachability is concentrated in a small set of sink nodes, whereas the undirected version yields more uniform coverage of the multi-step neighbourhood. (b) shows similar behaviour for the Fly Larva dataset: reachability under D accumulates disproportionately in sinks, while U promotes broader dispersal.

3.1 Local sink analysis

As seen in Section 2, the elements of the reachability matrix \mathbf{R} represent random walk transition probabilities, with each column describing a smoothing over multi-step neighbourhoods for a given node. However, local sinks in digraphs can significantly disrupt this smoothing process.

As highlighted in Figure 1b, local sinks are nodes, or small sets of nodes, without directed paths to the rest of the graph. Random walks initiated at sink nodes remain trapped, preventing them from gathering additional neighbourhood information. Moreover, transition probabilities of walks starting at other nodes tend to accumulate in sinks, as random walkers can enter but cannot escape. This accumulation disproportionately amplifies the influence of sink nodes on the reachability filter, introducing bias into the embedding model.

Figure 4 illustrates these effects using ridiculograms of real-world graphs. In Figure 4a, the Cora-ML citation graph shows strong accumulation in sink nodes when using default edge directions, i.e., $\mathbf{R}(\mathbf{A}_D)$, while the undirected variant, i.e., $\mathbf{R}(\mathbf{A}_U)$, achieves broader coverage over the multi-step neighbourhood. Figure 4b shows a similar pattern in the denser Fly Larva brain connectome (Winding et al., 2023).

The columns of the reachability matrix reflect how local sinks restrict information propagation. Consider a node j , whose random walk probabilities are given by the column $\mathbf{R}_{:,j}$. If j is a sink node, or if its transition probabilities are largely absorbed by a sink, the reachability will be concentrated in a few nonzero values in $\mathbf{R}_{:,j}$; a property we refer to as low *dispersal*. Conversely, high dispersal indicates that probabilities are more evenly distributed across j 's multi-step neighbourhood.

This notion of dispersal is naturally captured by *Shannon entropy* (Cover & Thomas, 2005, Ch. 2.1). We define the reachability entropy of node j as $H(j; \mathbf{R}) = -\sum_{i=1}^n R_{i,j} \log_2 R_{i,j}$. Entropy is always non-negative, reaching zero when all probability is concentrated in a single node. Conversely, it attains its maximum value, $\log_2 n$, when probabilities are uniformly distributed across all n nodes. Thus, higher entropy corresponds to reduced bias toward sink nodes and improved dispersal.

In Figure 5, we plot the entropy for the Cora-ML and Fly Larva graphs. The y-axes represent entropy, $H(j; \mathbf{R})$, while the x-axes correspond to nodes j , sorted by entropy values. We use $P_w = \text{Pois}(\tau)$, with each coloured line representing a distinct τ value. The dashed line indicates the entropy of the asymptotic reachability matrix as $\tau \rightarrow \infty$ for the undirected graph (see Appendix C.1).

Transitioning from the directed to the undirected graph yields a notable increase in entropy for Cora-ML (cf. Figure 5a and Figure 5b), reflecting improved dispersal in the absence of sink nodes. The effect of the neighbourhood scale parameter τ also differs: in the undirected case, entropy grows steadily toward its asymptotic limit, while in the directed case, it eventually declines as probability mass accumulates in sink nodes.

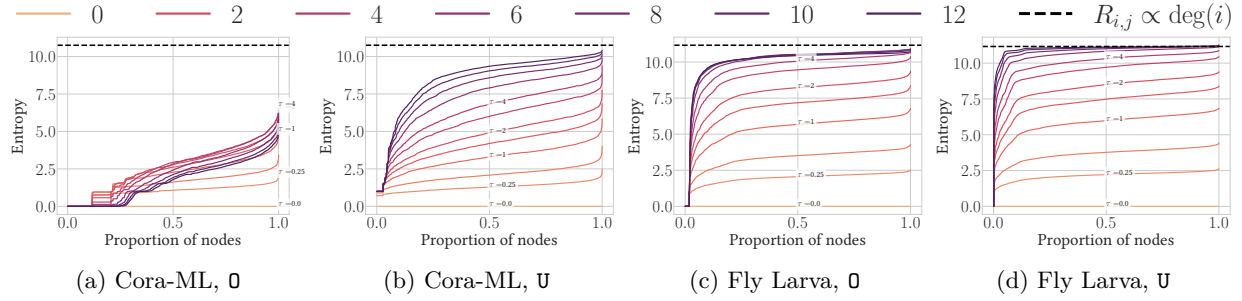


Figure 5: Reachability under $P_w = \text{Pois}(\tau)$, using A_0 and A_U . The x-axes correspond to nodes sorted by their entropy value. The colours represent various values of τ . The black dashed line indicates the entropy of a reachability distribution proportional to the node degrees, which corresponds to the limiting and uninformative distribution as $\tau \rightarrow \infty$. See Appendix C.1 for further details.

For the Fly Larva graph, entropy differences between the directed and undirected cases are smaller, as expected given the graph’s higher density. Still, entropy is consistently higher in the undirected case, with the largest gap appearing at the highest τ , as reachability in the directed graph again concentrates in sinks.

It is important to note that high dispersal alone does not guarantee embedding quality. This becomes evident when considering that maximum entropy is achieved when $R_{i,j} = \frac{1}{n}$, meaning smoothing is performed uniformly across the graph. While this maximizes dispersal, it renders the embeddings uninformative. Instead, dispersal must be complemented by another key property, which we discuss in the next section: expressivity.

3.2 Embedding model expressivity and neighbourhood multiplicity

As discussed in Section 2.3, ReachNEs respects node equivalence under mild conditions: proximity ReachNEs yields identical embeddings for structurally equivalent nodes, and message-passing ReachNEs does so for automorphically equivalent ones. However, this alone is not sufficient for producing high-quality embeddings. We also require model *expressivity*, which we define as the ability to *distinguish nonequivalent nodes*.

A fully expressive model would assign identical embeddings to equivalent nodes, and distinct embeddings to nonequivalent ones. Achieving this for automorphic equivalence is particularly challenging; in fact, any model that perfectly encodes automorphic equivalence would effectively solve the graph isomorphism problem (Grohe & Schweitzer, 2020). The expressivity of message-passing models is therefore an active area of research (Xu et al., 2019; Morris et al., 2019) (see Section 5.4 for further discussion).

Structural equivalence is simpler to handle. Since structurally equivalent nodes have identical rows and columns in the adjacency matrix A , one can construct a fully expressive, albeit impractical, embedding matrix by concatenating A and its transpose: $Z = [A \ A^T]$. This model captures structural distinctions perfectly but ignores multi-step relationships and produces a prohibitively large embedding space. Consequently, expressivity remains a key concern for proximity ReachNEs, as both the random-walk smoothing in R and the information loss from SVD reduction can limit the model’s ability to distinguish between nonequivalent nodes.

Model expressivity is especially important for digraphs due to the existence of *neighbourhood multiplicity*. As illustrated in Figure 1, the neighbourhood of a node in a digraph is not uniquely defined. At each step of a random walk, the walker can follow either incoming or outgoing edges, leading to 2^K possible neighbourhoods for a walk of length K . Embedding models that fail to capture the distinctions among these directed neighbourhoods inherently suffer from reduced expressivity.

We formalize this idea using a truncated K -step reachability matrix. Let the *edge orientation specifier* be a string σ of length K , where each character $\sigma_l \in \{U, 0, I\}$ indicates whether the l -th step uses undirected, outgoing, or incoming edges, respectively. Given a specifier σ , the corresponding reachability matrix is:

$$R^{(K)}(A; P_w, \sigma) = P_w(0)I_n + \sum_{k=1}^K P_w(k) \prod_{l=1}^k A_{\sigma_l}, \quad \prod_{l=1}^k A_{\sigma_l} = A_{\sigma_k} \cdots A_{\sigma_2} A_{\sigma_1}. \quad (8)$$

Like the reachability definition in Equation 2, the matrix $\mathbf{R}^{(K)}(\mathbf{A}; \mathbf{P}_w, \sigma)$ is a column-stochastic transition matrix, provided that \mathbf{P}_w is properly normalized over K steps. It captures transition probabilities for multi-directional random walks defined by the edge orientation specifier σ . For instance, $\sigma = \mathbf{0}\mathbf{I}\mathbf{0}$ corresponds to walks that follow outgoing, then incoming, then outgoing edges. Our definition also allows walk steps over undirected edges (\mathbf{U}), which will be relevant for defining DirSwitch in Section 4.

Equation 8 defines how to compute multi-directional reachability matrices, but it remains to determine how best to incorporate them into node embeddings. A widely used and effective strategy for this purpose is *embedding concatenation* (Jin et al., 2019; Corso et al., 2020). This approach involves extracting multiple sets of embeddings that capture different structural aspects of the graph, e.g., $\mathbf{Z}^{(1)} \in \mathbb{R}^{n \times p_1}$ and $\mathbf{Z}^{(2)} \in \mathbb{R}^{n \times p_2}$, and combining them into a joint embedding space, $\mathbf{Z} = [\mathbf{Z}^{(1)} \ \mathbf{Z}^{(2)}] \in \mathbb{R}^{n \times (p_1 + p_2)}$. Concatenation is particularly advantageous in unsupervised settings, as it ensures that all constituent embeddings contribute to the final representation.

Concatenation is particularly popular for producing *multi-scale* embeddings (Rozenberczki et al., 2021). The goal of multi-scale embeddings is to represent the inherent hierarchical structure of real-world graphs (Newman, 2018, Ch. 14.7.2) by concatenating embeddings that capture different neighbourhood scales. The ReachNEs framework enables this by incorporating multiple walk length distributions, \mathbf{P}_w . For example, using $\mathbf{P}_w^{(1)} = \text{Geom}(\tau = 1)$ and $\mathbf{P}_w^{(2)} = \text{Pois}(\tau = 3)$ results in two reachability matrices, $\mathbf{R}(\mathbf{P}_w^{(1)})$ and $\mathbf{R}(\mathbf{P}_w^{(2)})$, which focus on immediate and 3-step neighbourhoods, respectively.

The same principle extends to directed neighbourhoods. By concatenating embeddings generated from different edge direction specifiers σ , we obtain *multi-directional* embeddings. This approach can be further combined with multi-scale embeddings to jointly encode both aspects of graph structure.

However, a key limitation of embedding concatenation is that the embedding dimensionality, p , grows with each additional concatenation. This increase can be rapid, given the exponential number of distinct directed neighbourhoods and their possible combinations with multi-scale embeddings. Excessively high dimensionality leads to increased memory and computational costs, as well as potential issues related to the curse of dimensionality. Therefore, it is beneficial to pair concatenation with dimensionality reduction techniques, such as PCA (Murphy, 2012, Ch. 12.2).

4 DirSwitch: Switching reachability filters

The previous section highlights a dual challenge in digraph embedding modelling. Local sinks obstruct information propagation, leading to embeddings biased toward sink nodes. This issue can be easily mitigated by treating the digraph as undirected, effectively ignoring edge directions. However, this approach reduces model expressivity by discarding information about the multiplicity of directed neighbourhoods.

To address these intertwined issues, we propose *switching random walks*. The key idea is to generate multi-directional embeddings by following directed edges for the first r steps of a random walk, thereby capturing local directed neighbourhood structure. After r steps, the graph is treated as undirected, allowing the walk to escape sinks and achieve broader dispersal.

We name this approach DirSwitch. To formally define the DirSwitch model, we introduce a shorthand notation for the edge direction specifier σ . For a K -step random walk, we define σ as a string of length $r + 1$, where $r < K$. The final character, σ_{r+1} , determines the direction for the last $K - r$ steps of the walk.

Rewriting Equation 8 using this notation and taking the limit as $K \rightarrow \infty$, we obtain:

$$\mathbf{R}(\mathbf{A}; \mathbf{P}_w, \sigma) = \mathbf{P}_w(0)\mathbf{I}_n + \sum_{k=1}^r \mathbf{P}_w(k) \prod_{l=1}^{\widehat{k}} \mathbf{A}_{\sigma_l} + \sum_{k=r+1}^{\infty} \mathbf{P}_w(k) \mathbf{A}_{\sigma_{r+1}}^{k-r} \prod_{l=1}^{\widehat{r}} \mathbf{A}_{\sigma_l}. \quad (9)$$

In this formulation, a DirSwitch model is characterized by the first r letters in σ being either $\mathbf{0}$ or \mathbf{I} , representing directed walk steps, while the final character, σ_{r+1} , is set to \mathbf{U} , ensuring undirected propagation for the remaining steps.

The switch to undirected edges addresses the issue of local sinks, but it remains necessary to incorporate representations of different directed neighbourhoods. To achieve this, we use the multi-directional concatenation approach described in Section 3.2. For a given value of r , we define the DirSwitch- r embedding model as the concatenation of all 2^r possible configurations of a σ sequence of length $r + 1$ ending with $\sigma_{r+1} = \mathbf{U}$. For example, DirSwitch-1 concatenates embeddings using $\sigma = \mathbf{0U}$ and $\sigma = \mathbf{IU}$, while DirSwitch-2 uses $\sigma \in \{\mathbf{00U}, \mathbf{0IU}, \mathbf{IUU}, \mathbf{IIU}\}$, and so on for any DirSwitch- r model.

As a baseline in our experiments, we define MultiDir- r using the same notation as above. Specifically, MultiDir-1 uses $\sigma = \mathbf{0}$ and $\sigma = \mathbf{I}$, while MultiDir-2 includes all two-step combinations $\sigma \in \{\mathbf{0}, \mathbf{0I}, \mathbf{I0}, \mathbf{I}\}$, and so on, covering all 2^r possible r -step direction sequences. Unlike DirSwitch, however, MultiDir does not transition to undirected edges after the specified sequence. Instead, the walk continues in the last specified direction. For example, with K total steps and $\sigma = \mathbf{0I}$, the walk takes one step in the $\mathbf{0}$ direction, one in the \mathbf{I} direction, and then continues the remaining $K - 2$ steps in the \mathbf{I} direction.

As noted in Section 3.2, multi-directional concatenation can be combined with multi-scale concatenation to improve embedding quality. However, this combination exhibits diminishing returns: using r steps of directionality with s walk length distributions for a fixed embedding dimension p means that each (σ, P_w) -pair contributes only $\frac{p}{2^r s}$ dimensions. We discuss these diminishing returns further in Appendix A.6, where we also present a practical guide to DirSwitch- r hyperparameter selection based on our experiments and experience. Additionally, increasing r and s raises computational costs, which we analyse in Appendix A.5.

5 Related Work

5.1 Relation to node embedding and graph learning frameworks

The ReachNEs framework is built around the reachability matrix, interpreted as random walk transition probabilities, and two embedding reduction strategies: matrix decomposition and message passing. Similar ideas appear throughout the graph learning literature. Below, we highlight key connections and distinctions.

Sampled random walks. The use of random walks for proximity-based embeddings was popularized by Perozzi et al. (2014) and Grover & Leskovec (2016), who combined sampled walks with negative sampling—originally introduced for word embeddings (Mikolov et al., 2013). As shown by Qiu et al. (2018), these methods asymptotically approximate a matrix factorization of $\log(\mathbf{R}) + b$, where \log is applied elementwise, b is a constant, and \mathbf{R} is the random walk transition probability matrix, as in ReachNEs.

This suggests that sampled walk models and proximity ReachNEs converge to similar embeddings in the limit of infinite samples. However, in practice, sampling is expensive, and only a small fraction of possible walks and non-edges are used. As a result, sampled random walk models and proximity ReachNEs typically yield different embeddings in finite regimes.

Additionally, sampled walk methods usually rely on a single walk length distribution; commonly uniform, though Zhou et al. (2017) use the geometric distribution. In contrast, ReachNEs supports multiple walk length distributions, enabling multi-scale embeddings.

Proximity matrix reduction. Extending the idea of proximity embeddings via matrix decomposition, Zhu et al. (2021a) introduced the PhUSION framework. Like ReachNEs, PhUSION generates embeddings by applying reduction functions to a node proximity matrix, and can encode either structural or automorphic equivalence depending on the chosen reduction. We compare this aspect of the two frameworks in Section 5.2.

In other aspects, PhUSION differs from ReachNEs in two key ways. First, it is restricted to undirected graphs. Second, it interprets the proximity matrix through the lens of matrix functions (Higham, 2008), treating the scalars $P_w(k)$ in Equation 2 as Taylor coefficients rather than probabilities. While this perspective is mathematically more general, it limits analysability. For example, our analysis of reachability entropy in Section 3.1 depends on interpreting \mathbf{R} as a transition probability matrix.

Graph convolutional filters. Like matrix functions, *graph convolutional filters* generalize the reachability matrix (Defferrard et al., 2016; Isufi et al., 2024). They are defined and applied similarly to message-passing ReachNEs, also assuming the presence of initial node features \mathbf{X} . However, the series coefficients, $P_w(k)$ in Equation 2, can take arbitrary real values, making random walk interpretations generally inapplicable.

Instead, convolutional filters are viewed as extensions of classical convolutions to graphs and are typically analysed using spectral methods based on eigenfunctions (Bronstein et al., 2017). Since asymmetric matrices lack a standard eigendecomposition, spectral analysis becomes difficult on digraphs, and most research in this area focuses on undirected graphs (Zhang et al., 2021b). In contrast, ReachNEs defines \mathbf{R} purely through random walk probabilities, making the framework naturally applicable to both directed and undirected graphs.

Message-passing graph neural networks. The dominant paradigm in node embedding learning is message-passing graph neural networks (GNNs) (Gilmer et al., 2017). GNNs are typically nonlinear and parameterized, requiring training via non-convex optimization. While this complicates their analysis compared to ReachNEs, their overall algorithmic structure is similar. In fact, only minor modifications to Equation 6 are needed to recover the GraphSAGE encoder (Hamilton et al., 2017). We exploit this connection in Section 6.5 to show how our results extend to the self-supervised GNN setting.

Interestingly, the binomial walk length distribution in ReachNEs is related to the self-loops added in convolutional GNNs (Kipf & Welling, 2017; Wu et al., 2019). We explore this connection further in Appendix C.5.

Graph kernels. While many embedding methods implicitly rely on walk length distributions, few works explicitly analyse or compare their effects as we do. In contrast, the graph kernel literature, particularly *random walk kernels*, has explored related formulations, including both geometric and Poisson walk length distributions (Vishwanathan et al., 2010). However, since graph kernels are designed for graph-level comparison rather than node-level tasks, it remains an open question which insights transfer to ReachNEs.

5.2 Structural and automorphic equivalence

Recent work has highlighted the distinction between embedding models that capture structural versus automorphic node equivalence (Jin et al., 2021; Zhu et al., 2021a; Yan et al., 2024). Models targeting structural equivalence are often called *positional embeddings*, while those targeting automorphic equivalence are referred to as *structural embeddings*. The terminology is admittedly confusing, as the term *structural* is used in two different senses.

ReachNEs draws inspiration from the PhUSION framework (Zhu et al., 2021a), which also produces both positional and structural embeddings by applying different reduction methods to a node proximity matrix (analogous to our reachability matrix). For positional embeddings, both ReachNEs and PhUSION use SVD reduction. However, PhUSION assumes undirected graphs and retains only the left singular vectors \mathbf{U} , while ReachNEs concatenates both \mathbf{U} and \mathbf{V} to handle directed graphs.

For structural embeddings, PhUSION applies the permutation-invariant empirical characteristic function to each column of \mathbf{R} , following the approach introduced by Donnat et al. (2018). In contrast, ReachNEs uses message-passing reduction, aligning more directly with the message-passing GNN literature.

Finally, while some argue that positional and structural embeddings represent fundamentally different objectives (Rossi et al., 2020), Srinivasan & Ribeiro (2020) established a statistical equivalence between them. However, the practical consequences of this result remain an open question.

5.3 Digraph models

Research on message-passing embedding models, particularly GNNs, has primarily focused on undirected graphs (Rossi et al., 2023). This is especially true in unsupervised and self-supervised learning, where recent state-of-the-art models remain constrained to undirected graphs (Hassani & Khasahmadi, 2020; Thakoor et al., 2022; Zhang et al., 2021a; Hou et al., 2023).

For (semi-)supervised learning, Rossi et al. (2023) proposed digraph extensions for common GNN embedding encoders. Unlike DirSwitch, which uses concatenation to preserve the contributions of different directed neighbourhoods in distinct embedding components, Rossi et al. (2023) employs parameterized mixing of these components. This approach requires the model to learn how to extract and preserve relevant information from multi-directional neighbourhoods. However, this filtering process is more challenging in unsupervised settings due to the weaker learning signal.

In contrast, research on unsupervised proximity embeddings for digraphs has progressed further than its message-passing counterpart. The models HOPE (Ou et al., 2016), APP (Zhou et al., 2017), and NERD (Khosla et al., 2020) follow a similar approach to ReachNEs and DirSwitch, leveraging random walks and matrix factorization of proximity matrices. However, none of these models explicitly address local sinks and neighbourhood multiplicity, which are the primary focus of our work.

Additionally, there are several notable algorithmic differences. Both HOPE and APP use the geometric walk length distribution, while NERD employs the uniform distribution. In contrast, ReachNEs treats the choice of distribution as a parameter. NERD also introduces *alternating walks*, where edge direction specifiers switch between outgoing and incoming edges (e.g., $\sigma = 0101\dots$ or $\sigma = 1010\dots$). While this pattern may help random walks escape some sinks, restricting the model to only these sequences introduces bias by excluding other directed neighbourhood definitions.

Furthermore, HOPE utilizes partial and generalized SVD (Hochstenbach, 2009) for matrix factorization, whereas ReachNEs employs the more scalable single-pass SVD (Yu et al., 2017). Both APP and NERD rely on Monte Carlo sampling of random walks and implicit matrix factorization via gradient descent and negative sampling, while ReachNEs directly computes the reachability matrix using matrix multiplication.

Finally, DGGAN (Zhu et al., 2021b) and BLADE (Virinchi & Saladi, 2023) are unsupervised digraph embedding methods that incorporate neural networks. While neither explicitly addresses the issue of local sinks, BLADE follows an embedding mixing strategy similar to Rossi et al. (2023), where a GNN is responsible for filtering information from distinct directed neighbourhoods.

5.4 Expressivity of message-passing models

The expressivity of message-passing models is a central topic in graph representation learning due to its connection with the graph isomorphism problem (Grohe & Schweitzer, 2020). Li & Leskovec (2022, Def. 5.5) offers a formal definition, characterizing model expressivity as the ability to produce distinct representations for non-isomorphic graphs. While their formulation is mathematically rigorous, our definition in Section 3.2 captures the same core idea in a more intuitive manner, with a focus on node equivalence rather than full graph isomorphism.

The foundational works of Xu et al. (2019) and Morris et al. (2019) relate the expressivity of message-passing graph neural networks (MP-GNNs) to the Weisfeiler-Leman (WL) colour refinement algorithm (Grohe et al., 2021), an iterative method used to test graph isomorphism. Rossi et al. (2023) extended these results to digraphs by incorporating both in- and out-neighbourhoods in the message-passing process, in a manner analogous to our MultiDir and DirSwitch approaches.

These studies demonstrate that the expressive power of MP-GNNs is upper-bounded by the WL algorithm. This also applies to message-passing ReachNEs, which can be viewed as a linear MP-GNN. Moreover, they show that MP-GNNs can match the expressivity of the WL test if their message aggregation functions are sufficiently flexible—specifically, if they can form injective mappings from multisets to vectors, mimicking the WL hashing mechanism. These assumptions, however, do not hold for ReachNEs, which is linear and parameterized solely through scalar-valued walk length distributions.

The use of flexible, non-linear aggregation functions is well-suited to the supervised, end-to-end learning context studied by Xu et al. (2019), Morris et al. (2019), and Rossi et al. (2023), where parameters can be optimized to task-specific objectives. In contrast, the unsupervised setting considered in this work demands greater interpretability and analytical tractability, especially for downstream tasks such as anomaly detection.

Table 2: Graph statistics for the datasets. The columns report the number of nodes (n) and edges (m), the median out- and in-degrees ($|\text{deg}_0|$, $|\text{deg}_I|$), the number of weakly and strongly connected components ($\#CC$ and $\#SCC$), the global clustering coefficient (C_G), and the average path length ($\langle l_{\text{path}} \rangle$) computed on the undirected version of the graph. Additionally, when available, we list the number of node attributes (d), the number of node label classes ($|\mathbb{Y}|$), and the homophily coefficient of the undirected graph (h_U). The final column specifies the graph type. We use the shorthand notation with $K=10^3$ and $M=10^6$ for large values.

DATASET	n	m	$ \text{deg}_0 $	$ \text{deg}_I $	$\#CC$	$\#SCC$	C_G	$\langle l_{\text{path}} \rangle$	d	$ \mathbb{Y} $	h_U	TYPE
ARXIV ¹	170K	1.2M	4	1	1	141K	0.017	5.7	128	40	0.64	CITATION
ARXIV-YEAR ²	170K	1.2M	4	1	1	141K	0.017	5.7	128	5	0.29	CITATION
CITeseer ³	4.2K	5.4K	1	0	515	4209	0.084	7.4	602	6	0.96	CITATION
CoCite ⁴	44K	20K	2	2	652	44K	0.081	5.5	–	15	0.42	CITATION
CORA ³	20K	65K	2	1	364	16K	0.14	6.2	8710	70	0.59	CITATION
CORA-ML ³	3K	8.4K	2	1	61	2603	0.12	5.3	2879	7	0.82	CITATION
CORA (SUBELJ) ⁴	23K	92K	3	1	1	18K	0.12	5.8	–	70	0.56	CITATION
ENRON ⁵	7.9K	142K	3	7	58	861	0.16	3.1	–	–	–	EMAIL
EU-EMAIL ⁶	1K	25K	14	17	20	203	0.29	2.6	–	42	0.47	EMAIL
FLY LARVA ⁷	3K	116K	33	33	5	136	0.30	2.7	–	93	0.14	CONNECTOME
POLBLOGS ⁸	1.5K	19K	4	2	268	688	0.25	2.7	–	2	0.91	HYPERLINK
POKEC ²	1.6M	31M	8	8	1	326K	0.058	4.7	65	3	0.43	DATING
PUBMED ⁴	20K	44K	0	1	1	20K	0.054	6.4	–	3	0.79	CITATION
ROMAN EMPIRE ²	23K	33K	1	1	1	23K	0.28	2.4K	300	18	0.05	TEXT
SNAP PATENTS ²	2.9M	14M	3	3	181K	2.9M	0.066	6.8	269	5	0.22	CITATION
WIKIVOTE ⁹	7K	103K	2	0	24	5816	0.14	3.2	–	–	–	VOTING

¹ HU ET AL. (2020A)² LIM ET AL. (2021)³ BOJCHEVSKI & GÜNNEMANN (2018)⁴ KHOSLA ET AL. (2020)⁵ KLIMT & YANG (2004)⁶ YIN ET AL. (2017)⁷ WINDING ET AL. (2023)⁸ ADAMIC & GLANCE (2005)⁹ LESKOVEC ET AL. (2010)

In this regard, our work aligns more closely with Jin et al. (2019) and Corso et al. (2020), who emphasize embedding concatenation as a strategy for increasing expressivity. While their approaches combine embeddings from different aggregation operators (e.g., mean, max, sum), we instead concatenate embeddings derived from multiple edge directions and walk length distributions. This preserves both interpretability and theoretical clarity, as each component of the embedding retains a well-defined meaning grounded in random-walk semantics.

6 Experiments

Our experiments focus on evaluating DirSwitch and are divided into three parts. First, we validate that DirSwitch mitigates low dispersal caused by local sinks in digraphs (Section 6.1) while preserving the ability to represent directed neighbourhoods (Section 6.2). Second, in Section 6.3, we demonstrate that these improvements lead to higher-quality embeddings by evaluating DirSwitch in combination with ReachNEs on 14 node classification benchmark datasets.

Third, we evaluate DirSwitch’s practical effectiveness by comparing it to recent digraph embedding models on node classification benchmarks. Specifically, we first assess DirSwitch with ReachNEs proximity embeddings against state-of-the-art unsupervised digraph proximity embedding approaches (Section 6.4). Then, in Section 6.5, we demonstrate DirSwitch’s flexibility by applying it beyond the ReachNEs framework. We extend self-supervised GNNs to digraphs using ReachNEs and compare this approach to the method proposed by Rossi et al. (2023), which generalizes GNNs to digraphs in the semi-supervised setting.

Our experiments use the graph learning datasets summarized in Table 2, which span a diverse range of graph types and properties, including variations in density, connectivity, and node attributes. Unless otherwise specified, we compute truncated reachability using $K = 12$ steps. All experiments were conducted in a Google Cloud environment with an Nvidia L4 24GB GPU, 32 vCPUs @ 2.20GHz, and 128GB of memory.

6.1 Improving dispersal

To verify that DirSwitch mitigates local sinks and improves dispersal, we measure the reachability entropy $H(j; \mathbf{R}) = -\sum_{i=1}^n R_{i,j} \log_2 R_{i,j}$ (see Section 3.1) for various reachability matrices $\mathbf{R}(\mathbf{P}_w, \sigma)$. Specifically, we

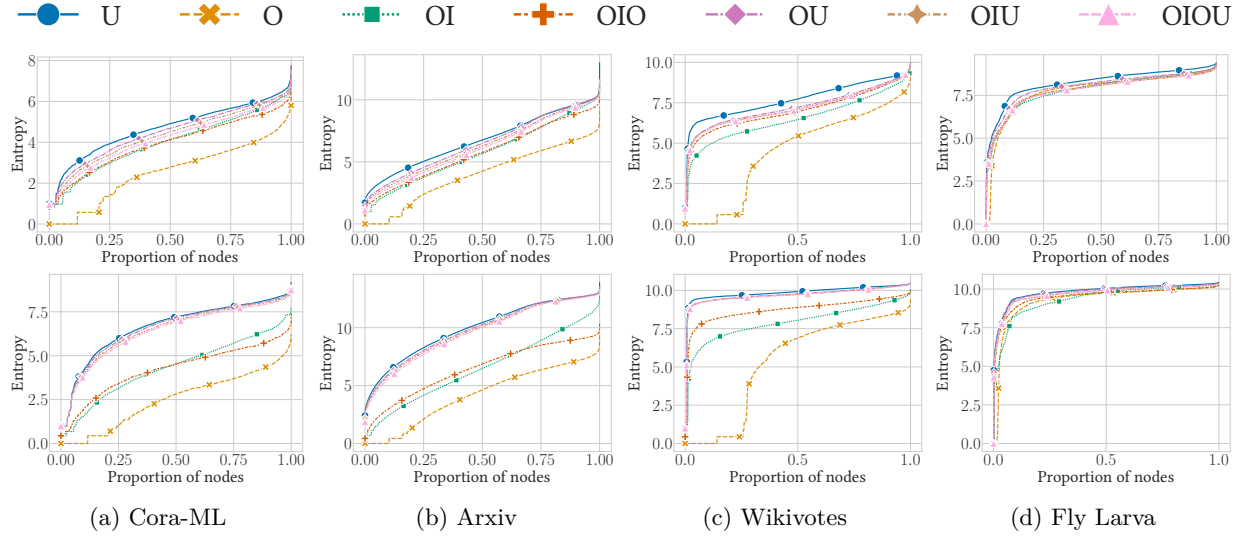


Figure 6: Neighbourhood dispersal evaluation for four graphs, measured via reachability entropy, $-\sum_{i=1}^n R_{i,j} \log_2 R_{i,j}$, computed for each node and sorted. Each curve corresponds to a different edge direction specifier σ , with the top row showing results for $P_w(k; \tau) = \text{Pois}(\tau = 2)$ (local dispersal) and the bottom row for $\mathcal{U}(\tau = 5)$ (long-range dispersal). DirSwitch variants (e.g., OU, OIU, OIOU) demonstrate high dispersal, comparable to U, while purely directed specifiers (O, OI, OIO) exhibit lower entropy due to sink effects.

compare DirSwitch edge direction specifiers $\sigma \in \{\text{OU}, \text{OIU}, \text{OIOU}\}$ against MultiDir specifiers $\sigma \in \{\text{O}, \text{OI}, \text{OIO}\}$ and purely undirected edges, $\sigma = \text{U}$. We use two walk length distributions: $\text{Pois}(\tau = 2)$, which highlights short-range differences, and $\mathcal{U}(\tau = 5)$, which captures long-range behaviour.

Figure 6 presents results for four representative graphs, with additional datasets provided in Appendix E. The top row shows results for $\text{Pois}(\tau = 2)$, while the bottom row corresponds to $\mathcal{U}(\tau = 5)$.

The entropy associated with DirSwitch specifiers closely approaches that of undirected edges, which consistently achieve the highest entropy due to the absence of sinks. This outcome is both expected and desirable, highlighting how the switch to undirected edges mitigates the dispersal-limiting effects of local sinks.

These effects are most pronounced for long-range walks using $\mathcal{U}(\tau = 5)$, shown in the bottom row. In Figures 6a, 6b, and 6c, there is a significant entropy gap between the DirSwitch curves and the MultiDir curves (O, OI, OIO). The Cora-ML, Arxiv, and Wikivote graphs are sparsely connected (see Table 2), leading to increased reachability concentration in local sinks and lower entropy. In contrast, the densely connected Fly Larva graph in Figure 6d exhibits a less pronounced but still noticeable entropy difference.

For local smoothing with $\text{Pois}(\tau = 2)$ (top row), the entropy increase for DirSwitch over OI and OIO is visible but small. This is because alternating O- and I-steps allows short-range walkers to partially escape sinks. Moreover, we observe a discernible gap between U and the DirSwitch directions. This gap reflects the reduced dispersal caused by using edge directions in the initial steps of DirSwitch random walks. As walk length increases, both effects vanish as the undirected regime of DirSwitch dominates, as seen in the bottom row.

6.2 Directed neighbourhood expressivity

We use node embedding graph alignment (Heimann et al., 2018) to verify that DirSwitch can represent the diversity of directed neighborhoods. Graph alignment is a generalized version of the graph isomorphism problem, where nodes in two graphs are matched based on structural similarity (Skitsas et al., 2023).

Following the embedding benchmark protocols of Heimann et al. (2018) and Jin et al. (2021), we construct a second graph $\mathcal{G}_2 = (\mathbb{N}, \mathbb{M}_2)$ from a given graph $\mathcal{G}_1 = (\mathbb{N}, \mathbb{M}_1)$ by removing 15% of the edges from \mathbb{M}_1 and randomly permuting the node indices. Node embeddings are then computed for both graphs, and each node

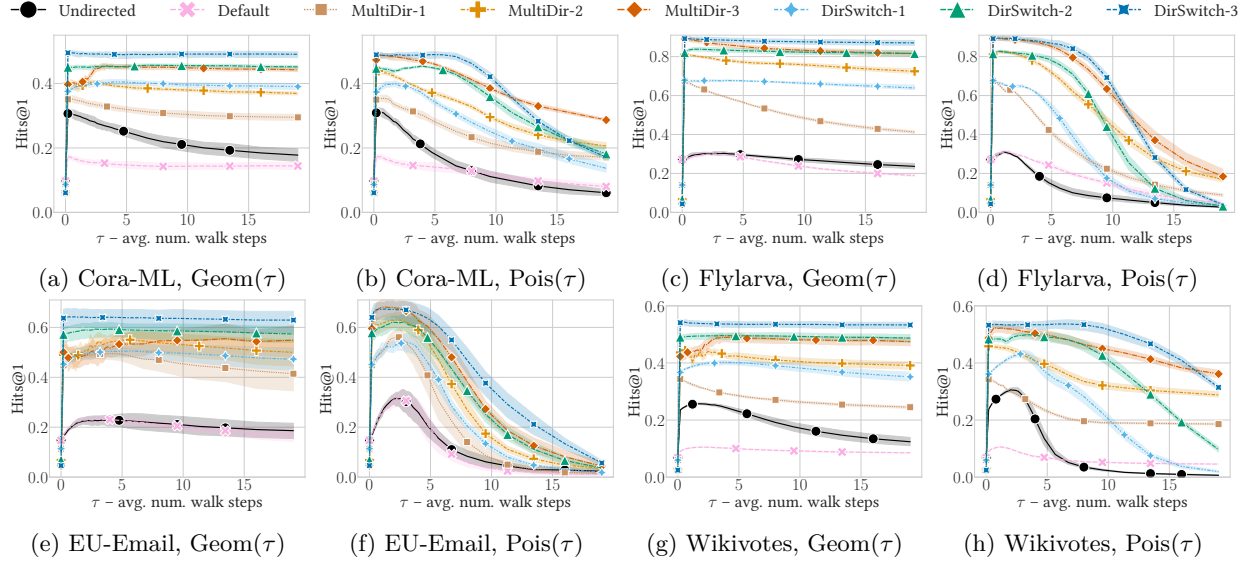


Figure 7: Evaluation of edge direction expressivity for four graphs using the geometric, $\text{Geom}(\tau)$, and Poisson, $\text{Pois}(\tau)$, walk length distributions. The y-axes represent graph alignment accuracy under 15% edge removal, while the x-axes correspond to τ , the average walk length. The curve colours and styles denote different sets of edge direction specifiers, σ .

in \mathcal{G}_2 is matched to a node in \mathcal{G}_1 based on the shortest Euclidean distance between embeddings. This process is repeated five times with different random seeds.

The graph alignment task benefits significantly from embedding expressivity, particularly the ability to distinguish distinct directed neighbourhoods, as this reduces the risk of erroneous matches. However, for matching to be feasible, embeddings must also capture local graph structure that generalizes across separate graph components. To achieve this, we use ReachNEs message-passing embeddings combined with node attributes derived from graph structure, including out- and in-degrees and the four digraph local clustering coefficients—*out*, *in*, *cycle*, and *middleman*—as defined by Fagiolo (2007).

We compare the multi-directional DirSwitch- r against MultiDir- r for $r \leq 3$, as well as undirected ($\sigma = \mathbf{U}$) and default directed ($\sigma = \mathbf{0}$) edges. To account for the influence of walk length distribution on embedding expressivity, we evaluate both $P_w = \text{Geom}(\tau)$ and $P_w = \text{Pois}(\tau)$ for $\tau \in (0, 20]$, using a reachability truncation of $K = 30$ steps.

Figure 7 presents results for four datasets, with τ values along the x-axes and graph alignment accuracy on the y-axes. Additional results for other datasets and walk length distributions are provided in Appendix F. DirSwitch- r and MultiDir- r achieve comparable accuracies for each r , with accuracy increasing as r grows. This similarity is expected, as both methods use concatenation to represent 2^r directed neighbourhoods.

A closer comparison reveals that DirSwitch- r consistently outperforms MultiDir- r , provided τ is not too large. This improvement stems from DirSwitch mitigating sink bias. When reachability is concentrated in small node sets, embeddings become overly sensitive to edge removal, as this can then radically change the reachability distribution. This sensitivity reduces the alignment accuracy. By alleviating this bias, DirSwitch provides more robust representations.

For the same reason, the undirected edges perform better than the default directed edges. However, overall, both these single-neighbourhood baselines achieve the lowest overall accuracies due to their inability to capture directed neighbourhood multiplicity.

Another notable observation is the impact of the walk length distribution $P_w(\tau)$. In all cases, accuracy collapses as $\tau \rightarrow 0$, corresponding to the use of node attributes without smoothing. Additionally, for the Poisson distribution, accuracy decreases as τ increases, whereas for the geometric distribution, it plateaus. This behaviour aligns with the expected asymptotic properties of both distributions.

Table 3: Abbreviations used when reporting results for multi-scale reachability embeddings. Visualizations are available in Figure 11 in the Appendix.

ABBREVIATION	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DISTRIBUTIONS, P_w	$\text{Geom}(k; \tau = 1)$	$\text{Pois}(k; \tau = 2)$	$\text{Geom}(k; \tau = 1),$ $\mathcal{U}(k - 1; \tau = 2)$	$\text{Binom}(k; \tau = 1),$ $\text{Binom}(k - 2; \tau = 2),$ $\text{Binom}(k - 5; \tau = 3)$	$\text{Geom}(k; \tau = 1),$ $\text{Geom}(k - 1; \tau = 2),$ $\text{Geom}(k - 2; \tau = 3),$ $\text{Geom}(k - 3; \tau = 4)$

In short, as τ increases, the mode of $\text{Pois}(\tau)$ shifts continuously, emphasizing long-range walks over local smoothing. This results in less distinguishable embeddings and lower alignment accuracy. In contrast, the geometric distribution maintains local neighbourhood smoothing, preventing significant accuracy drops at higher τ values. See Appendix C for further details.

The difference between Pois and Geom is also evident in the comparison between $\text{DirSwitch-}r$ and $\text{MultiDir-}r$. With the geometric distribution, DirSwitch consistently achieves higher alignment accuracy across varying τ values, as discussed above. However, for the Poisson distribution, DirSwitch ’s accuracy declines more rapidly than that of MultiDir as τ increases in some datasets. For instance, this effect is observed for DirSwitch-3 and MultiDir-3 in Figures 7b and 7d.

This behaviour stems from the shifting mode of Pois . In the limit $\tau \rightarrow \infty$, DirSwitch embeddings converge to identical representations within a weakly connected component. In contrast, MultiDir embeddings retain finer granularity, as equivalence groups in the directed case remain more distinct, see our discussion at the end of Appendix C.1. This explains why MultiDir can surpass DirSwitch in alignment accuracy at large τ values.

6.3 DirSwitch embedding quality evaluation

In this section, we demonstrate that DirSwitch ’s increased dispersal and multi-directionality lead to higher-quality embeddings. To assess this, we use node classification as a representative downstream task, comparing ReachNEs embeddings generated with $\text{DirSwitch-}r$, $\text{MultiDir-}r$, undirected edges (\mathcal{U}), and default directed edges (\mathcal{O}).

Setup Table 2 lists the node classification benchmark datasets. For datasets with node attributes, we evaluate message-passing ReachNEs , while for those without attributes, we use proximity ReachNEs .

Node embeddings are computed in a fully unsupervised manner, after which a logistic regression classifier is trained on frozen embeddings. We employ a 3x repeated 5-fold cross-validation strategy to obtain mean performance and standard deviations.

To assess the impact of both multi-directional and multi-scale concatenation, we include both single- and multi-scale embeddings in our analysis. Table 3 details the walk length distributions used for single- and multi-scale embeddings. We consider two single-scale distributions and three multi-scale distributions, incorporating two, three, and four scales, respectively. The embedding dimensions are set to $p = 1024$ and $p = 512$.

Results Overview Table 4 reports classification accuracies for message-passing ReachNEs using datasets with node attributes, while Table 6 presents results for proximity ReachNEs for datasets without attributes, both using $p = 1024$. Across all datasets and embedding types, $\text{DirSwitch-}r$ consistently achieves the highest accuracies or performs within one to two standard deviations of the best results. In the tables, we highlight the best results for each scale in bold blue and results within one standard deviation in light blue. The consistent top performance of DirSwitch across datasets and multi-scale choices underscores its superior embedding quality compared to non-switching approaches. Similar trends are observed for $p = 512$, with further analysis provided in Appendix G.

A notable trend is that $\sigma = 0$ (default) yields some of the lowest accuracies across all datasets and scales, as highlighted in orange. This suggests that embedding quality improves both by using undirected edges (\mathcal{U}) and by concatenating $\sigma = \mathcal{I}$ embeddings, as done in MultiDir-1 . These results further support the claim that mitigating local sinks and capturing different directed neighbourhoods enhances embedding quality.

Next, we conduct a deeper analysis of these results. To facilitate this, we use the relative accuracy improvements in Tables 5 and 7. These values represent the highest accuracy in each row of Tables 4 and 6, relative to the accuracy of the default. For default edge directions themselves, absolute accuracy values are reported.

Table 4: Node classification accuracy for message-passing ReachNEs with $p = 1024$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. Each entry reports the mean accuracy and standard deviation. Bold blue highlights the highest accuracy in each column, with light blue indicating results within one standard deviation of the best. Similarly, bold orange denotes the lowest accuracy, and light orange represents values within one standard deviation of the worst.

EDGE DIRECTIONS	ARXIV					ARXIV YEAR				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	55.0±0.2	42.4±0.3	55.5±0.3	54.2±0.2	55.4±0.3	36.5±0.3	35.4±0.5	36.7±0.2	37.8±0.2	38.2±0.2
UNDIRECTED	61.6±0.2	64.9±0.3	67.0±0.2	69.8±0.2	69.1±0.2	36.4±0.2	36.5±0.2	37.1±0.2	37.6±0.3	38.3±0.3
MULTIDIR-1	59.8±0.3	59.4±0.3	60.2±0.3	60.5±0.2	60.3±0.2	37.0±0.3	35.3±0.8	38.1±0.2	39.3±0.2	39.7±0.3
MULTIDIR-2	61.2±0.3	65.3±0.3	65.2±0.3	64.9±0.3	65.3±0.3	38.5±0.3	37.9±0.4	40.9±0.2	42.0±0.2	41.7±0.2
MULTIDIR-3	61.5±0.3	65.5±0.3	65.9±0.3	65.9±0.3	65.7±0.3	40.0±0.3	40.6±0.3	42.3±0.2	43.2±0.3	42.4±0.3
DIRSWITCH-1	61.5±0.3	65.9±0.2	67.3±0.3	69.7±0.2	68.4±0.3	38.2±0.3	37.8±0.2	39.9±0.3	40.2±0.3	41.5±0.2
DIRSWITCH-2	61.5±0.3	65.8±0.3	66.9±0.3	69.6±0.3	68.4±0.2	39.6±0.2	39.5±0.3	41.4±0.3	41.9±0.3	42.4±0.2
DIRSWITCH-3	61.5±0.3	65.8±0.3	66.7±0.3	69.0±0.3	67.8±0.3	40.0±0.3	40.9±0.4	42.4±0.3	43.3±0.3	43.3±0.3

EDGE DIRECTIONS	CORA-ML					CITeseer					CORA				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	31.5±1.6	35.8±2.4	46.1±2.1	58.4±2.2	63.2±2.3	67.5±2.0	56.5±1.9	60.8±2.7	70.0±2.0	72.9±2.1	52.4±0.6	51.2±0.7	57.6±0.7	59.3±0.7	60.2±0.8
UNDIRECTED	31.2±1.9	41.0±2.2	54.5±1.7	72.1±2.2	76.1±1.8	76.2±1.8	83.7±1.4	72.6±1.8	88.0±1.7	89.0±1.2	59.4±0.6	64.9±0.8	67.5±0.7	69.3±0.7	69.6±0.7
MULTIDIR-1	55.1±2.1	53.2±2.0	66.9±2.1	72.2±2.1	75.4±2.3	66.1±2.0	60.8±1.8	73.8±1.3	78.6±1.7	81.0±1.6	64.3±0.6	63.3±0.6	65.8±0.8	65.3±1.0	64.8±1.0
MULTIDIR-2	71.3±1.6	71.1±2.3	78.4±1.8	79.8±2.1	82.3±1.4	79.4±1.2	78.4±1.2	84.3±1.4	85.0±1.2	86.0±1.4	67.1±0.7	68.4±0.7	67.8±0.9	66.5±0.9	66.0±0.9
MULTIDIR-3	78.9±1.9	80.0±1.0	84.0±1.4	83.9±1.4	84.5±1.5	84.7±1.7	87.2±1.2	87.0±1.5	87.3±1.3	87.5±1.2	66.7±0.9	68.5±0.9	66.8±0.9	65.3±0.8	64.4±0.9
DIRSWITCH-1	56.9±2.9	57.5±2.4	72.0±2.0	79.4±1.6	82.3±2.0	73.9±1.2	76.5±1.9	85.6±1.5	91.4±0.8	91.3±1.5	65.6±0.7	67.5±0.8	69.5±0.8	70.1±0.9	69.7±0.8
DIRSWITCH-2	71.6±1.8	72.4±2.4	80.7±1.6	84.3±1.6	84.4±1.9	80.4±1.3	85.1±1.4	88.6±1.0	90.9±1.0	90.1±1.0	67.5±0.6	68.9±0.7	68.9±0.9	69.0±0.7	68.0±0.7
DIRSWITCH-3	79.1±2.0	80.6±1.6	84.0±1.5	85.6±1.6	85.7±1.4	84.7±1.3	88.0±0.9	87.7±1.3	89.6±0.8	88.5±1.2	66.4±0.8	68.7±0.8	67.0±0.8	67.1±0.6	65.6±0.9

EDGE DIRECTIONS	ROMAN EMPIRE					POKEC					SNAP PATENTS				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	69.3±0.9	32.5±0.4	68.5±0.6	51.2±0.6	67.2±0.8	60.1±0.2	59.0±0.8	61.7±0.1	60.7±0.1	65.4±0.1	35.7±0.1	41.7±0.1	39.2±0.1	42.2±0.2	43.1±0.3
UNDIRECTED	67.2±0.8	46.9±0.7	70.9±0.7	65.9±0.8	70.0±0.4	60.3±0.1	59.6±0.5	62.0±0.1	62.3±0.1	71.7±0.1	30.8±0.1	29.9±0.1	31.8±0.1	33.2±0.1	34.7±0.4
MULTIDIR-1	71.8±0.8	42.0±0.8	66.7±0.8	56.6±0.7	69.9±0.9	61.3±0.1	60.1±0.1	62.2±0.1	61.7±0.1	67.3±0.1	41.0±0.1	43.9±0.4	45.7±0.1	47.8±0.2	46.9±0.2
MULTIDIR-2	71.9±0.5	60.4±0.7	71.1±0.5	65.2±0.6	75.5±0.6	61.5±0.1	61.1±0.1	63.2±0.1	63.0±0.1	69.1±0.1	45.2±0.1	46.5±0.2	46.8±0.2	49.8±0.2	48.1±0.2
MULTIDIR-3	73.1±0.6	68.0±0.6	73.7±0.7	69.2±0.7	76.0±0.5	61.8±0.1	63.5±0.0	64.8±0.1	64.8±0.1	70.7±0.1	45.9±0.1	49.1±0.1	47.4±0.1	49.8±0.7	46.5±0.1
DIRSWITCH-1	73.0±0.6	58.4±0.8	72.4±0.9	67.5±0.5	74.6±0.6	61.2±0.1	60.3±0.3	62.8±0.1	62.8±0.1	70.0±0.1	40.8±0.1	40.8±0.1	42.7±0.1	43.8±0.1	44.8±0.3
DIRSWITCH-2	72.0±0.5	62.8±0.7	74.2±0.5	70.8±0.6	76.4±0.7	61.4±0.1	60.9±0.1	63.5±0.1	64.1±0.1	71.5±0.1	46.6±0.1	48.1±0.2	46.8±0.1	47.2±0.2	47.6±0.4
DIRSWITCH-3	72.8±0.6	67.9±0.6	75.1±0.6	72.5±0.6	76.4±0.6	61.6±0.1	62.7±0.1	64.7±0.1	64.4±0.1	70.3±0.1	46.3±0.1	50.8±0.3	49.2±0.1	49.4±0.3	48.8±0.3

Table 5: Relative improvements in node classification accuracy for message-passing ReachNEs with $p = 1024$. The values reflect the maximum accuracy for per row and dataset in Table 4. The top row displays absolute accuracies for the default edge directions, with standard deviations expressed as percentages. The subsequent rows present the relative improvements compared to the top row. The table is structured with the four homophilic datasets on the left and the four heterophilic datasets on the right.

EDGE DIRECTIONS	CORA-ML	CITeseer	CORA	ARXIV	ROMAN EMPIRE	ARXIV YEAR	POKEC	SNAP PATENTS
DEFAULT	63.2±3.4%	72.9±3.0%	60.2±1.1%	55.5±0.5%	69.3±0.9%	38.2±0.8%	65.4±0.4%	43.1±0.4%
UNDIRECTED	+20.5%	+22.0%	+15.5%	+25.8%	+2.2%	+0.2%	+9.5%	-19.3%
MULTIDIR-1	+19.4%	+11.1%	+9.2%	+9.0%	+3.5%	+3.9%	+2.8%	+11.0%
MULTIDIR-2	+30.3%	+18.0%	+13.5%	+17.7%	+8.9%	+9.9%	+5.6%	+15.6%
MULTIDIR-3	+33.7%	+20.0%	+13.6%	+18.9%	+9.6%	+12.9%	+8.0%	+15.6%
DIRSWITCH-1	+30.3%	+25.4%	+16.3%	+25.7%	+7.5%	+8.5%	+6.9%	+4.1%
DIRSWITCH-2	+33.6%	+24.7%	+14.5%	+25.4%	+10.2%	+10.9%	+9.3%	+11.7%
DIRSWITCH-3	+35.7%	+22.9%	+14.0%	+24.5%	+10.2%	+13.4%	+7.5%	+17.9%

Homophilic vs heterophilic datasets for message-passing embeddings The node classification datasets in Table 2 exhibit varying levels of *homophily* (Zhu et al., 2020), i.e., the tendency of nodes with the same class label to connect. The homophily level is quantified in the h_v column, representing the average fraction of a node’s neighbours that share its label (Rossi et al., 2023, Eq. 1).

Table 6: Node classification accuracy for proximity ReachNEs with $p = 1024$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. The values denote average accuracies with standard deviations. Bold blue highlights the best results in each column, with light blue indicating results within one standard deviation. Similarly, bold orange marks the worst results, with light orange showing values within one standard deviation of the lowest performance.

EDGE DIRECTIONS	FLY LARVA					EU-EMAIL					POLBLOGS				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	42.2±2.0	45.2±2.0	49.0±1.9	49.5±1.7	50.1±2.0	30.3±3.3	36.6±2.9	54.0±3.0	65.0±2.8	66.0±3.2	61.8±2.7	70.2±2.4	73.0±2.7	79.1±1.9	82.0±2.1
UNDIRECTED	45.0±1.7	50.9±2.0	53.7±2.1	54.6±1.7	55.4±1.7	29.8±3.0	35.4±3.7	57.2±3.5	70.6±2.9	72.5±3.1	64.4±3.3	76.8±2.0	79.1±2.4	84.3±2.2	85.7±1.9
MULTIDIR-1	47.3±1.6	49.1±1.8	50.8±1.8	51.9±1.9	52.0±1.9	50.2±3.1	59.7±3.7	68.1±2.9	72.2±2.8	73.5±3.0	73.7±2.6	74.7±2.6	83.6±1.9	85.2±1.7	86.7±2.1
MULTIDIR-2	50.3±2.3	50.7±2.1	52.4±2.2	54.2±2.0	53.6±1.7	66.8±3.2	69.5±3.3	74.9±3.0	75.2±2.8	74.7±2.6	83.4±2.3	84.4±2.1	88.2±2.0	88.7±1.6	89.2±1.5
MULTIDIR-3	51.8±2.6	52.6±2.2	53.8±2.1	53.6±1.8	53.0±1.5	73.5±3.4	74.2±2.4	74.9±2.3	73.1±3.3	70.4±2.9	88.0±1.7	88.0±2.0	89.2±1.4	89.7±1.5	89.8±1.6
DIRSWITCH-1	48.8±1.7	51.6±1.6	53.8±2.4	55.8±1.7	55.8±2.4	50.9±3.2	59.8±3.7	70.5±3.1	74.8±2.8	75.6±2.4	75.6±2.6	75.3±2.2	85.5±2.4	86.6±2.1	86.2±2.2
DIRSWITCH-2	50.3±1.7	51.9±2.1	54.1±1.5	56.7±2.0	56.2±1.9	66.8±2.7	69.8±3.5	74.6±2.9	75.7±2.5	74.8±2.5	83.8±2.2	85.2±1.7	88.4±2.1	87.1±2.1	89.7±1.6
DIRSWITCH-3	52.4±1.7	52.8±1.6	55.6±1.8	55.6±1.8	55.6±1.8	73.3±2.6	73.8±2.7	75.3±2.6	72.2±2.9	70.6±2.8	87.9±1.8	88.4±1.6	89.4±1.6	87.1±2.0	89.8±1.5

EDGE DIRECTIONS	CoCITE					PUBMED					CORa (SUBELJ)				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	38.6±0.4	39.0±0.4	39.9±0.5	39.8±0.4	40.2±0.6	72.4±0.7	72.6±0.7	73.3±0.8	73.1±0.7	72.6±0.6	58.2±0.7	58.2±1.0	59.0±0.8	59.2±0.9	59.3±0.8
UNDIRECTED	45.3±0.5	45.8±0.6	46.6±0.4	47.9±0.5	47.6±0.5	81.9±0.5	82.2±0.6	82.2±0.5	82.4±0.6	82.2±0.6	65.7±0.7	66.0±0.8	66.7±0.7	67.1±0.8	67.0±0.8
MULTIDIR-1	39.8±0.5	40.0±0.5	40.8±0.5	41.3±0.6	41.3±0.5	73.9±0.7	73.5±0.8	72.9±0.6	72.7±0.8	71.5±0.7	58.7±0.8	58.5±0.8	58.7±0.7	58.3±0.8	58.3±0.7
MULTIDIR-2	42.5±0.6	43.2±0.5	43.9±0.6	45.0±0.6	45.0±0.5	77.8±0.6	77.4±0.7	77.0±0.6	76.3±0.6	75.7±0.9	61.3±0.8	61.5±0.7	61.7±0.6	61.6±0.7	61.8±0.8
MULTIDIR-3	44.2±0.5	44.7±0.5	44.6±0.4	46.4±0.4	46.1±0.5	79.3±0.6	79.0±0.6	78.7±0.6	78.0±0.6	77.7±0.6	61.6±0.8	62.1±0.7	62.4±0.7	62.2±0.7	62.5±0.7
DIRSWITCH-1	45.3±0.4	45.9±0.5	46.9±0.5	48.4±0.4	47.9±0.5	81.6±0.7	81.9±0.6	82.0±0.7	82.3±0.6	81.9±0.7	65.2±0.7	65.7±0.8	66.0±0.8	66.4±0.7	66.4±0.7
DIRSWITCH-2	44.9±0.5	45.9±0.5	46.8±0.5	48.3±0.5	48.0±0.4	81.8±0.7	81.8±0.7	81.8±0.6	81.9±0.7	81.1±0.6	64.0±0.9	64.7±0.8	65.1±0.8	65.6±0.9	65.4±0.7
DIRSWITCH-3	43.8±0.5	45.1±0.5	45.8±0.5	48.3±0.5	47.5±0.5	81.3±0.6	81.4±0.6	80.9±0.7	81.2±0.7	80.4±0.6	61.8±0.8	62.9±0.8	63.1±0.9	64.4±0.8	63.9±0.9

We observe distinct differences between homophilic and heterophilic datasets in the message-passing results in Table 5. The three homophilic datasets—Cora-ML, Citeseer, and Cora (leftmost in the table)—experience a 15–26% accuracy improvement when using undirected edges. In contrast, the heterophilic Roman Empire and Arxiv Year datasets show no significant improvement over default directions, while Snap Patents exhibits a 19% accuracy drop. Pokec is a partial exception, where undirected edges yield benefits with Geom-4, a case analysed in Appendix G.

On the other hand, heterophilic datasets such as Arxiv Year and Snap Patents see substantial accuracy gains as r increases in DirSwitch- r and MultiDir- r . For instance, DirSwitch-3 improves accuracy by +17.9% on Snap Patents, compared to +4.1% for DirSwitch-1. Meanwhile, in homophilic datasets, DirSwitch’s accuracy remains stable across different r values.

Together, these results indicate that multi-directional representations are more important for heterophilic datasets, whereas using undirected edges is beneficial in homophilic cases. These findings align with prior observations in supervised graph neural networks for digraphs (Rossi et al., 2023).

This contrasting behaviour can be explained through the lens of reachability. In homophilic datasets, nodes with the same label tend to cluster locally. Removing edge directions improves dispersal and promotes local smoothing, resulting in more similar embeddings for same-label nodes and improved classification performance. This effect is illustrated in Figure 2.

In heterophilic graphs, however, capturing local connectivity patterns that repeat across different regions is essential for distinguishing node roles. Here, multi-directional approaches play a crucial role, as discarding edge direction information can introduce spurious similarities between nodes with distinct, direction-dependent neighbourhoods.

Sparse vs dense graphs for proximity embeddings We next analyse classification accuracies for proximity embeddings in Table 7, observing distinct trends based on graph density. Denser graphs (Fly Larva, EU-Email, and Polblogs) benefit more from multi-directional embeddings, whereas sparser graphs (CoCite, Pubmed, and Cora (subelj)) achieve higher accuracies when dispersal is increased using undirected edges. For instance, the sparse CoCite graph improves by 19% with undirected edges, with a similar 20% improvement

Table 7: Relative improvements in node classification accuracy for proximity ReachNEs with $p = 1024$. The values reflect the maximum accuracy for per row and dataset in Table 6. The top row displays absolute accuracies for the default edge directions, with standard deviations expressed as percentages. The subsequent rows present the relative improvements compared to the top row. The table is structured with the three denser graphs on the left and the three sparser graphs on the right.

EDGE DIRECTIONS	FLY LARVA	EU-EMAIL	POLBLOGS	CoCITE	PUBMED	CORA (SUBELJ)
DEFAULT	50.1±3.9%	66.0±4.6%	82.0±2.9%	40.2±1.2%	73.3±1.0%	59.3±1.4%
UNDIRECTED	+10.6%	+9.8%	+4.5%	+19.0%	+12.4%	+13.2%
MULTIDIR-1	+3.9%	+11.3%	+5.8%	+2.7%	+0.8%	-1.0%
MULTIDIR-2	+8.2%	+14.0%	+8.8%	+11.9%	+6.1%	+4.3%
MULTIDIR-3	+7.5%	+13.4%	+9.5%	+15.5%	+8.2%	+5.4%
DIRSWITCH-1	+11.5%	+14.5%	+5.6%	+20.2%	+12.3%	+12.1%
DIRSWITCH-2	+13.2%	+14.7%	+9.4%	+20.2%	+11.7%	+10.7%
DIRSWITCH-3	+11.1%	+14.0%	+9.5%	+20.2%	+11.0%	+8.6%

for each DirSwitch model. Conversely, the dense Polblogs graph gains only 4.5% from undirected edges but sees a 9.5% boost with DirSwitch-3.

This difference between dense and sparse graphs can be explained by reachability entropy. Even with purely directed edges, entropy remains high for the dense Fly Larva graph (see Figure 6d), indicating that local sinks do not significantly hinder dispersal. A similar pattern is observed for EU-Email and Polblogs in Appendix E, explaining why undirected edges provide less benefit in these graphs.

Instead, increasing embedding distinguishability by representing directed neighbourhoods becomes more important. Comparing default edges and DirSwitch-1 in Table 6, we observe significant accuracy gains when transitioning from single-scale to multi-scale walk length distributions. For example, Polblogs achieves only 61% accuracy with single-scale Geom but improves to 82% with multi-scale Geom-4. Conversely, sparser datasets show only minor improvements with multi-scale embeddings, with dispersal remaining the primary factor influencing embedding quality.

Further analysis Additional discussion of these results is provided in Appendix G, where we examine the accuracy improvements from multi-scale embeddings, the diminishing returns of repeated concatenation and perform further analysis of the results for the Pokec dataset. We also present evaluation results using unsupervised community detection as an alternative to node classification.

6.4 Comparison to state-of-the-art digraph proximity embeddings

We compare DirSwitch with proximity ReachNEs embeddings against state-of-the-art unsupervised digraph proximity embedding approaches discussed in Section 5.3. Hyperparameters for each model are optimized via cross-validation grid search, with the search grids and best values reported in Appendix D.2.

Table 8: Node classification accuracies for proximity node embedding models. Bold indicate the top accuracy, and results within one standard deviation, for each dataset. Average and standard deviations are calculated over 3x repeated 5-fold cross validations, and 5 different random seeds.

MODEL	FLYLARVA	EU-EMAIL	POLBLOGS	CoCITE	PUBMED	CORA (SUBELJ)
HOPE ¹	35.8 ± 1.9	49.3 ± 3.6	83.2 ± 3.1	30.4 ± 1.2	65.3 ± 0.9	34.0 ± 0.8
APP ²	42.8 ± 2.3	76.3 ± 2.8	88.9 ± 1.7	46.7 ± 0.5	81.2 ± 0.6	65.6 ± 0.7
NERD ³	40.3 ± 2.2	71.3 ± 3.2	89.4 ± 1.7	28.1 ± 0.5	76.5 ± 0.7	45.7 ± 0.6
DGGAN ⁴	16.5 ± 1.4	11.7 ± 2.5	54.6 ± 2.9	17.7 ± 0.4	41.5 ± 0.9	6.2 ± 0.3
BLADE ⁵	49.2 ± 2.2	66.7 ± 3.0	89.3 ± 1.9	24.2 ± 0.5	60.2 ± 1.2	27.7 ± 1.0
DIRSWITCH	56.8 ± 2.2	75.6 ± 2.7	90.0 ± 1.7	48.3 ± 0.5	82.2 ± 0.6	66.6 ± 0.7

¹OU ET AL. (2016) ²ZHOU ET AL. (2017) ³KHOSLA ET AL. (2020) ⁴ZHU ET AL. (2021B) ⁵VIRINCHI & SALADI (2023)

Table 9: Node classification accuracy for self-supervised GraphSAGE embeddings using training via GraphMAEv2 and CCA-SSG losses. Bold indicate the top accuracy for each dataset. Average and standard deviations are calculated over 3x repeated 5-fold cross validations and 5 different seeds.

(a) GraphMAEv2

EDGE DIRECTIONS	ROMAN EMPIRE	ARXIV YEAR	POKEC	SNAP PATENTS	CORA-ML	CITSEER	CORA	ARXIV
DEFAULT	66.09±0.79	44.49±0.24	64.62±0.20	45.68±0.26	79.61±1.76	80.13±1.28	59.48±0.61	59.81±0.25
UNDIRECTED	68.17±0.74	40.08±0.26	65.99±0.52	33.82±0.07	86.03±1.42	90.71±0.92	68.54±0.74	70.49±0.23
ROSSI ET AL. (2023)	67.05±1.85	46.38±0.34	57.93±0.55	61.34±0.28	81.88±3.94	87.53±2.30	34.35±4.42	65.53±0.52
DIRSWITCH-1	71.76±0.76	41.28±0.28	67.16±0.34	39.49±0.30	85.79±1.48	89.03±1.19	68.01±0.83	69.95±0.21
DIRSWITCH-3	71.90±0.78	46.65±0.31	67.57±0.25	57.99±0.21	85.24±1.37	86.67±1.06	66.31±0.80	68.40±0.23

(b) CCA-SSG

EDGE DIRECTIONS	ROMAN EMPIRE	ARXIV YEAR	POKEC	SNAP PATENTS	CORA-ML	CITSEER	CORA	ARXIV
DEFAULT	53.73±0.94	45.01±0.26	58.07±0.27	46.37±0.25	61.70±3.82	69.72±1.77	39.70±0.94	36.73±1.46
UNDIRECTED	53.85±2.45	35.82±0.28	60.38±0.74	38.60±0.13	82.64±1.59	80.71±1.36	52.35±1.30	45.93±1.18
ROSSI ET AL. (2023)	53.74±0.75	46.83±0.44	60.54±0.50	65.55±0.28	60.44±2.94	64.64±2.41	31.28±2.21	34.92±0.54
DIRSWITCH-1	56.77±2.36	45.19±0.35	66.23±0.70	49.84±0.60	82.00±1.76	77.35±2.53	52.70±1.28	61.63±0.38
DIRSWITCH-3	66.39±1.24	52.64±0.30	67.25±0.21	64.08±0.19	82.27±1.42	85.26±1.32	61.07±0.80	63.58±0.30

Following hyperparameter tuning, test accuracy is evaluated using the best hyperparameter setting on new 3x repeated 5-fold CV splits, with five different random seeds per model. We use the official implementations for all baselines except BLADE, which we reimplement as its source code is not publicly available.

Table 8 presents the average test accuracies and standard deviations for each proximity embedding dataset. DirSwitch consistently achieves the highest average accuracy across datasets or performs within one standard deviation of the best model (EU-Email). Notably, DirSwitch outperforms competing approaches by a significant margin on Fly Larva, achieving 57% accuracy compared to 49% for BLADE. These results highlight DirSwitch’s ability to leverage both multi-directional and multi-scale embeddings while avoiding local sinks, distinguishing it from prior proximity embedding models.

6.5 Applying DirSwitch to self-supervised message-passing graph neural networks

DirSwitch is not limited to the ReachNEs framework and can also be applied to message-passing graph neural networks (GNNs). The most natural integration is with the GraphSAGE model (Hamilton et al., 2017), given its structural similarity to ReachNEs. By defining $\mathbf{H}^{(0)} = \mathbf{Z}^{(0)} = \mathbf{X}$, we can represent the DirSwitch-GraphSAGE model using the following iterative formulation, analogous to Equation 6:

$$\text{For } k \in \{1, 2, \dots\}, \quad \begin{cases} \mathbf{H}^{(k)} &= \mathbf{A}_*^\top \mathbf{H}^{(k-1)}, \\ \mathbf{Z}^{(k)} &= \sigma \left(\mathbf{Z}^{(k-1)} \mathbf{W}_1^{(k-1)} + \mathbf{H}^{(k)} \mathbf{W}_2^{(k-1)} \right), \end{cases} \quad (10)$$

The key differences are the weight matrices \mathbf{W}_1 and \mathbf{W}_2 , which replace the walk length probability coefficients, and the non-linear activation function σ .

For training the weight matrices, we employ self-supervised learning, the current state-of-the-art approach for unsupervised GNN training. These loss functions typically combine reconstruction and contrastive learning objectives. In our experiments, we use two recent and efficient methods: GraphMAEv2 (Hou et al., 2023) and CCA-SSG (Zhang et al., 2021a). We adopt the default hyperparameters for both loss functions and optimizers, as recommended by the original implementations. We train DirSwitch-1 and DirSwitch-3 models but exclude DirSwitch-2 due to the high computational cost of GNN training.

As baselines, we compare against GraphSAGE with default (0) and undirected (U) edges, as well as the digraph extension of GraphSAGE proposed by Rossi et al. (2023). Like DirSwitch, Rossi et al. (2023) incorporates multi-directional neighbourhoods but does not include a switch to undirected edges.

Unlike ReachNEs, GNNs typically use a small number of aggregation steps to balance computational efficiency and mitigate over-smoothing (Chen et al., 2020). We set $K = 4$ aggregation steps, following the default in GraphMAEv2 (Hou et al., 2023).

Table 9 presents the results for GraphMAEv2 and CCA-SSG. As observed in the ReachNEs experiments (Section 6.3), performance varies qualitatively between heterophilic and homophilic datasets, with undirected edges generally performing better in homophilic cases. This is particularly evident in Cora, where DirSwitch and undirected edges significantly outperform both the default edges and Rossi et al. (2023).

The only dataset where Rossi et al. (2023) surpasses DirSwitch in accuracy is Snap Patents. As shown in Table 5, using undirected edges causes a substantial accuracy drop for Snap Patents in ReachNEs, and the same trend is observed with GraphSAGE.

7 Limitations and future work

This work addresses two key challenges in digraph embedding: the dispersal-obstructing effects of local sinks and the need for expressive embeddings that capture multiple directed neighbourhoods. A limitation of our approach is its reliance on embedding concatenation to achieve both multi-directionality and multi-scale representation, which incurs notable computational overhead. While our results suggest that few directed neighbourhoods and scales are often sufficient—due to the typically short characteristic path lengths of real-world networks (Watts & Strogatz, 1998)—further work is needed to better quantify diminishing returns and to explore more efficient alternatives to concatenation. Additionally, other structural properties of digraphs, such as the role of cycles and self-loops, may be equally important and deserve further investigation.

For tractability, our analysis considered graphs with a single type of supplementary information: node attributes. In reality, graphs are often far more complex, incorporating multiple edge types, attributed edges (as in knowledge graphs (Wang et al., 2017a)), or temporal dynamics where edges evolve over time. While we expect the DirSwitch principle, i.e., separating short-range and long-range behaviours, to generalize to these cases, addressing expressivity in multi-modal graphs will likely require specialized techniques beyond our current concatenation approach.

Another limitation lies in the embedding reduction methods explored. We focused on two widely used techniques: node attribute smoothing for message-passing embeddings and SVD for proximity embeddings. However, other reduction methods, such as structural embeddings, remain unexplored in our analysis. Investigating how local sinks and multi-directional embeddings affect alternative reduction techniques is an important avenue for future research, especially since much of the existing work on structural embeddings is limited to undirected graphs (Donnat et al., 2018; Zhu et al., 2021a; Jin et al., 2021).

Finally, while our benchmark evaluation demonstrates DirSwitch’s effectiveness, its real-world utility in unsupervised tasks such as node clustering or anomaly detection remains underexplored. These tasks often lack well-defined ground-truth labels, making evaluation challenging and necessitating further research. Expanding DirSwitch to these domains presents exciting opportunities for future work.

8 Conclusion

In this paper, we analysed unsupervised node embedding learning on digraphs through our reachability-based random walk filter framework, ReachNEs. Our analysis identified two key challenges: local sinks obstruct information propagation, while embeddings must also represent directed neighbourhoods to be expressive.

To address these issues, we introduced DirSwitch. By decoupling local and global smoothing, DirSwitch preserves fine-grained directed structure through multi-directional embedding concatenation, while treating the graph as undirected for long-range interactions to mitigate the impact of local sinks.

We demonstrated that DirSwitch significantly enhances embedding quality, achieving higher accuracy on standard node classification benchmarks. Additionally, we showcased its practical effectiveness by outperforming state-of-the-art unsupervised proximity embedding models. Finally, we highlighted DirSwitch’s broad applicability by using it to generalize self-supervised graph neural networks to digraphs, illustrating its flexibility across different embedding paradigms.

Acknowledgments

This work was partially supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation. We would also like to thank Petra Poklukar, Marcus Klasson, Aniss Medbouhi, and Miguel Vasco for their valuable feedback and support in improving this paper.

References

- Lada A. Adamic and Natalie Glance. The Political Blogosphere and the 2004 U.S. Election: Divided They Blog. In *LinkKDD'05*, pp. 36–43, 2005. URL <https://doi.org/10.1145/1134271.1134277>.
- David Arthur and Sergei Vassilvitskii. K-means++: The advantages of careful seeding. In *SODA'07*, pp. 1027–1035, 2007. URL <https://dl.acm.org/doi/10.5555/1283383.1283494>.
- Albert-László Barabási and Réka Albert. Emergence of Scaling in Random Networks. *Science*, 286(5439): 509–512, 1999. URL <https://www.science.org/doi/abs/10.1126/science.286.5439.509>.
- Aleksandar Bojchevski and Stephan Günnemann. Deep Gaussian Embedding of Graphs: Unsupervised Inductive Learning via Ranking. In *ICLR'18*, 2018. URL <https://openreview.net/forum?id=r1ZdKJ-OW>.
- Stephen P. Borgatti and Martin G. Everett. Notions of Position in Social Network Analysis. *Sociological methodology*, 22:1–35, 1992. URL <https://doi.org/10.2307/270991>.
- Michael M. Bronstein, Joan Bruna, Yann LeCun, Arthur Szlam, and Pierre Vandergheynst. Geometric Deep Learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 34(4):18–42, 2017. URL <https://doi.org/10.1109/MSP.2017.2693418>.
- Ciwan Ceylan, Kambiz Ghooarchian, and Danica Kragic. Digraphwave: Scalable Extraction of Structural Node Embeddings via Diffusion on Directed Graphs, 2022. URL <https://doi.org/10.48550/arXiv.2207.10149>.
- Sudhanshu Chanpuriya and Cameron Musco. Infinitewalk: Deep Network Embeddings as Laplacian Embeddings with a Nonlinearity. In *KDD'20*, pp. 1325–1333, 2020. URL <https://doi.org/10.1145/3394486.3403185>.
- Deli Chen, Yankai Lin, Wei Li, Peng Li, Jie Zhou, and Xu Sun. Measuring and Relieving the Over-Smoothing Problem for Graph Neural Networks from the Topological View. *AAAI'20*, 34(04):3438–3445, 2020. URL <https://doi.org/10.1609/aaai.v34i04.5747>.
- Fan R.K. Chung. *Spectral graph theory*. Number 92 in Regional conference series in mathematics. American Mathematical Society, 1997. URL <https://doi.org/10.1090/cbms/092>.
- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Velickovic. Principal Neighbourhood Aggregation for Graph Nets. In *NeurIPS'20*, 2020. URL <https://doi.org/10.48550/arXiv.2004.05718>.
- Thomas M. Cover and Joy A. Thomas. *Elements of Information Theory*. John Wiley & Sons, Ltd, 2nd edition, 2005. URL <https://doi.org/10.1002/047174882X>.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *NeurIPS'16*, 2016. URL <https://proceedings.neurips.cc/paper/2016/file/04df4d434d481c5bb723be1b6df1ee65-Paper.pdf>.
- Claire Donnat, Marinka Zitnik, David Hallac, and Jure Leskovec. Learning Structural Node Embeddings Via Diffusion Wavelets. In *KDD'18*, pp. 1320–1329, 2018. URL <https://doi.org/10.1145/3219819.3220025>.
- Giorgio Fagiolo. Clustering in complex directed networks. *Physical Review E*, 76(2), 2007. URL <https://doi.org/10.1103/PhysRevE.76.026107>.

- Michalis Faloutsos, Petros Faloutsos, and Christos Faloutsos. On power-law relationships of the internet topology. *ACM SIGCOMM Computer Communication Review*, 29(4):251–262, 1999. URL <https://doi.org/10.1145/316194.316229>.
- Matthias Fey and Jan E. Lenssen. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019. URL <https://doi.org/10.48550/arXiv.1903.02428>.
- Michael Fire and Carlos Guestrin. The rise and fall of network stars: Analyzing 2.5 million graphs to reveal how high-degree vertices emerge over time. *Information Processing and Management: an International Journal*, 57(2), 2020. URL <https://doi.org/10.1016/j.ipm.2019.05.002>.
- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural Message Passing for Quantum Chemistry. In *ICML’17*, pp. 1263–1272, 2017. URL <https://proceedings.mlr.press/v70/gilmer17a.html>.
- Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, 4th edition, 2013. URL <https://epubs.siam.org/doi/book/10.1137/1.9781421407944>.
- Martin Grohe and Pascal Schweitzer. The graph isomorphism problem. *Communications of the ACM*, 63(11):128–134, 2020. URL <https://doi.org/10.1145/3372123>.
- Martin Grohe, Kristian Kersting, Martin Mladenov, and Pascal Schweitzer. Color refinement and its applications. In *An Introduction to Lifted Probabilistic Inference*. The MIT Press, 08 2021. URL <https://doi.org/10.7551/mitpress/10548.003.0023>.
- Aditya Grover and Jure Leskovec. node2vec: Scalable Feature Learning for Networks. In *KDD’16*, pp. 855–864, 2016. URL <https://doi.org/10.1145/2939672.2939754>.
- Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding Structure with Randomness: Probabilistic Algorithms for Constructing Approximate Matrix Decompositions. *SIAM Review*, 53(2):217–288, 2011. URL <https://doi.org/10.1137/090771806>.
- Brian C. Hall. The Matrix Exponential. In *Lie Groups, Lie Algebras, and Representations*. Springer Cham, 2015. URL <https://doi.org/10.1007/978-3-319-13467-3>.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive Representation Learning on Large Graphs. In *NeurIPS’17*, volume 30, 2017. URL <https://doi.org/10.48550/arXiv.1706.02216>.
- Kaveh Hassani and Amir Hosein Khasahmadi. Contrastive Multi-View Representation Learning on Graphs. In *ICML’20*, pp. 4116–4126, 2020. URL <https://proceedings.mlr.press/v119/hassani20a.html>.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer New York, 2nd edition, 2009. URL <https://doi.org/10.1007/978-0-387-84858-7>.
- Mark Heimann, Haoming Shen, Tara Safavi, and Danai Koutra. Regal: Representation Learning-Based Graph Alignment. In *CIKM’18*, pp. 117–126, 2018. URL <https://doi.org/10.1145/3269206.3271788>.
- Nicholas J. Higham. *Functions of Matrices*. Society for Industrial and Applied Mathematics, 2008. URL <https://doi.org/10.1137/1.9780898717778>.
- Michiel E. Hochstenbach. A Jacobi–Davidson type method for the generalized singular value problem. *Linear Algebra and its Applications*, 431(3):471–487, 2009. URL <https://doi.org/10.1016/j.laa.2009.03.003>.
- Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2nd edition, 2012. URL <https://doi.org/10.1017/9781139020411>.
- Zhenyu Hou, Yufei He, Yukuo Cen, Xiao Liu, Yuxiao Dong, Evgeny Kharlamov, and Jie Tang. GraphMAE2: A Decoding-Enhanced Masked Self-Supervised Graph Learner. In *WWW’23*, pp. 737–746, 2023. URL <https://doi.org/10.1145/3543507.3583379>.

- Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. Open Graph Benchmark: Datasets for Machine Learning on Graphs. In *NeurIPS'20*, pp. 22118–22133, 2020a. URL <https://doi.org/10.48550/arXiv.2005.00687>.
- Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for Pre-training Graph Neural Networks. In *ICLR'20*, 2020b. URL <https://openreview.net/forum?id=HJlWWJSFDH>.
- Elvin Isufi, Fernando Gama, David I Shuman, and Santiago Segarra. Graph Filters for Signal Processing and Machine Learning on Graphs. *IEEE Transactions on Signal Processing*, 72:4745–4781, 2024. URL <https://doi.org/10.1109/TSP.2024.3349788>.
- Di Jin, Ryan A. Rossi, Eunye Koh, Sungchul Kim, Anup Rao, and Danai Koutra. Latent Network Summarization: Bridging Network Embedding and Summarization. In *KDD'19*, pp. 987–997, 2019. URL <https://doi.org/10.1145/3292500.3330992>.
- Junchen Jin, Mark Heimann, Di Jin, and Danai Koutra. Toward Understanding and Evaluating Structural Node Embeddings. *ACM TKDD*, 16(3), 2021. URL <https://doi.org/10.1145/3481639>.
- Megha Khosla, Jurek Leonhardt, Wolfgang Nejdl, and Avishek Anand. Node Representation Learning for Directed Graphs. In *ECML PKDD 2019*, pp. 395–411, 2020. URL https://doi.org/10.1007/978-3-030-46150-8_24.
- Thomas N. Kipf and Max Welling. Semi-supervised Classification with Graph Convolutional Networks. In *ICLR'17*. OpenReview, 2017. URL <https://openreview.net/pdf?id=SJU4ayYgl>.
- Bryan Klimt and Yiming Yang. The Enron Corpus: A New Dataset for Email Classification Research. In *ECML'04*. Springer Berlin Heidelberg, 2004. URL https://doi.org/10.1007/978-3-540-30115-8_22.
- Yehuda Koren, Robert Bell, and Chris Volinsky. Matrix Factorization Techniques for Recommender Systems. *Computer*, 42(8):30–37, 2009. URL <https://doi.org/10.1109/MC.2009.263>.
- Jure Leskovec, Daniel Huttenlocher, and Jon Kleinberg. Signed Networks in Social Media. In *CHI'10*, 2010. URL <https://doi.org/10.1145/1753326.1753532>.
- Pan Li and Jure Leskovec. The Expressive Power of Graph Neural Networks. In Lingfei Wu, Peng Cui, Jian Pei, and Liang Zhao (eds.), *Graph Neural Networks: Foundations, Frontiers, and Applications*, chapter 5. Springer Singapore, 1st edition, 2022. URL https://doi.org/10.1007/978-981-16-6054-2_5.
- Derek Lim, Felix Hohne, Xiuyu Li, Sijia Linda Huang, Vaishnavi Gupta, Omkar Bhalerao, and Ser-Nam Lim. Large Scale Learning on Non-Homophilous Graphs: New Benchmarks and Strong Simple Methods. In *NeurIPS'21*, 2021. URL <https://openreview.net/forum?id=DfGu8WwT0d>.
- Andrew Kachites McCallum, Kamal Nigam, Jason Rennie, and Kristie Seymore. Automating the construction of internet portals with machine learning. *Information Retrieval*, 3:127–163, 2000. URL <https://doi.org/10.1023/A:1009953814988>.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed Representations of Words and Phrases and their Compositionality. In *NeurIPS'13*, 2013. URL https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L. Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and Leman Go Neural: Higher-Order Graph Neural Networks. In *AAAI'19*, 2019. URL <https://doi.org/10.1609/aaai.v33i01.33014602>.
- Kevin P. Murphy. *Machine Learning: A Probabilistic Perspective*. MIT press, 2012. URL <https://mitpress.mit.edu/9780262018029/>.
- Mark Newman. *Networks*. Oxford university press, 2nd edition, 2018. URL <https://doi.org/10.1093/oso/9780198805090.001.0001>.

- Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. Asymmetric Transitivity Preserving Graph Embedding. In *KDD'16*, pp. 1105–1114, 2016. URL <https://doi.org/10.1145/2939672.2939751>.
- Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The PageRank Citation Ranking: Bringing Order to the Web. Technical Report 1999-66, Stanford InfoLab, 1999. URL <http://ilpubs.stanford.edu:8090/422/>.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *NeurIPS'19*, pp. 8024–8035, 2019. URL <https://dl.acm.org/doi/10.5555/3454287.3455008>.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. DeepWalk: online learning of social representations. In *KDD'14*, pp. 701–710. Association for Computing Machinery, 2014. URL <https://doi.org/10.1145/2623330.2623732>.
- Kaare Brandt Petersen and Michael Syskind Pedersen. The Matrix Cookbook, 2012. URL <http://www2.compute.dtu.dk/pubdb/pubs/3274-full.html>. Version 20121115.
- Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. Network Embedding as Matrix Factorization: Unifying DeepWalk, LINE, PTE, and Node2vec. In *WSDM'18*, pp. 459–467, 2018. URL <https://doi.org/10.1145/3159652.3159706>.
- Emanuele Rossi, Bertrand Charpentier, Francesco Di Giovanni, Fabrizio Frasca, Stephan Günnemann, and Michael M Bronstein. Edge Directionality Improves Learning on Heterophilic Graphs. In *LoG'23*, 2023. URL <https://doi.org/10.48550/arXiv.2305.10498>.
- Ryan A. Rossi, Di Jin, Sungchul Kim, Nesreen K. Ahmed, Danai Koutra, and John Boaz Lee. On Proximity and Structural Role-Based Embeddings in Networks: Misconceptions, Techniques, and Applications. *ACM TKDD*, 14(5), 2020. URL <https://doi.org/10.1145/3397191>.
- Benedek Rozemberczki, Carl Allen, and Rik Sarkar. Multi-Scale attributed node embedding. *Journal of Complex Networks*, 9(2), 2021. URL <https://doi.org/10.1093/comnet/cnab014>.
- Tim Sainburg, Leland McInnes, and Timothy Q Gentner. Parametric UMAP Embeddings for Representation and Semisupervised learning. *Neural Computation*, 33(11):2881–2907, 2021. URL https://doi.org/10.1162/neco_a_01434.
- Konstantinos Skitsas, Karol Orlowski, Judith Hermanns, Davide Mottin, and Panagiotis Karras. Comprehensive Evaluation of Algorithms for Unrestricted Graph Alignment. In *EDBT'23*, 2023. URL <https://dx.doi.org/10.48786/edbt.2023.21>.
- Balasubramaniam Srinivasan and Bruno Ribeiro. On the equivalence between positional node embeddings and structural graph representations. In *ICLR'20*, 2020. URL <https://openreview.net/forum?id=SJxzFySKwH>.
- Lubos Takac and Michal Zabovsky. Data analysis in public social networks. In *International scientific conference and international workshop present day trends of innovations*, volume 6, 2012. URL <https://snap.stanford.edu/data/soc-Pokec.html>.
- Shantanu Thakoor, Corentin Tallec, Mohammad Gheshlaghi Azar, Mehdi Azabou, Eva L. Dyer, Remi Munos, Petar Veličković, and Michal Valko. Large-Scale Representation Learning on Graphs via Bootstrapping. In *ICLR'22*, 2022. URL <https://doi.org/10.48550/arXiv.2102.06514>.
- Zekun Tong, Yuxuan Liang, Changsheng Sun, Xinke Li, David Rosenblum, and Andrew Lim. Digraph Inception Convolutional Networks. In *NeurIPS'20*, volume 33, pp. 17907–17918, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/cffb6e2288a630c2a787a64ccc67097c-Abstract.html>.

- J. J. Peter Veerman and Robert Lyons. A Primer on Laplacian Dynamics in Directed Graphs, 2020. URL <https://doi.org/10.48550/arXiv.2002.02605>.
- Srinivas Virinchi and Anoop Saladi. BLADE: Biased Neighborhood Sampling based Graph Neural Network for Directed Graphs. In *WSDM'23*, pp. 42–50, 2023. URL <https://dl.acm.org/doi/10.1145/3539597.3570430>.
- S.V.N. Vishwanathan, Nicol N. Schraudolph, Risi Kondor, and Karsten M. Borgwardt. Graph Kernels. *Journal of Machine Learning Research*, 11(40):1201–1242, 2010. URL <http://jmlr.org/papers/v11/vishwanathan10a.html>.
- Ulrike von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007. URL <https://doi.org/10.1007/s11222-007-9033-z>.
- Haonan Wang, Jieyu Zhang, Qi Zhu, Wei Huang, Kenji Kawaguchi, and Xiaokui Xiao. Single-Pass Contrastive Learning Can Work for Both Homophilic and Heterophilic Graph. *Transactions on Machine Learning Research*, 2023. URL <https://openreview.net/forum?id=244KePn09i>.
- Quan Wang, Zhendong Mao, Bin Wang, and Li Guo. Knowledge Graph Embedding: A Survey of Approaches and Applications. *IEEE Transactions on Knowledge and Data Engineering*, 29(12):2724–2743, 2017a. URL <https://doi.org/10.1109/TKDE.2017.2754499>.
- Xiao Wang, Peng Cui, Jing Wang, Jian Pei, Wenwu Zhu, and Shiqiang Yang. Community Preserving Network Embedding. *AAAI'17*, 35(1), 2017b. URL <https://ojs.aaai.org/index.php/AAAI/article/view/10488>.
- Duncan J. Watts and Steven H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 393(6684): 440–442, jun 1998. ISSN 1476-4687. doi: 10.1038/30918. URL <https://www.nature.com/articles/30918>.
- Michael Winding, Benjamin D. Pedigo, Christopher L. Barnes, Heather G. Patsolic, Youngser Park, Tom Kazimiers, Akira Fushiki, Ingrid V. Andrade, Avinash Khandelwal, Javier Valdes-Aleman, Feng Li, Nadine Randel, Elizabeth Barsotti, Ana Correia, Richard D. Fetter, Volker Hartenstein, Carey E. Priebe, Joshua T. Vogelstein, Albert Cardona, and Marta Zlatić. The connectome of an insect brain. *Science*, 379 (6636), 2023. URL <https://doi.org/10.1126/science.add9330>.
- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying Graph Convolutional Networks. In *ICML'19*, pp. 6861–6871, 2019. URL <https://doi.org/10.48550/arXiv.1902.07153>.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How Powerful are Graph Neural Networks? In *International Conference on Learning Representations, ICLR'19*. OpenReview.net, 2019. URL <https://openreview.net/forum?id=ryGs6iA5Km>.
- Yuchen Yan, Yongyi Hu, Qinghai Zhou, Lihui Liu, Zhichen Zeng, Yuzhong Chen, Menghai Pan, Huiyuan Chen, Mahashweta Das, and Hanghang Tong. PaCER: Network Embedding From Positional to Structural. In *WWW'24*, pp. 2485–2496. Association for Computing Machinery, 2024. URL <https://doi.org/10.1145/3589334.3645516>.
- Hao Yin, Austin R. Benson, Jure Leskovec, and David F. Gleich. Local Higher-Order Graph Clustering. In *KDD'17*, pp. 555–564, 2017. URL <https://doi.org/10.1145/3097983.3098069>.
- Wenjian Yu, Yu Gu, Jian Li, Shenghua Liu, and Yaohang Li. Single-Pass PCA of Large High-Dimensional Data. In *IJCAI'17*, pp. 3350–3356, 2017. URL <https://doi.org/10.24963/ijcai.2017/468>.
- Hengrui Zhang, Qitian Wu, Junchi Yan, David Wipf, and S Yu Philip. From Canonical Correlation Analysis to Self-supervised Graph Neural Networks. In *NeurIPS'21*, 2021a. URL <https://doi.org/10.48550/arXiv.2106.12484>.

- Xitong Zhang, Yixuan He, Nathan Brugnone, Michael Perlmuter, and Matthew Hirn. MagNet: A Neural Network for Directed Graphs. In *NeurIPS'21*, volume 34, pp. 27003–27015, 2021b. URL <https://doi.org/10.48550/arXiv.2102.11391>.
- Chang Zhou, Yuqiong Liu, Xiaofei Liu, Zhongyi Liu, and Jun Gao. Scalable Graph Embedding for Asymmetric Proximity. *AAAI'17*, 31(1), 2017. URL <https://doi.org/10.1609/aaai.v31i1.10878>.
- Jing Zhu, Xingyu Lu, Mark Heimann, and Danai Koutra. Node Proximity is All You Need: Unified Structural and Positional Node and Graph Embedding. In *SDM'21*, pp. 163–171, 2021a. URL <https://doi.org/10.1137/1.9781611976700.19>.
- Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Beyond Homophily in Graph Neural Networks: Current Limitations and Effective Designs. In *NeurIPS'20*, 2020. URL <https://doi.org/10.48550/arXiv.2006.11468>.
- Shijie Zhu, Jianxin Li, Hao Peng, Senzhang Wang, and Lifang He. Adversarial Directed Graph Embedding. *AAAI'21*, 35(5):4741–4748, 2021b. URL <https://doi.org/10.1609/aaai.v35i5.16605>.

A Additional details on ReachNEs and DirSwitch

A.1 Verification of column stochasticity of the random walk normalized adjacency matrices

In Section 2.1, we defined the random walk normalized adjacency matrices as follows:

$$A_{0i,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg_0(j) = 0, \\ \frac{A_{i,j}}{\deg_0(j)} & \text{otherwise,} \end{cases} \quad A_{Ii,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg_I(j) = 0, \\ \frac{A_{j,i}}{\deg_I(j)} & \text{otherwise,} \end{cases}$$

$$A_{Uj,j} = \begin{cases} 1 & \text{if } i = j \text{ and } \deg(j) = 0, \\ \frac{A_{undir i,j}}{\deg(j)} & \text{otherwise.} \end{cases}$$

These definitions ensure that each matrix is column stochastic, meaning that each column sums to one. We verify this property below:

$$\sum_{k=1}^n A_{0k,j} = \begin{cases} 1 & \text{if } \deg_0(j) = 0 \\ \frac{1}{\deg_0(j)} \sum_{k=1}^n A_{k,j} & \text{if } \deg_0(j) > 0 \end{cases} = \begin{cases} 1 & \text{if } \deg_0(j) = 0 \\ \frac{\deg_0(j)}{\deg_0(j)} & \text{if } \deg_0(j) > 0 \end{cases} = 1,$$

$$\sum_{k=1}^n A_{Ik,j} = \begin{cases} 1 & \text{if } \deg_I(j) = 0 \\ \frac{1}{\deg_I(j)} \sum_{k=1}^n A_{j,k} & \text{if } \deg_I(j) > 0 \end{cases} = \begin{cases} 1 & \text{if } \deg_I(j) = 0 \\ \frac{\deg_I(j)}{\deg_I(j)} & \text{if } \deg_I(j) > 0 \end{cases} = 1,$$

$$\sum_{k=1}^n A_{Uk,j} = \begin{cases} 1 & \text{if } \deg(j) = 0 \\ \frac{1}{\deg(j)} \sum_{k=1}^n A_{undirk,j} & \text{if } \deg(j) > 0 \end{cases} = \begin{cases} 1 & \text{if } \deg(j) = 0 \\ \frac{\deg(j)}{\deg(j)} & \text{if } \deg(j) > 0 \end{cases} = 1.$$

A.2 Message-passing embeddings for directed edges

In Section 2.3.2, we highlighted the connection between the reachability message-passing reduction using $\mathbf{R}(\mathbf{A}_U)$ and the undirected message-passing graph neural network GraphSAGE (Hamilton et al., 2017) in Equation 7. Here, we present the corresponding formulas for directed cases:

$$\mathbf{A}_0^\top \mathbf{H} = \mathbf{D}_0^{-1} \mathbf{A}^\top \mathbf{H}, \quad [\mathbf{A}_0^\top \mathbf{H}]_{i,:} = \sum_{j=1}^n A_{0j,i} \mathbf{H}_{j,:} = \frac{1}{\deg_0(i)} \sum_{j:(i,j) \in \mathbb{M}} \mathbf{H}_{j,:}, \quad (11)$$

$$\mathbf{A}_I^\top \mathbf{H} = \mathbf{D}_I^{-1} \mathbf{A} \mathbf{H}, \quad [\mathbf{A}_I^\top \mathbf{H}]_{i,:} = \sum_{j=1}^n A_{Ij,i} \mathbf{H}_{j,:} = \frac{1}{\deg_I(i)} \sum_{j:(j,i) \in \mathbb{M}} \mathbf{H}_{j,:}. \quad (12)$$

Note that in Equation 11, messages are aggregated over *outgoing* edges of node i , meaning that they flow opposite to edge directions. Conversely, in Equation 12, messages are aggregated along the *incoming*, following the edge directions.

A.3 Proximity embeddings and structural equivalence

We extend Section 2.3.1 by deriving Equation 4, which relates the distances between reachability vectors in \mathbf{R} to distances between embedding vectors in $\mathbf{Z} = [\mathbf{Z}_U \ \mathbf{Z}_V]$. For convenience, we restate the equation:

$$\|\mathbf{R}_{i,:} - \mathbf{R}_{j,:}\|_2^2 + \|\mathbf{R}_{:,i} - \mathbf{R}_{:,j}\|_2^2 = \left\| (\mathbf{Z}_{i,:} - \mathbf{Z}_{j,:}) \sqrt{\hat{\Sigma}} \right\|_2^2, \quad \hat{\Sigma} = \begin{bmatrix} \Sigma_{:,q,:q} & \mathbf{0} \\ \mathbf{0} & \Sigma_{:,q,:q} \end{bmatrix}. \quad (13)$$

Assuming $\text{rank}(\mathbf{R}) = q$, the $n - q$ smallest singular values vanish, leading to:

$$\mathbf{R} = \mathbf{U}_{:,q} \Sigma_{:,q,:q} (\mathbf{V}_{:,q})^\top = \mathbf{U}_{:,q} \sqrt{\Sigma_{:,q,:q}} \left(\mathbf{V}_{:,q} \sqrt{\Sigma_{:,q,:q}} \right)^\top = \mathbf{Z}_U \mathbf{Z}_V^\top.$$

We can then express each row of \mathbf{R} as the inner product between the corresponding row in \mathbf{Z}_U and \mathbf{Z}_V , and similarly, each column of \mathbf{R} can be expressed as the inner product between a row in \mathbf{Z}_V and \mathbf{Z}_U :

$$\mathbf{R}_{i,:} = \mathbf{Z}_{U i,:} \mathbf{Z}_V^\top, \quad \mathbf{R}_{:,i} = \mathbf{Z}_U (\mathbf{Z}_{V i,:})^\top.$$

We now derive Equation 13, using the simplified notation $\mathbf{U} = \mathbf{U}_{:,q}$, $\mathbf{\Sigma} = \mathbf{\Sigma}_{:,q,q}$, and $\mathbf{V} = \mathbf{V}_{:,q}$:

$$\|\mathbf{R}_{i,:} - \mathbf{R}_{j,:}\|_2^2 + \|\mathbf{R}_{:,i} - \mathbf{R}_{:,j}\|_2^2 = \|\mathbf{Z}_{U i,:} \mathbf{Z}_V^\top - \mathbf{Z}_{U j,:} \mathbf{Z}_V^\top\|_2^2 + \|\mathbf{Z}_U (\mathbf{Z}_{V i,:})^\top - \mathbf{Z}_U (\mathbf{Z}_{V j,:})^\top\|_2^2 \quad (14)$$

$$= \|(\mathbf{Z}_{U i,:} - \mathbf{Z}_{U j,:}) \mathbf{Z}_V^\top\|_2^2 + \|\mathbf{Z}_U (\mathbf{Z}_{V i,:} - \mathbf{Z}_{V j,:})^\top\|_2^2 \quad (15)$$

$$= \|(\mathbf{Z}_{U i,:} - \mathbf{Z}_{U j,:}) \sqrt{\mathbf{\Sigma}} \mathbf{V}^\top\|_2^2 + \|\mathbf{U} \sqrt{\mathbf{\Sigma}} (\mathbf{Z}_{V i,:} - \mathbf{Z}_{V j,:})^\top\|_2^2 \quad (16)$$

$$= \|(\mathbf{Z}_{U i,:} - \mathbf{Z}_{U j,:}) \sqrt{\mathbf{\Sigma}}\|_2^2 + \|\sqrt{\mathbf{\Sigma}} (\mathbf{Z}_{V i,:} - \mathbf{Z}_{V j,:})^\top\|_2^2 \quad (17)$$

$$= \sum_{k=1}^q \sqrt{\Sigma_{k,k}} (Z_{U i,k} - Z_{U j,k})^2 + \sum_{k=1}^q \sqrt{\Sigma_{k,k}} (Z_{V i,k} - Z_{V j,k})^2 \quad (18)$$

$$= \sum_{k=1}^{2q} \sqrt{\hat{\Sigma}_{k,k}} (Z_{i,k} - Z_{j,k})^2 \quad (19)$$

$$= \|(\mathbf{Z}_{i,:} - \mathbf{Z}_{j,:}) \sqrt{\hat{\mathbf{\Sigma}}}\|_2^2 \quad (20)$$

The transition from Equation 16 to Equation 17 follows from the orthogonality of \mathbf{U} and \mathbf{V} , i.e., $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_q$ and $\mathbf{V}^\top \mathbf{V} = \mathbf{I}_q$.

A.4 Computational complexity of ReachNEs

In this section, we discuss the time and memory complexity of ReachNEs and how to efficiently implement it for embedding learning on large graphs. The primary computational bottleneck is the $n \times n$ reachability matrix \mathbf{R} , which is generally dense. Storing \mathbf{R} explicitly is infeasible for even moderately large graphs—for instance, a graph with one million nodes would require over 3TB of memory.

For message-passing embeddings, the iterative algorithm described in Equation 6 offers a practical way to bypass the need to store \mathbf{R} explicitly. For convenience, we restate it below:

$$\text{For } k \in \{1, 2, \dots\}, \quad \begin{cases} \mathbf{H}^{(k)} &= \mathbf{A}_*^\top \mathbf{H}^{(k-1)}, \\ \mathbf{Z}^{(k)} &= \mathbf{Z}^{(k-1)} + \mathbf{P}_w(k) \mathbf{H}^{(k)}, \end{cases} \quad \mathbf{H}^{(0)} = \mathbf{X}, \quad \mathbf{Z}^{(0)} = \mathbf{P}_w(0) \mathbf{X}. \quad (21)$$

Rather than computing \mathbf{R} in full and multiplying it with the node feature matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$, this iterative formulation directly constructs the embedding matrix $\mathbf{Z}^{(K)} = \mathbf{R}^\top \mathbf{X}$ through sequential summation. This reduces the memory complexity from quadratic to linear, specifically $\mathcal{O}(nd+m)$, where nd accounts for the embedding matrices $\mathbf{H}^{(k)}$ and $\mathbf{Z}^{(k)}$, and m is the number of non-zero entries in the sparse adjacency matrix \mathbf{A}_* .

The time complexity of the iterative message-passing algorithm in Equation 21 is $\mathcal{O}(Kmd + Knd)$, where md arises from the sparse-dense matrix multiplication $\mathbf{A}_*^\top \mathbf{H}^{(k-1)}$, and nd comes from the embedding summation update. Both operations are efficiently supported by existing GPU libraries (Paszke et al., 2019; Fey & Lenssen, 2019), making the method highly practical in large-scale settings.

By contrast, proximity ReachNEs presents greater memory challenges. Both standard and randomized SVD methods (Halko et al., 2011) require multiple passes over all elements of the input matrix, which is problematic since storing \mathbf{R} in memory is intractable, and recomputing it on-the-fly is computationally prohibitive.

To address this, we adopt the *single-pass randomized SVD* (SP-rSVD) algorithm proposed by Yu et al. (2017). SP-rSVD follows the general structure of randomized SVD (rSVD) (Halko et al., 2011), but avoids multiple passes over the input. First, the reachability matrix \mathbf{R} is projected into a lower-dimensional subspace via a

random matrix $\Omega \in \mathbb{R}^{n \times p}$ to compute $\mathbf{G} = \mathbf{R}\Omega$. Then, \mathbf{G} is used as a proxy to approximate the left singular subspace of \mathbf{R} .

Unlike standard rSVD, SP-rSVD also computes a second projection $\mathbf{F} = \mathbf{R}^\top \mathbf{G} \in \mathbb{R}^{n \times p}$ in the same pass, enabling accurate SVD approximation without re-accessing \mathbf{R} . Once \mathbf{G} and \mathbf{F} are computed, the remaining SVD steps proceed efficiently on these compressed matrices. We refer the reader to the original paper by Yu et al. (2017) for implementation details and theoretical guarantees.

To compute \mathbf{G} and \mathbf{F} efficiently in a single pass over \mathbf{R} , we combine the SP-rSVD algorithm with the sequential summation strategy described in Equation 21. In particular, a single row of \mathbf{R} can be obtained by setting $\mathbf{X} = \mathbf{I}_{n:,i}$, yielding $\mathbf{R}_{i,:}^\top = \mathbf{R}^\top \mathbf{I}_{n:,i}$, where $\mathbf{I}_{n:,i}$ is the i th column on \mathbf{I}_n .

Using this approach, the matrices \mathbf{G} and \mathbf{F} can be constructed incrementally, as shown in Equation 22, starting with initializations $\mathbf{G} = \mathbf{0}_{n \times p}$ and $\mathbf{F} = \mathbf{0}_{n \times p}$:

$$\begin{aligned}
 &\text{For } i \in \{1, 2, \dots, n\} : \\
 &\quad \text{Compute } \mathbf{R}_{i,:} \text{ via Equation 21 with } \mathbf{X} = \mathbf{I}_{n:,i} \\
 &\quad \mathbf{g} = \mathbf{R}_{i,:} \Omega \\
 &\quad \mathbf{G}_{i,:} = \mathbf{g} \\
 &\quad \mathbf{F} = \mathbf{F} + \mathbf{R}_{i,:}^\top \mathbf{g}.
 \end{aligned} \tag{22}$$

In practice, it is advantageous to process multiple rows of \mathbf{R} in parallel batches to maximize GPU throughput. Letting b denote the batch size, the memory complexity of the algorithm becomes $\mathcal{O}(np + nb + m)$, where m is the number of non-zero entries in the adjacency matrix.

The total time complexity of proximity ReachNEs is $\mathcal{O}(nm + n^2p + np^2)$. The first term, nm , accounts for computing all n rows of \mathbf{R} using sparse matrix multiplications. The n^2p term arises from evaluating the remaining steps in Equation 22, and np^2 corresponds to the final SVD step in SP-rSVD after computing \mathbf{G} and \mathbf{F} (Yu et al., 2017).

To accelerate this process, our implementation leverages the PyTorch framework, with built-in support for both single- and multi-GPU environments. A further advantage over the original MATLAB implementation of SP-rSVD is PyTorch’s support for automatic differentiation, which opens up new research directions for self-supervised proximity embedding methods based on differentiable reachability matrices.

A.5 Computational complexity of DirSwitch and run time results

A.5.1 Computational complexity

We begin by analysing the time and space complexity of DirSwitch- r , followed by empirical run time measurements of our GPU-accelerated implementation.

The computational complexity of DirSwitch depends on the following variables:

- n : number of nodes.
- m : number of edges.
- p : embedding dimensionality.
- r : number of multi-directional steps.
- s : number of walk length distributions.
- K : number of walk steps (i.e., reachability terms).
- b : batch size used in SP-rSVD (proximity version only; see Section A.4).
- d : number of node features (message-passing version only).

We separately analyse the complexity of proximity and message-passing variants. We begin with the message-passing version, which is simpler due to the absence of batching.

Message-passing DirSwitch. Message-passing DirSwitch- r is implemented using the algorithm in Equation 6, which performs K sparse-dense matrix multiplications, each with time complexity $\mathcal{O}(md)$. This results in an overall cost of $\mathcal{O}(Kmd)$ per walk.

Each resulting $n \times d$ matrix is then reduced to $\frac{p}{s2^r}$ dimensions using PCA via SVD, with time complexity $\mathcal{O}(nd^2)$ (Golub & Van Loan, 2013, Ch. 5.5.6). This reduction is performed for all $s \cdot 2^r$ combinations of walk length distributions and directed neighbourhoods, resulting in a total time complexity of $\mathcal{O}(s2^r(Kmd + nd^2))$.

In terms of space complexity, we must store the m nonzero entries of the normalized adjacency matrix, the intermediate s reachability representations of shape $n \times d$, and the final embedding matrix $\mathbf{Z} \in \mathbb{R}^{n \times p}$. Our implementation computes the s multi-scale walks in parallel using tensors of shape $s \times n \times d$, while the 2^r directed neighbourhoods are processed sequentially. The resulting space complexity is $\mathcal{O}(m + snd + np)$.

Proximity DirSwitch. Proximity DirSwitch- r is implemented using the algorithm in Equation 22. Assuming a batch size b , the algorithm runs in $\frac{n}{b}$ iterations. Each iteration computes a slice of the reachability matrix via K sparse-dense matrix multiplications with cost $\mathcal{O}(Kmb)$, followed by projection onto $\mathbf{\Omega} \in \mathbb{R}^{n \times p}$ and the computation of \mathbf{F} , both requiring $\mathcal{O}(bnp)$. Repeating for all $\frac{n}{b}$ batches yields a total cost of $\mathcal{O}(Kmn + n^2p)$.

After constructing the sketch matrices, SP-rSVD computes the final embeddings in time $\mathcal{O}(np^2)$ (Yu et al., 2017), which is negligible when $p \ll n$. As with message-passing DirSwitch, this process is repeated for all $s \cdot 2^r$ combinations of walk length distributions and edge direction patterns, resulting in total time complexity $\mathcal{O}(s2^r(Kmn + n^2p))$.

In terms of space, the algorithm requires storage of the m edges, one $n \times b$ slice of the reachability matrix at a time, and the $n \times p$ matrices $\mathbf{\Omega}$, \mathbf{G} , and \mathbf{F} used by SP-rSVD. Our implementation stores s copies of $\mathbf{\Omega}$, one per walk length distribution, computed in parallel. The total space complexity is $\mathcal{O}(snp + snb + m)$.

Discussion. The use of SP-rSVD assumes $p \ll n$, and its numerical accuracy may degrade as p approaches n . For small graphs, we recommend computing the full reachability matrix directly and applying SVD without batching. This yields a time complexity of $\mathcal{O}(s2^r(Kmn + n^3))$ and space complexity of $\mathcal{O}(sn^2 + m)$.

While SP-rSVD helps improve the scalability of proximity DirSwitch- r by mitigating memory usage (see Section A.4), the method still incurs a quadratic time cost in n . In contrast, message-passing DirSwitch- r has linear time complexity in n , making it more scalable for large graphs.

A potential alternative to improve proximity scalability is a hybrid approach: apply SVD to the sparse adjacency matrix to generate initial proximity embeddings, then use these as node attributes in a message-passing DirSwitch- r model. Investigating the effectiveness and scalability of this hybrid method is a promising direction for future work.

A.5.2 Run time results

To evaluate practical run times, we conducted a small-scale scalability experiment. We generated synthetic undirected graphs using the Barabási–Albert model (Barabási & Albert, 1999), with an average degree of 2. We then measured the run time of DirSwitch while varying one parameter at a time— n , r , s , p , or d —while keeping the others fixed. For proximity DirSwitch, we fixed $n = 2^{16} = 65,536$, $r = 2$, $s = 3$, and $p = 256$. For message-passing DirSwitch, we used $n = 2^{20} = 1,048,576$, $r = 2$, $s = 3$, $p = 512$, and $d = 128$. Each experiment was repeated three times.

We tested proximity DirSwitch on three computational environments in Google Cloud: (i) 2 Nvidia L4 (24GB) GPUs, 24 vCPUs @ 2.20GHz, and 96GB RAM; (ii) 1 Nvidia L4 (24GB) GPU, 16 vCPUs @ 2.20GHz, and 64GB RAM; and (iii) CPU-only with 16 vCPUs @ 2.20GHz and 64GB RAM. These are denoted as GPU-2, GPU-1, and CPU-0, respectively. As our current implementation of message-passing DirSwitch does not support multi-GPU computation, we only evaluated it in the GPU-1 and CPU-0 environments.

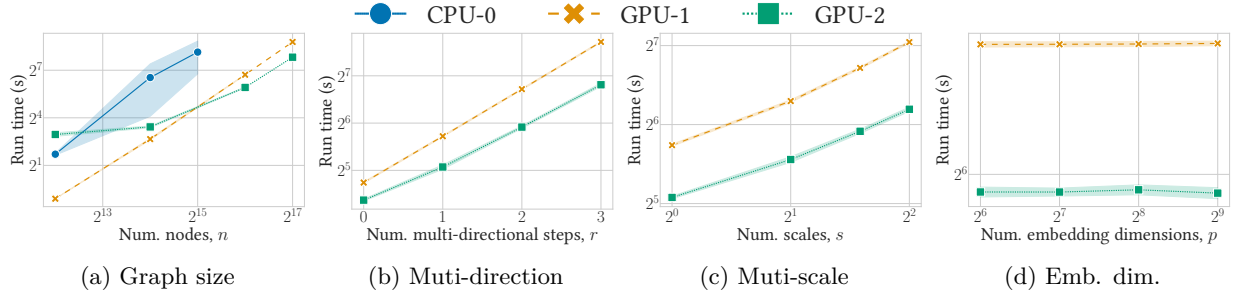


Figure 8: Scalability results for proximity DirSwitch- r . Each subfigure plots run time (in seconds) against a key parameter in the DirSwitch time complexity. Markers indicate the mean over three runs, with shaded areas showing standard deviation. Line colour and style correspond to the number of GPUs used: 0, 1, or 2.

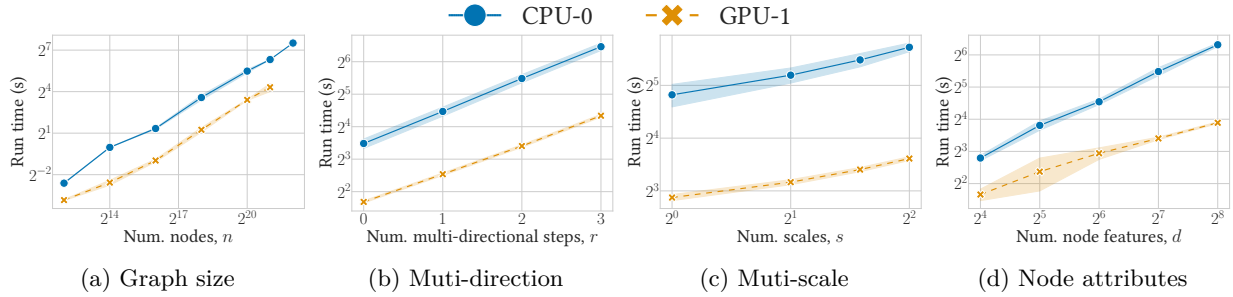


Figure 9: Scalability results for message-passing DirSwitch- r . Each subfigure plots run time (in seconds) against a key parameter in the DirSwitch time complexity. Markers show the mean over three runs; shaded areas indicate standard deviation. Line colour and style correspond to the number of GPUs used: 0 or 1.

Figures 8 and 9 show the average run times for proximity and message-passing DirSwitch, respectively, with line colours indicating the compute environment. These results clearly demonstrate the advantage of GPU acceleration, as both models run significantly faster on GPU than on CPU. For proximity DirSwitch, CPU runtimes were prohibitively slow at $n = 2^{16}$, and we therefore aborted those runs—hence the missing data points. Although multi-GPU setups introduce some overhead, they become advantageous for larger graphs.

Most trends align with theoretical expectations. Figure 8a confirms the quadratic time complexity in n for proximity DirSwitch, while Figure 9a shows linear scaling for the message-passing variant. The exponential growth in the number of directed neighbourhoods, 2^r , is reflected in the roughly exponential growth in run times shown in Figures 8b and 9b. In Figures 8c and 9c, we observe sublinear growth in run time with respect to s , due to parallelization across the multi-scale computations.

One unexpected result appears in Figure 8d: the run time of proximity DirSwitch remains largely unaffected by the embedding dimensionality p . This is counterintuitive, as the theoretical time complexity includes an n^2p term. However, this term is dominated in practice by the Kmn cost of computing reachability. A similar pattern is observed in Figure 9d, where run time increases linearly with d for message-passing DirSwitch—indicating that the Kmd term outweighs the nd^2 term. These findings suggest that optimizing the reachability computation is the most promising avenue for improving the overall efficiency of DirSwitch.

Tables 10 and 11 report measured run times for the experiments in Section 6.3. These were run on varying computational environments, so runtimes may differ slightly across datasets. We only report results for the slowest datasets. On the others, runtime typically falls between 1–3 seconds.

The results confirm the trends observed in the scalability experiments. Since the implementation parallelizes over the number of scales (s) while computing different direction sequences (r) sequentially, r has a significantly greater impact on runtime than s .

The Cora datasets yield the longest runtimes, primarily because they contain the largest number of node features ($d = 8710$; see Table 2).

Table 10: Embedding computation duration for message-passing ReachNEs with $p = 1024$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers.

EDGE DIRECTIONS	CORA-ML					CORA				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	2s	2s	4s	6s	8s	1m 24s	1m 25s	2m 49s	5m 33s	6m 2s
UNDIRECTED	2s	4s	4s	7s	8s	1m 23s	1m 24s	2m 49s	4m 26s	5m 40s
MULTIDIR-1	3s	4s	8s	13s	17s	2m 47s	2m 49s	5m 36s	9m 43s	11m 36s
MULTIDIR-2	7s	7s	16s	27s	35s	5m 33s	5m 37s	11m 10s	19m 13s	22m 56s
MULTIDIR-3	19s	19s	43s	1m 9s	1m 30s	10m 54s	11m 2s	21m 57s	37m 11s	44m 48s
DIRSWITCH-1	3s	3s	8s	14s	18s	2m 47s	2m 47s	5m 35s	8m 52s	11m 18s
DIRSWITCH-2	7s	7s	17s	29s	38s	5m 35s	5m 37s	11m 10s	17m 58s	22m 40s
DIRSWITCH-3	14s	15s	33s	59s	1m 15s	11m 6s	11m 13s	22m 16s	36m 14s	45m 21s

EDGE DIRECTIONS	POKEC					SNAP-PATENTS				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	18s	22s	23s	21s	21s	26s	26s	45s	57s	56s
UNDIRECTED	22s	18s	19s	21s	21s	31s	37s	42s	53s	1m 39s
MULTIDIR-1	44s	44s	47s	57s	1m 3s	59s	1m 24s	2m 8s	2m 45s	2m 2s
MULTIDIR-2	1m 48s	1m 30s	1m 54s	1m 41s	1m 46s	1m 58s	1m 59s	4m 57s	6m 27s	7m 34s
MULTIDIR-3	4m 10s	4m 11s	4m 21s	4m 37s	4m 48s	4m 56s	4m 56s	6m 39s	8m 23s	9m 52s
DIRSWITCH-1	34s	43s	46s	48s	42s	1m 54s	1m 1s	1m 23s	2m 53s	3m 34s
DIRSWITCH-2	1m 26s	1m 22s	1m 28s	1m 17s	1m 22s	2m 4s	2m 2s	5m 7s	3m 59s	7m 53s
DIRSWITCH-3	2m 54s	2m 54s	3m 0s	3m 8s	3m 26s	6m 32s	7m 18s	5m 29s	11m 40s	8m 15s

Table 11: Embedding computation duration for proximity ReachNEs with $p = 1024$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers.

EDGE DIRECTIONS	CoCITE					PUBMED					CORA (SUBELJ)				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	3s	3s	4s	5s	6s	1s	1s	1s	1s	1s	1s	1s	1s	1s	2s
UNDIRECTED	4s	5s	6s	7s	8s	1s	1s	2s	3s	3s	1s	1s	1s	2s	2s
MULTIDIR-1	6s	6s	8s	9s	11s	1s	1s	2s	2s	3s	1s	1s	2s	2s	3s
MULTIDIR-2	12s	12s	15s	19s	22s	3s	3s	4s	4s	5s	3s	3s	4s	5s	6s
MULTIDIR-3	37s	37s	57s	1m 17s	1m 37s	7s	7s	10s	14s	18s	10s	10s	15s	21s	26s
DIRSWITCH-1	8s	8s	10s	13s	15s	2s	2s	3s	4s	4s	2s	2s	2s	3s	3s
DIRSWITCH-2	14s	15s	19s	25s	28s	3s	3s	4s	6s	7s	3s	3s	4s	6s	7s
DIRSWITCH-3	28s	28s	36s	47s	53s	6s	6s	8s	11s	13s	6s	6s	9s	11s	13s

Finally, Table 12 shows the run times of the models used in the node classification benchmark in Section 6.4. DirSwitch ranks among the fastest methods, alongside BLADE and NERD. Note that graph size and hyperparameter settings vary across these results.

A.6 Guide to hyperparameter selection for DirSwitch

DirSwitch- r embeddings are governed by four key hyperparameters: the embedding dimensionality p , the multi-directional neighbourhood radius r , the number of embedding scales s , and the choice of walk length distributions $P_w(k; \tau)$. Each of these influences both embedding quality and computational cost. This section offers practical guidelines for selecting these hyperparameters, based on our experimental findings and experience with DirSwitch.

Walk length distribution. We begin with the walk length distribution $P_w(k - l_o; \tau)$, where l_o denotes an offset. For a single-scale setup ($s = 1$), we recommend using $Pois(k - l_o; \tau = 2)$, with $l_o = 0$ for message-passing and $l_o = 1$ for proximity embeddings, as explained in Section 2.3.1. This setting emphasizes the $k = 1$ and $k = 2$ neighbourhoods, which typically offer the most structural diversity, while still incorporating

Table 12: Embedding computation duration for proximity node embedding models. Average and standard deviations are computed using 5 random seeds.

MODEL	FLY LARVA	EU-EMAIL	POLBLOGS	CoCITE	PUBMED	CORA (SUBELJ)
HOPE ¹	58s	9s	7s	41M 47s	10M 46s	4M 45s
APP ²	15s	3s	1s	8M 41s	1M 54s	2M 14s
NERD ³	9s	28s	4s	12s	18s	11s
DGGAN ⁴	29M 10s	6M 19s	2M 24s	1H 33M	7M 25s	51M 47s
BLADE ⁵	5s	1s	1s	23s	11s	21s
DIRSWITCH	7s	677MS	287MS	2M 31s	36s	34s

¹OU ET AL. (2016)²ZHOU ET AL. (2017)³KHOSLA ET AL. (2020)⁴ZHU ET AL. (2021B)⁵VIRINCHI & SALADI (2023)

information from $k = 0, 3$, and 4 (see Figure 11a). Moreover, our graph alignment results suggest that $\tau \approx 2$ tends to correspond to peak expressivity (see Figures 7 and 15). Lower values of τ risk overly localized embeddings, while higher values may lead to over-smoothing.

When using multiple distributions ($s > 1$), it is important to avoid excessive overlap between them, as this leads to redundant embedding dimensions. We find that Geometric distributions with small τ and varying offsets are particularly effective, as they place their modes at distinct distances from the source node. This configuration performed especially well on the Pokec dataset (Section G).

As a general rule for $s \geq 3$, we recommend the following family of walk distributions:

$$\{\text{Geom}(k - s_i; \tau = 1 + 0.5s_i) \mid s_i \in \{0, 1, \dots, s - 1\}\}. \quad (23)$$

For $s = 2$, it is preferable to combine one short-range and one long-range distribution, such as $\text{Geom}(k; \tau = 2)$ and $\text{Pois}(k - 2; \tau = 3)$.

Expressivity and computation trade-offs. We now discuss the parameters p , r , and s , which govern the trade-off between model expressivity and computational cost. Increasing any of these parameters typically enhances embedding expressivity and improves downstream task performance, as demonstrated in Section 6. However, this comes at the cost of increased computation, as analysed in Section A.5.

Notably, increasing r and s leads to diminishing returns for two key reasons. First, larger values result in embeddings that capture increasingly broader neighbourhoods. Since most real-world graphs exhibit the small-world property (Watts & Strogatz, 1998)—i.e., short diameters and small average path lengths (see Table 2)—larger neighbourhoods quickly converge to covering the full graph, reducing the embeddings’ ability to distinguish local structure.

Our graph alignment experiments support this: using the Poisson distribution, the most discriminative and noise-robust embeddings arise from walks of 1 to 5 steps (see Figures 7 and 15). We therefore recommend keeping r and the support of the walk length distribution within 5–6 steps, unless the graph’s structure specifically warrants longer ranges. Nonetheless, further research is needed to quantify these diminishing returns and understand how they vary across graph types and downstream tasks.

A second source of diminishing returns arises from the constraint of fixed embedding dimensionality p . Since DirSwitch- r concatenates embeddings across s scales and 2^r directed neighbourhoods, each (σ, P_w) pair contributes only $\frac{p}{s2^r}$ dimensions. If s or r are set too high relative to p , each component becomes overly compressed, degrading the overall embedding quality.

Increasing p generally improves expressivity, but it also raises computational demands. In particular, excessive values may exceed GPU memory limits, forcing computations to fall back to CPUs—significantly increasing run time, as shown in Section A.5.

Based on our experimental findings with $p = 512$ in Section G, where we observed diminishing returns, we recommend ensuring $p \geq 32s2^r$ as a practical guideline. For instance, with $p = 1024$, we suggest $r = 3$ and $s = 3$ or $s = 4$, which matches the best-performing models in Section 6.3. If hardware permits, p , r , and s can be increased proportionally while respecting the $p \geq 32s2^r$ rule. For resource-constrained settings, we recommend $r = 1$, $s = 2$, and $p = 256$ as an efficient configuration that still benefits from both multi-directionality and multi-scale expressivity.

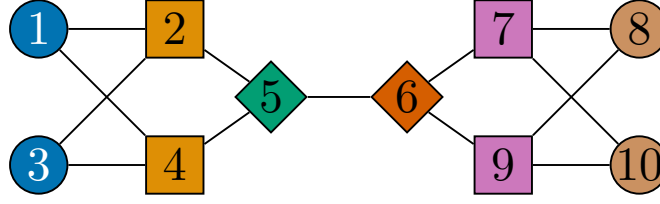


Figure 10: Illustration of structural and automorphic node equivalence. Nodes with the same colour (e.g., nodes 1 and 3) are structurally equivalent, meaning they share identical connections to the same set of nodes. Nodes with the same shape (e.g., nodes 2, 4, 7, and 9) are automorphically equivalent, meaning they are indistinguishable based solely on graph structure without node labels.

B Message-passing embeddings and automorphic node equivalence

This section provides additional detail on automorphic node equivalence and proves that message-passing ReachNEs embeddings assign identical representations to automorphic nodes.

Figure 10 illustrates both structural and automorphic node equivalence. Nodes with the same colour are structurally equivalent, while nodes with the same shape are automorphically equivalent. As shown, structural equivalence always implies automorphic equivalence, but not vice versa. Figure 10 also highlights how structural equivalence implies node proximity, whereas automorphically equivalent nodes can be far apart, such as nodes 1 and 10. In fact, automorphic nodes can even exist in separate weakly connected components.

To prove that message-passing ReachNEs embeddings preserve automorphic equivalence, we begin by formalizing the definition. Let $\pi : \mathbb{N} \rightarrow \mathbb{N}$ be a permutation of the node set in a graph $\mathcal{G} = (\mathbb{N}, \mathbb{M})$, and let $\mathbf{P}_\pi \in \{0, 1\}^{n \times n}$ denote the corresponding permutation matrix, where $\mathbf{P}_\pi(\pi(i), i) = 1$ for all i , and all other entries are zero. An *automorphism* is any such permutation π that preserves the edge structure of the graph, i.e., $(i, j) \in \mathbb{M} \iff (\pi(i), \pi(j)) \in \mathbb{M}$. Equivalently, this condition holds if the adjacency matrix satisfies $\mathbf{A} = \mathbf{P}_\pi \mathbf{A} \mathbf{P}_\pi^\top$. Two nodes i and j are *automorphically equivalent* if there exists an automorphism π such that $\pi(i) = j$, i.e., $\mathbf{P}_\pi(j, i) = 1$.

Next, we prove that message-passing ReachNEs produces identical embeddings for automorphic nodes, i.e., $\mathbf{P}_\pi \mathbf{Z} = \mathbf{Z}$, under the *node attribute equivalence* assumption, $\mathbf{X} = \mathbf{P}_\pi \mathbf{X}$, which is guaranteed if \mathbf{X} consists of graph-derived features.

We start by showing $\mathbf{R} = \mathbf{P}_\pi \mathbf{R} \mathbf{P}_\pi^\top$. To see this, note that by the definition of automorphism, we have $\mathbf{A} = \mathbf{P}_\pi \mathbf{A} \mathbf{P}_\pi^\top$. The same holds for the normalized adjacency matrices, $\mathbf{A}_\sigma = \mathbf{P}_\pi \mathbf{A}_\sigma \mathbf{P}_\pi^\top$, and their powers, $\mathbf{A}_\sigma^k = \mathbf{P}_\pi \mathbf{A}_\sigma^k \mathbf{P}_\pi^\top$. This follows from:

$$\prod_{l=1}^k \mathbf{A}_{\sigma_l} = \prod_{l=1}^k \mathbf{P}_\pi \mathbf{A}_{\sigma_l} \mathbf{P}_\pi^\top = \mathbf{P}_\pi \mathbf{A}_{\sigma_k} \dots \mathbf{A}_{\sigma_2} \underbrace{\mathbf{P}_\pi^\top \mathbf{P}_\pi}_{\mathbf{I}_n} \mathbf{A}_{\sigma_1} \mathbf{P}_\pi^\top = \mathbf{P}_\pi \prod_{l=1}^k \mathbf{A}_{\sigma_l} \mathbf{P}_\pi^\top.$$

Since the reachability matrix \mathbf{R} is a sum of power matrices, see Equation 8, it follows that \mathbf{R} is also invariant under automorphism $\mathbf{R} = \mathbf{P}_\pi \mathbf{R} \mathbf{P}_\pi^\top$.

To complete the proof, recall that $\mathbf{Z} = \mathbf{R}^\top \mathbf{X}$ represents message-passing embeddings (see Section 2.3.2). We then simply insert $\mathbf{R} = \mathbf{P}_\pi \mathbf{R} \mathbf{P}_\pi^\top$ and $\mathbf{X} = \mathbf{P}_\pi \mathbf{X}$:

$$\mathbf{P}_\pi \mathbf{Z} = \mathbf{P}_\pi \mathbf{R}^\top \mathbf{X} = \mathbf{P}_\pi \mathbf{R}^\top \underbrace{\mathbf{I}_n}_{\mathbf{P}_\pi^\top \mathbf{P}_\pi} \mathbf{X} = \underbrace{\mathbf{P}_\pi \mathbf{R}^\top \mathbf{P}_\pi^\top}_{\mathbf{R}^\top} \underbrace{\mathbf{P}_\pi \mathbf{X}}_{\mathbf{X}} = \mathbf{R}^\top \mathbf{X} = \mathbf{Z}.$$

Unlike message-passing ReachNEs, proximity embeddings are obtained via SVD reduction, a global operation that does not necessarily preserve automorphism invariance. Thus, it is generally the case for proximity embeddings that $\mathbf{P}_\pi \mathbf{Z} \neq \mathbf{Z}$. A formal proof for undirected graphs is provided by Zhu et al. (2021a).

C Asymptotic behaviour of random walk length distributions

This section expands on the analysis in Section 2.2, focusing on the asymptotic behaviour of different random walk length distributions and their implications for the reachability matrix \mathbf{R} as $\tau \rightarrow \infty$.

We begin with the Poisson distribution, deriving its associated diffusion differential equation and establishing its well-known connection to the normalized Laplacian. Next, we derive a similar equation for the geometric distribution, demonstrating that, unlike the Poisson case, it does not necessarily converge to a non-informative stationary distribution as $\tau \rightarrow \infty$.

We then examine the relationships between distributions in the K -truncated setting, proving that:

$$\begin{aligned}\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Geom}(\tau)) &\rightarrow \mathbf{R}^{(K)}(\mathbf{A}_*; \mathcal{U}(\tau)), \quad \text{as } \tau \rightarrow \infty, \\ \mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\tau)) &\rightarrow \mathbf{R}^{(K)}(\mathbf{A}_*; \text{Pois}(\tau)), \quad \text{as } K \rightarrow \infty.\end{aligned}$$

Finally, we highlight the relation between the binomial distribution and random walks on graphs with self-loops added at each node.

C.1 The Poisson Distribution

We now establish the identity:

$$\mathbf{R}(\mathbf{A}_*; \text{Pois}(\tau)) = e^{-\tau \mathbf{L}_*}, \quad (24)$$

where $\mathbf{L}_* = \mathbf{I}_n - \mathbf{A}_*$ is the normalized graph Laplacian.

Starting from the definition in Equation 2, we rewrite the reachability matrix as:

$$\mathbf{R}(\mathbf{A}_*; \text{Pois}(\tau)) = \sum_{k=0}^{\infty} e^{-\tau} \frac{\tau^k}{k!} \mathbf{A}_*^k = e^{-\tau} e^{\tau \mathbf{A}_*} = e^{-\tau \mathbf{L}_*},$$

using properties of the matrix exponential (Hall, 2015).

Next, we show that for an initial probability distribution $\mathbf{p}^{(0)} \in [0, 1]^n$, the solution:

$$\mathbf{p}(\tau) = \mathbf{R}(\mathbf{A}_*, \text{Pois}(\tau)) \mathbf{p}^{(0)} \quad (25)$$

satisfies the differential equation:

$$\frac{d\mathbf{p}}{d\tau} = -\mathbf{L}_* \mathbf{p}(\tau). \quad (26)$$

Proof. From Equation 24, we recall that:

$$\mathbf{R}(\mathbf{A}_*, \text{Pois}(\tau)) = e^{-\tau \mathbf{L}_*}.$$

Setting $\tau = 0$ gives $\mathbf{R}(\mathbf{A}_*, \text{Pois}(0)) = \mathbf{I}_n$, ensuring that the initial condition is satisfied:

$$\mathbf{p}(0) = \mathbf{p}^{(0)}.$$

Differentiating Equation 25 with respect to τ :

$$\frac{d\mathbf{p}(\tau)}{d\tau} = \frac{d}{d\tau} e^{-\tau \mathbf{L}_*} \mathbf{p}^{(0)} = -\mathbf{L}_* e^{-\tau \mathbf{L}_*} \mathbf{p}^{(0)} = -\mathbf{L}_* \mathbf{p}(\tau).$$

Thus, $\mathbf{p}(\tau)$ satisfies the differential equation as required. \square

The differential relation in Equation 26 is useful for analysing the dynamics of $\mathbf{R}(\mathbf{A}_*, \text{Pois}(\tau))$ with respect to τ , allowing us to establish its asymptotic behavior.

Focusing first on the undirected case, the eigenvalue multiplicity of 0 in the normalized Laplacian \mathbf{L}_U equals the number of weakly connected components in the graph (Chung, 1997; von Luxburg, 2007). The corresponding eigenvectors \mathbf{u} take values:

$$u_i = \frac{\deg(i)}{\sum_{k \in \mathcal{C}(j)} \deg(k)}$$

for each node i within a weakly connected component $\mathcal{C}(j)$. These eigenvectors are the only stationary solutions to Equation 26, as all other eigenvalues of \mathbf{L}_U are strictly positive.

Thus, as $\tau \rightarrow \infty$, the reachability matrix converges to:

$$R(\mathbf{A}_U, \text{Pois}(\tau))_{i,j} \rightarrow \frac{\deg(i)}{\sum_{k \in \mathcal{C}(j)} \deg(k)}. \quad (27)$$

Importantly, all nodes within the same weakly connected component have identical asymptotic reachability values, making them indistinguishable under the Poisson walk length distribution for undirected graphs. Consequently, downstream tasks relying on embedding distinguishability, such as graph alignment, are expected to degrade as τ increases. This effect was also observed in Figures 7d and 7b.

The above analysis provides a good first-order approximation of the Poisson distribution's behaviour on digraphs. However, a precise theoretical analysis using \mathbf{L}_0 and \mathbf{L}_1 is more involved and beyond the scope of this paper. Instead, we refer to Veerman & Lyons (2020) for valuable insights. Specifically, the multiplicity of the eigenvalue 0 corresponds to the number of *reaches* in the graph (Veerman & Lyons, 2020, Theorem 4.6). Unlike weakly connected components, multiple reaches can overlap. Consequently, the eigenvectors exhibit a more complex structure (Veerman & Lyons, 2020, Theorem 5.1), and nodes within a reach do not necessarily converge to the same reachability vector as $\tau \rightarrow \infty$.

C.2 The geometric distribution

We now analyse the asymptotic behavior of walk lengths following the geometric distribution. We start with the parameterization $\alpha \in [0, 1)$, related to τ by $\alpha = \frac{\tau}{1+\tau}$, so that as $\tau \rightarrow \infty$, we have $\alpha \rightarrow 1$.

As before, we express the reachability matrix as a matrix function:

$$\mathbf{R}(\mathbf{A}_*; \text{Geom}(\alpha)) = (1 - \alpha) \sum_{k=0}^{\infty} \alpha^k \mathbf{A}_*^k = (1 - \alpha) (\mathbf{I}_n - \alpha \mathbf{A}_*)^{-1}. \quad (28)$$

Substituting $\alpha = \frac{\tau}{1+\tau}$ and simplifying, we obtain:

$$\begin{aligned} \mathbf{R}(\mathbf{A}_*; \text{Geom}(\tau)) &= \frac{1}{1+\tau} \left(\mathbf{I}_n - \frac{\tau}{1+\tau} \mathbf{A}_* \right)^{-1} \\ &= (\mathbf{I}_n + \tau(\mathbf{I}_n - \mathbf{A}_*))^{-1} \\ &= (\mathbf{I}_n + \tau \mathbf{L}_*)^{-1}. \end{aligned}$$

Using this identity, we show that $\mathbf{p}(\tau) = \mathbf{R}(\mathbf{A}_*; \text{Geom}(\tau)) \mathbf{p}^{(0)}$ satisfies the differential equation:

$$\frac{d\mathbf{p}(\tau)}{d\tau} = -(\mathbf{I}_n + \tau \mathbf{L}_*)^{-1} \mathbf{L}_* \mathbf{p}(\tau), \quad \mathbf{p}(0) = \mathbf{p}^{(0)}. \quad (29)$$

Proof. The initial condition is satisfied since $\mathbf{R}(\mathbf{A}_*; \text{Geom}(0)) = \mathbf{I}_n$, so $\mathbf{p}(0) = \mathbf{p}^{(0)}$. To compute the derivative $\frac{d\mathbf{p}(\tau)}{d\tau}$, we apply the matrix inversion derivative rule (Petersen & Pedersen, 2012, Sec. 2.2):

$$\begin{aligned} \frac{d\mathbf{p}(\tau)}{d\tau} &= \frac{d}{d\tau} \left((\mathbf{I}_n + \tau \mathbf{L}_*)^{-1} \right) \mathbf{p}^{(0)} \\ &= -(\mathbf{I}_n + \tau \mathbf{L}_*)^{-1} \mathbf{L}_* (\mathbf{I}_n + \tau \mathbf{L}_*)^{-1} \mathbf{p}^{(0)} \\ &= -(\mathbf{I}_n + \tau \mathbf{L}_*)^{-1} \mathbf{L}_* \mathbf{p}(\tau). \end{aligned}$$

Thus, the differential equation is verified. \square

The differential equation in Equation 29 offers insight into the similarities and differences between the dynamics of the geometric and Poisson distributions. Both equations contain the term $\mathbf{L}_* \mathbf{p}(\tau)$, implying that any stationary distribution under Poisson dynamics is also stationary for the geometric case. However, the presence of the prefactor $(\mathbf{I}_n + \tau \mathbf{L}_*)^{-1}$ in Equation 29 suggests that this stationary distribution may not be reached under geometric dynamics.

To understand this, let λ be an eigenvalue of \mathbf{L}_* . Then $\frac{1}{1+\tau\lambda}$ is an eigenvalue of $(\mathbf{I}_n + \tau \mathbf{L}_*)^{-1}$. For any $\lambda > 0$, we see that $\frac{1}{1+\tau\lambda} \rightarrow 0$ as $\tau \rightarrow \infty$. Consequently, many eigenvectors will yield near-zero gradients as τ increases, effectively stalling the dynamics.

As a result, $\mathbf{p}(\tau)$ is likely to become “stuck” before converging to the non-informative stationary distribution that characterizes Poisson dynamics. This phenomenon is also observed empirically in Figures 5 and 7, where both dispersal entropy and graph alignment accuracy exhibit this same saturation effect.

C.3 Relation between the geometric and the uniform distribution

An interesting asymptotic behavior of the geometric distribution is that its K -truncation approaches the uniform distribution as $\tau \rightarrow \infty$, or equivalently, as $\alpha \rightarrow 1$.

To demonstrate this, we begin with the expression for the K -truncated reachability of the geometric distribution using the α parameterization:

$$\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Geom}(\alpha)) = \frac{(1-\alpha)}{Z(\alpha, K)} \sum_{k=0}^{\infty} \alpha^k \mathbf{A}_*^k, \quad (30)$$

where $Z(\alpha, K)$ is the normalization factor required due to the truncation.

Next, we reparameterize using $\varepsilon = 1 - \alpha$, yielding:

$$\mathbf{R}^{(K)}(\mathbf{A}_*; \varepsilon) = \frac{\varepsilon}{Z(\varepsilon, K)} \sum_{k=0}^K (1-\varepsilon)^k \mathbf{A}_*^k. \quad (31)$$

We can compute the normalization factor $Z(\varepsilon, K)$ as:

$$Z(\varepsilon, K) = \sum_{k=0}^K \varepsilon (1-\varepsilon)^k = \varepsilon \frac{1 - (1-\varepsilon)^{K+1}}{1 - (1-\varepsilon)} = 1 - (1-\varepsilon)^{K+1}.$$

As $\varepsilon \rightarrow 0$, we can use the approximation $(1-\varepsilon)^k = 1 - k\varepsilon + \mathcal{O}(\varepsilon^2)$ for both $Z(\varepsilon, K)$ and the summands in Equation 31 to show that $\mathbf{R}^{(K)}(\mathbf{A}_*; \varepsilon)$ approaches the reachability matrix under the uniform distribution:

$$\begin{aligned} \lim_{\varepsilon \rightarrow 0} \mathbf{R}^{(K)}(\mathbf{A}_*; \varepsilon) &= \lim_{\varepsilon \rightarrow 0} \frac{\varepsilon}{(K+1)\varepsilon} \sum_{k=0}^K (1 - k\varepsilon) \mathbf{A}_*^k \\ &= \lim_{\varepsilon \rightarrow 0} \frac{1}{K+1} \sum_{k=0}^K \mathbf{A}_*^k - \frac{\varepsilon}{K+1} \sum_{k=0}^K k \mathbf{A}_*^k \\ &= \frac{1}{K+1} \sum_{k=0}^K \mathbf{A}_*^k. \end{aligned}$$

C.4 Relation between the binomial and the Poisson distribution

The final asymptotic behavior we highlight is that the reachability under the binomial walk length distribution, $\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\tau))$, approaches that of the Poisson distribution $\mathbf{R}(\mathbf{A}_*; \text{Pois}(\tau))$ as $K \rightarrow \infty$.

We begin by expressing $\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\tau))$ as a matrix function using the Binomial Theorem:

$$\begin{aligned}\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\tau)) &= \sum_{k=0}^K \binom{K}{k} \left(1 - \frac{\tau}{K}\right)^{K-k} \left(\frac{\tau}{K}\right)^k \mathbf{A}_*^k \\ &= \left(\left(1 - \frac{\tau}{K}\right) \mathbf{I}_n + \frac{\tau}{K} \mathbf{A}_*\right)^K \\ &= \left(\mathbf{I}_n - \frac{\tau}{K} \mathbf{L}_*\right)^K.\end{aligned}$$

Next, we perform the substitution $K_\tau = \frac{K}{\tau}$, and as we take the limit $K_\tau \rightarrow \infty$, we obtain:

$$\lim_{K_\tau \rightarrow \infty} \left(\mathbf{I}_n - \frac{1}{K_\tau} \mathbf{L}_*\right)^{\tau K_\tau} = e^{-\tau \mathbf{L}_*} = \mathbf{R}(\mathbf{A}_*; \text{Pois}(\tau)). \quad (32)$$

Here, we recover the matrix exponential using its limit definition (Hall, 2015), which, together with Equation 24, proves the asymptotic equality to the Poisson-based reachability.

C.5 Relation between the binomial distribution and self-loops

Graph Convolutional Networks (GCNs) (Kipf & Welling, 2017) popularized the common pre-processing step of adding self-loops to each node, a practice that has since been widely adopted in self-supervised GNN methods (Wu et al., 2019; Zhang et al., 2021a; Thakoor et al., 2022). As discussed in Section 2.2, there is a close connection between the binomial walk length distribution and graphs with added self-loops, which we explore in this section.

This connection becomes evident when we consider the random surfer interpretation of the binomial distribution. In this view, a random walker makes a total of K decisions. At each step, the walker moves to a neighbouring node with probability α , or remains at its current node (i.e., performs a self-loop) with probability $1 - \alpha$. The probability of taking exactly k steps to neighbouring nodes (and $K - k$ self-loops) then follows the binomial distribution:

$$P_{\mathbf{w}}(k; \alpha) = \binom{K}{k} (1 - \alpha)^{K-k} \alpha^k.$$

The corresponding reachability matrix is given by:

$$\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\alpha)) = \sum_{k=0}^K \binom{K}{k} (1 - \alpha)^{K-k} \alpha^k \mathbf{A}_*^k. \quad (33)$$

We can rewrite this expression using the Binomial Theorem and the fact that the diagonal matrix $(1 - \alpha)\mathbf{I}_n$ commutes with any matrix under multiplication:

$$\begin{aligned}\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\alpha)) &= \sum_{k=0}^K \binom{K}{k} (1 - \alpha)^{K-k} \alpha^k \mathbf{A}_*^k \\ &= \sum_{k=0}^K \binom{K}{k} (1 - \alpha)^{K-k} \mathbf{I}_n^{K-k} \alpha^k \mathbf{A}_*^k \\ &= \sum_{k=0}^K \binom{K}{k} ((1 - \alpha)\mathbf{I}_n)^{K-k} (\alpha \mathbf{A}_*)^k \\ &= ((1 - \alpha)\mathbf{I}_n + \alpha \mathbf{A}_*)^K \\ &= \mathbf{A}_{*+}^K.\end{aligned}$$

Here, we define $\mathbf{A}_{*+} = (1 - \alpha)\mathbf{I}_n + \alpha\mathbf{A}_*$ as a self-loop enhanced normalized adjacency matrix. We see that the binomial reachability matrix is simply the K th power of \mathbf{A}_{*+} .

This matrix consists of two parts: a diagonal component $(1 - \alpha)\mathbf{I}_n$, representing the probability of staying at the current node (i.e., self-loops), and a weighted adjacency component $\alpha\mathbf{A}_*$, representing the transition probabilities to neighbouring nodes.

A very similar construction arises in message-passing models that explicitly add self-loops. In what follows, we formalize this for random walks using the outgoing edge transition matrix \mathbf{A}_0 . The analysis extends directly to the incoming edge matrix \mathbf{A}_1 and the undirected matrix \mathbf{A}_u .

First, let $\mathbf{S} = \mathbf{A} + \mathbf{I}_n$ be the adjacency matrix with added self-loops, and let $\mathbf{D}_{0+} = \mathbf{D}_0 + \mathbf{I}_n$ denote the corresponding out-degree matrix. Then, the normalized transition matrix for the self-loop-enhanced graph is given by $\mathbf{S}_0 = \mathbf{S}\mathbf{D}_{0+}^{-1}$.

We decompose \mathbf{S}_0 into two components: $\mathbf{S}_0 = \mathbf{P}_\circ + \mathbf{P}_\rightarrow$, where $\mathbf{P}_\circ = \text{diag}(\mathbf{S}_0)$ contains the self-loop (stay-in-place) probabilities, and \mathbf{P}_\rightarrow contains the off-diagonal transition probabilities to neighbouring nodes.

This structure closely mirrors that of $\mathbf{A}_{*+} = (1 - \alpha)\mathbf{I}_n + \alpha\mathbf{A}_*$. In both cases, the stay-in-place behavior is captured by a diagonal matrix— \mathbf{P}_\circ for \mathbf{S}_0 and $(1 - \alpha)\mathbf{I}_n$ for \mathbf{A}_{*+} —while the remaining transitions are modelled via \mathbf{P}_\rightarrow and $\alpha\mathbf{A}_*$, respectively.

Consequently, the K th power of \mathbf{S}_0 produces a matrix polynomial with binomial coefficients, analogous to the binomial reachability matrix in Equation 33. This observation aligns with practice: in a K -step linear GCN, node embeddings are computed as $\mathbf{Z} = (\mathbf{S}_0^K)^\top \mathbf{X}$ (Wu et al., 2019), which matches the message-passing form of the binomial reachability matrix $\mathbf{R}^{(K)}(\mathbf{A}_*; \text{Binom}(\alpha)) = \mathbf{A}_{*+}^K$ described in Section 2.3.2.

Thus, both approaches share the same underlying random surfer interpretation and exhibit similar mathematical structures. However, there is one key difference: in the binomial reachability model, the stay-in-place probability is uniform across all nodes, controlled by the scalar $(1 - \alpha)$. In contrast, for GCNs with added self-loops, the stay-in-place probability varies by node and is given by \mathbf{P}_\circ , where $P_{\circ,i,i} = \frac{1}{1 + \deg_0(i)}$. This means that nodes with lower out-degree have higher probability of remaining in place.

Whether this degree-dependent behavior improves embedding quality is ultimately an empirical question and may depend on the characteristics of the graph and the downstream task.

D Additional experiment setup information

D.1 Single- and multi-scale ReachNEs embedding distributions

This section provides further details related to the experimental setup described in Section 6. Figure 11 shows the random walk length distributions used in the ReachNEs experiments, as listed in Table 3.

The first plot, Figure 11a, displays the two single-scale distributions, referred to as **Geom** and **Pois**. The Geom distribution corresponds to $\text{Geom}(k; \tau = 1)$, emphasizing the node itself ($k = 0$) and short-range walks (1–2 steps). In contrast, Poiss uses $\text{Pois}(k; \tau = 2)$, which down-weights self-loops and focuses on medium-range neighbourhoods (2–4 steps).

The **Geom-U** setting creates two-scale embeddings by combining two distributions designed to capture both short- and medium-range structure. The first is $\text{Geom}(k; \tau = 1)$ (same as in Geom), while the second is the shifted uniform distribution $\mathcal{U}(k - 1; \tau = 2)$, which excludes the $k = 0$ term and concentrates probability mass over $k \in [1, 5]$.

The **Binom-3** and **Geom-4** settings are multi-scale distributions designed to capture short-, medium-, and long-range reachability. Binom-3 is inspired by Gaussian mixture models (Murphy, 2012, Ch. 11.2.1) and combines three binomial components: $\text{Binom}(k; \tau = 1)$ centred at $k = 0$, $\text{Binom}(k - 2; \tau = 1)$ centred at $k = 2$, and $\text{Binom}(k - 5; \tau = 1)$ centred at $k = 5$. To account for the increasing range, we increase the parameter τ for longer walks, thereby widening the support of the distribution.

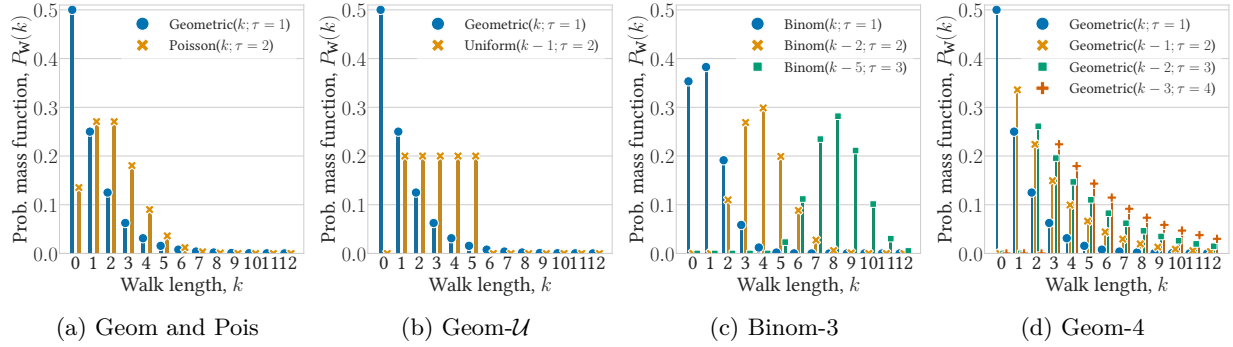


Figure 11: Plots of the walk length distributions P_w listed in Table 3. (a) show the two distributions used to create single-scale ReachNEs embeddings, while (b), (c) and (d) show the distributions used for multi-scale embeddings, with 2, 3 and 4 components respectively.

In contrast, Geom-4 uses four geometric components with progressively increasing mode and tail width: $\text{Geom}(k; \tau = 1)$, $\text{Geom}(k - 1; \tau = 2)$, $\text{Geom}(k - 2; \tau = 3)$, and $\text{Geom}(k - 3; \tau = 4)$. Unlike Binom-3, these distributions are not centred on long-range steps but exhibit heavy tails, meaning they still provide smoothing over distant nodes. The advantage of Geom-4 is that it places distinct modes over steps $k \in \{0, 1, 2, 3\}$, a range where most node-distinguishing information typically resides. This focus on short- and mid-range steps aligns with the goal of preserving local structure while retaining some long-range awareness.

D.2 Hyperparameters for proximity embedding comparison

Table 13 shows the hyperparameter search grids for each digraph proximity method compared in Section 6.4. For all models, we tune the number of embedding dimensions using $p \in [64, 1024]$. Note that for HOPE, we report halved values, as the implementation internally doubles the embedding size.

For DirSwitch- r , we tune both the number of directed steps r and the choice of walk length distributions P_w . For the other baselines, we tune their respective hyperparameters that influence proximity scale. These include:

- the β parameter for HOPE,
- the number of walk steps and jump factor for APP,
- the number of walk steps for NERD,
- the number of layers for BLADE.

We also tune selected learning-related hyperparameters: the number of sampled walks for NERD, the loss trade-off parameter λ for DGGAN, and the learning rate and number of epochs for BLADE.

With this setup, the total grid sizes are approximately comparable across methods. DirSwitch uses 75 different hyperparameter combinations. Among the baselines, APP and BLADE have the largest grids, with 120 and 300 combinations respectively. DGGAN has the smallest grid (45 combinations), as it is significantly more computationally expensive to train, making extensive tuning infeasible.

The best-performing hyperparameter settings for each model and dataset, based on cross-validation accuracy, are listed in Table 14. For DirSwitch- r , we observe that multi-scale embeddings with three or four components consistently yield the best results, aligning with prior findings on the benefits of multi-scale embedding strategies (Rozemberczki et al., 2021).

Table 13: Hyperparameter grids used for the comparison of proximity embedding models in Section 6.4.

MODEL	HYPERPARAMETER	VALUES
DIRSWITCH- r	EMB. DIM. p	64, 128, 256, 512, 1024
	r	1, 2, 3
	P_w	GEOM, POIS, GEOM- \mathcal{U} , BINOM-3, GEOM-4
HOPE	β	0.001, 0.01, 0.1, 0.2, 0.35, 0.5, 0.65, 0.8, 0.9, 1.0
	EMB. DIM. p	32, 64, 128, 256, 512
APP	EMB. DIM. p	64, 128, 256, 512, 1024
	JUMP FACTOR	0.1, 0.25, 0.5, 0.75
	WALK STEPS	2, 3, 4, 5, 6, 7
NERD	EMB. DIM. p	64, 128, 256, 512, 1024
	NUM. SAMPLES	1, 2, 3
	WALK STEPS	2, 3, 4, 5, 6
DGGAN	EMB. DIM. p	64, 128, 256, 512, 1024
	λ	5E-06, 1E-05, 5E-05
	LEARNING RATE	5E-05, 0.0001, 0.0005
BLADE	EMB. DIM. p	64, 128, 256, 512, 1024
	LEARNING RATE	0.0001, 0.001, 0.01
	NUM. EPOCHS	10, 30, 50, 100
	NUM. LAYERS	2, 3, 4, 5, 6

Table 14: The best hyperparameter values for each proximity embedding model and dataset.

MODEL	FLY LARVA	EU-EMAIL	POLBLOGS	CoCite	PUBMED	CORA (SUBELJ)
DirSWITCH	EMB. DIM. p 256	EMB. DIM. p 1024	EMB. DIM. p 128	EMB. DIM. p 1024	EMB. DIM. p 1024	EMB. DIM. p 1024
	r 3	r 2	r 1	r 1	r 1	r 1
	P_w GEOM-4	P_w BINOM-3	P_w BINOM-3	P_w BINOM-3	P_w BINOM-3	P_w BINOM-3
HOPE	β 1	β 1	β 0.9	β 0.9	β 0.9	β 1
	EMB. DIM. p 512	EMB. DIM. p 128	EMB. DIM. p 32	EMB. DIM. p 512	EMB. DIM. p 512	EMB. DIM. p 512
APP	EMB. DIM. p 1024	EMB. DIM. p 512	EMB. DIM. p 64	EMB. DIM. p 512	EMB. DIM. p 1024	EMB. DIM. p 1024
	JUMP FACTOR 0.75	JUMP FACTOR 0.75	JUMP FACTOR 0.5	JUMP FACTOR 0.25	JUMP FACTOR 0.75	JUMP FACTOR 0.75
	WALK STEPS 5	WALK STEPS 4	WALK STEPS 5	WALK STEPS 5	WALK STEPS 5	WALK STEPS 7
NERD	EMB. DIM. p 256	EMB. DIM. p 1024	EMB. DIM. p 128	EMB. DIM. p 256	EMB. DIM. p 512	EMB. DIM. p 256
	NUM. SAMPLES 3	NUM. SAMPLES 3	NUM. SAMPLES 2	NUM. SAMPLES 3	NUM. SAMPLES 3	NUM. SAMPLES 3
	WALK STEPS 5	WALK STEPS 5	WALK STEPS 5	WALK STEPS 6	WALK STEPS 6	WALK STEPS 6
DGGAN	EMB. DIM. p 256	EMB. DIM. p 512	EMB. DIM. p 128	EMB. DIM. p 128	EMB. DIM. p 128	EMB. DIM. p 256
	λ 5E-06	λ 1E-05	λ 5E-05	λ 1E-05	λ 5E-05	λ 5E-06
	LEARNING RATE 0.0005	LEARNING RATE 0.0005	LEARNING RATE 5E-05	LEARNING RATE 5E-05	LEARNING RATE 0.0001	LEARNING RATE 5E-05
BLADE	EMB. DIM. p 1024	EMB. DIM. p 1024	EMB. DIM. p 512	EMB. DIM. p 1024	EMB. DIM. p 1024	EMB. DIM. p 1024
	LEARNING RATE 0.0001	LEARNING RATE 0.0001	LEARNING RATE 0.001	LEARNING RATE 0.0001	LEARNING RATE 0.0001	LEARNING RATE 0.0001
	NUM. EPOCHS 10	NUM. EPOCHS 10	NUM. EPOCHS 10	NUM. EPOCHS 10	NUM. EPOCHS 10	NUM. EPOCHS 30
	NUM. LAYERS 3	NUM. LAYERS 2	NUM. LAYERS 6	NUM. LAYERS 4	NUM. LAYERS 6	NUM. LAYERS 2

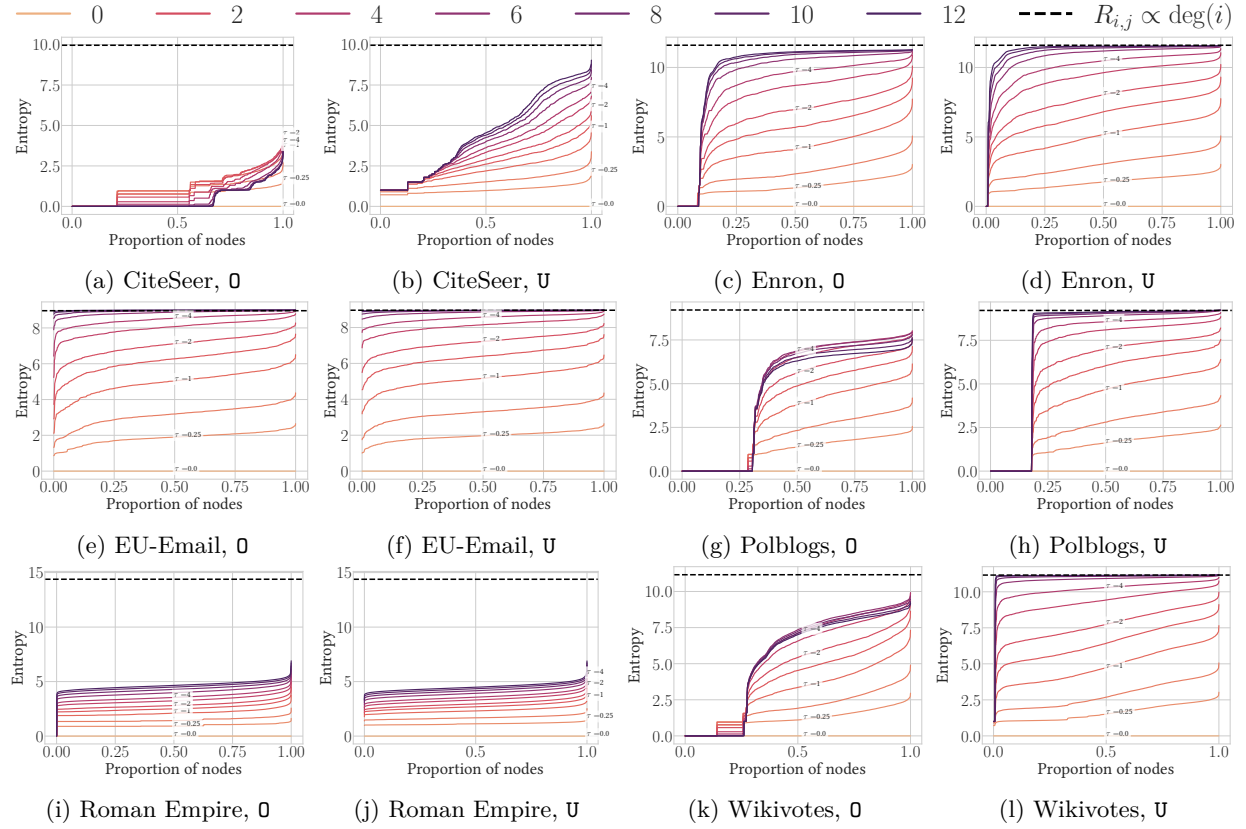


Figure 12: Reachability under $P_w = \text{Pois}(\tau)$, using \mathbf{A}_0 and \mathbf{A}_U . The x-axes correspond to nodes sorted by their entropy value. The colours represent various values of τ . The black dashed line indicates the entropy of a reachability distribution proportional to the node degrees, which corresponds to the limiting and uninformative distribution as $\tau \rightarrow \infty$. See Appendix C.1 for further details.

E Additional analysis of local sinks and dispersal

This section provides additional analysis of local sinks and reachability dispersal for graph datasets not included in the main paper due to space constraints. As before, we use *reachability entropy* to quantify dispersal. Recall that the reachability entropy of node j is defined as:

$$H(j; \mathbf{R}) = - \sum_{i=1}^n R_{i,j} \log_2 R_{i,j}, \quad (34)$$

where $R_{i,j}$ is the (i, j) -th element of the reachability matrix $\mathbf{R}(P_w, \sigma)$, computed under a walk length distribution P_w and edge direction specifier σ .

Figure 12 complements Figure 5 by presenting results for six additional datasets, using the Poisson walk length distribution $P_w = \text{Pois}(\tau)$. Each coloured curve corresponds to a different value of τ . The figure focuses on comparing the dispersal effects of default directed edges ($\sigma = 0$) versus undirected edges ($\sigma = U$).

Consistent with the examples in Section 3.1, we observe substantial entropy gains when using undirected edges on sparse and moderately dense graphs. For example, the sparse CiteSeer graph shows a pronounced increase in entropy, while medium-density graphs such as Enron, Polblogs, and Wikivotes also exhibit clear improvements. These gains become more pronounced at larger τ values, as reachability under directed edges increasingly concentrates in local sinks.

In contrast, the dense EU-Email graph, with an average node degree of 25, shows only minor differences between the directed and undirected cases. While small increases in entropy are still visible (cf. Figure 12e and 12f), the graph’s high density mitigates sink formation and enables more uniform information propagation.

Figure 13 extends Figure 6 by presenting results for eleven additional datasets. We again use $P_w = \text{Pois}(\tau = 2)$ to model short-range walks and $P_w = \mathcal{U}(\tau = 5)$ for longer-range behaviour.

The figure compares DirSwitch edge direction specifiers with those of MultiDir. The results closely follow the trends observed in Figure 5, with DirSwitch consistently achieving higher entropy than MultiDir—especially under the long-range distribution $P_w = \mathcal{U}(\tau)$. As before, the differences are most pronounced on sparse graphs, including Arxiv, CoCite, Pubmed, and Snap Patents.

The Roman Empire graph constitutes a counterexample to the trends described above. As shown in Figure 12i and 12j, there is virtually no difference in entropy between $\sigma = 0$ and $\sigma = \mathcal{U}$, despite the graph’s sparsity. In fact, close inspection reveals that entropy is slightly *lower* for the undirected case, particularly for large values of τ . This effect is even more pronounced in Figure 13j, where $\sigma = 0$ yields the highest entropy under the uniform walk length distribution. This is the only dataset where such behaviour is observed.

To understand this counterintuitive result, we must examine the specific structural properties of the Roman Empire graph. As shown in Table 2, this graph stands out with an unusually long average path length—around 2400 steps. The underlying reason is its chain-like topology: the graph primarily consists of a long linear sequence of nodes, where each node typically has an in-degree and out-degree of one. Occasionally, nodes branch off from the main chain and later reconnect.

Figure 14 visualizes a subgraph of the Roman Empire graph, clearly illustrating this structure. It also highlights why entropy may decrease in the undirected case. The node colours represent reachability values from the same starting node under $\sigma = 0$ on the left, and $\sigma = \mathcal{U}$ on the right, both using the uniform walk length distribution. As shown in Figure 14a, reachability in the directed case is spread more uniformly along the chain. In contrast, Figure 14b shows that the undirected case concentrates reachability around the starting node.

This is a consequence of the chain structure: in the directed case, random walks proceed almost exclusively forward along the chain. In the undirected case, however, each step introduces a 33–50% chance of stepping backward. This backward drift increases the likelihood of returning to or remaining near the starting node, reducing entropy.

F Additional graph alignment results

Figure 15 presents graph alignment results for three additional datasets, complementing the four datasets shown in Figure 7. As in the main paper, the y-axes indicate graph alignment accuracy, while the x-axes show values of τ . The top row corresponds to $P_w = \text{Geom}(\tau)$, and the bottom row to $P_w = \text{Pois}(\tau)$.

The results in Figure 15 are consistent with the findings from the main paper. For both DirSwitch- r and MultiDir- r , alignment accuracy increases with r , reflecting the benefit of representing a larger number of directed neighbourhoods. In contrast, models using $\sigma = 0$ and $\sigma = \mathcal{U}$ consistently achieve the lowest accuracies, as they lack the capacity to capture multiple directed neighbourhoods. DirSwitch- r generally outperforms MultiDir- r , highlighting that mitigating the effects of local sinks contributes to improved alignment accuracy.

We again observe the distinct behaviours of the Poisson and geometric walk length distributions: increasing τ tends to degrade accuracy under $\text{Pois}(\tau)$, whereas accuracy plateaus under $\text{Geom}(\tau)$, as previously discussed.

A notable exception is the Pubmed dataset, where DirSwitch- r with $P_w = \text{Pois}(\tau)$ reaches peak accuracy around $\tau \approx 9$, indicating that a relatively large receptive field is required for optimal performance. This contrasts with the other datasets, which generally attain their highest accuracy for $\tau < 5$.

This behaviour can be attributed to the structure of the Pubmed graph. As shown in Table 2, Pubmed has a very sparse local structure—its median out-degree is 0, and its median in-degree is 1—meaning that the short-range neighbourhoods of many nodes are nearly indistinguishable. Additionally, Pubmed exhibits the longest average path length among the alignment datasets. These factors necessitate a broader receptive field to generate more expressive and distinguishable embeddings.

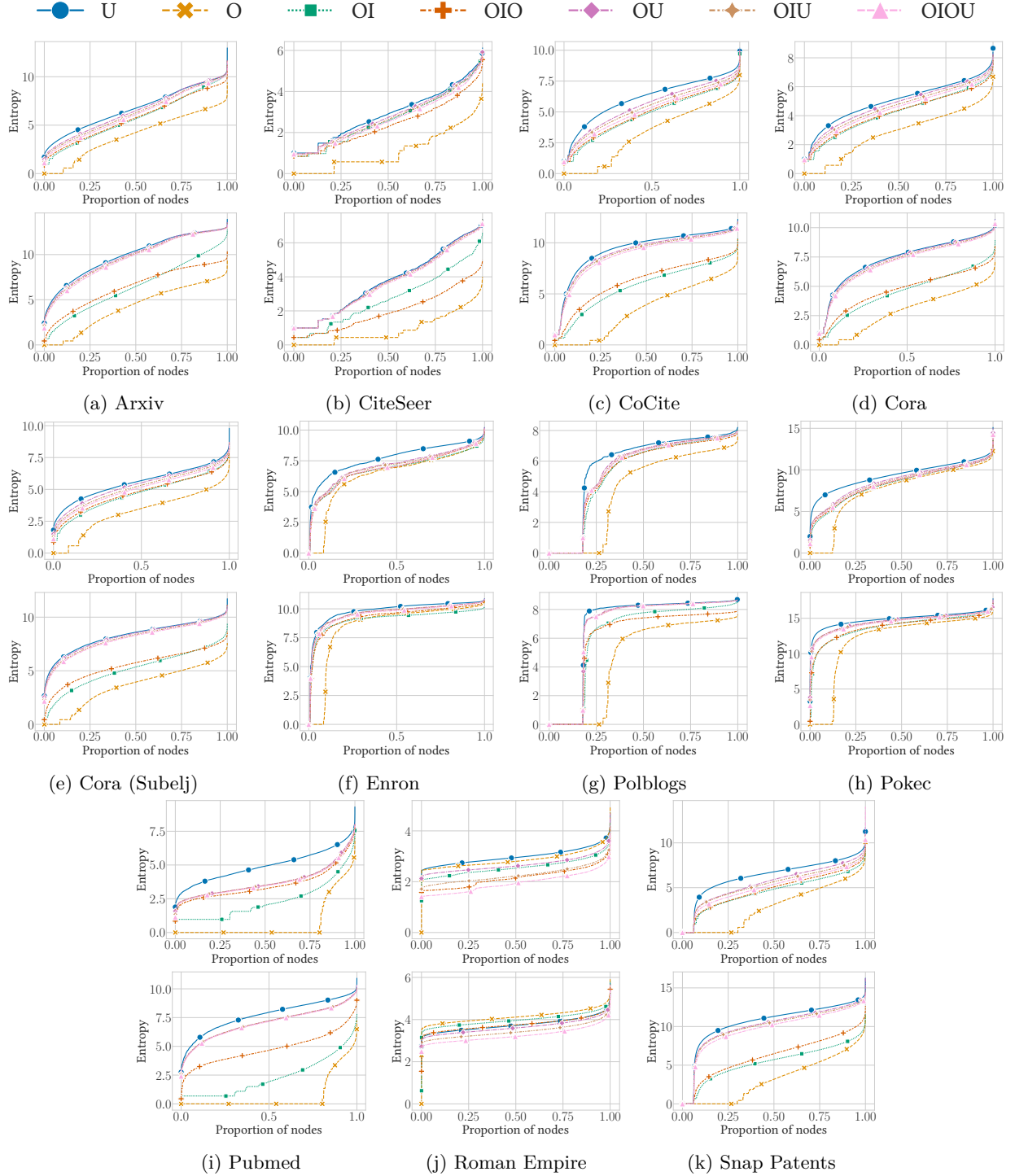


Figure 13: Neighbourhood dispersal evaluation for 11 graphs, measured via reachability entropy, $-\sum_{i=1}^n R_{i,j} \log_2 R_{i,j}$, computed for each node and sorted. Each curve corresponds to a different edge direction specifier σ , with the top row showing results for $P_w(k; \tau) = \text{Pois}(\tau = 2)$ (local dispersal) and the bottom row for $\mathcal{U}(\tau = 5)$ (long-range dispersal). DirSwitch variants (e.g., OU, OIU, OIOU) demonstrate high dispersal, comparable to U, while purely directed specifiers (O, OI, OIO) exhibit lower entropy due to sink effects.

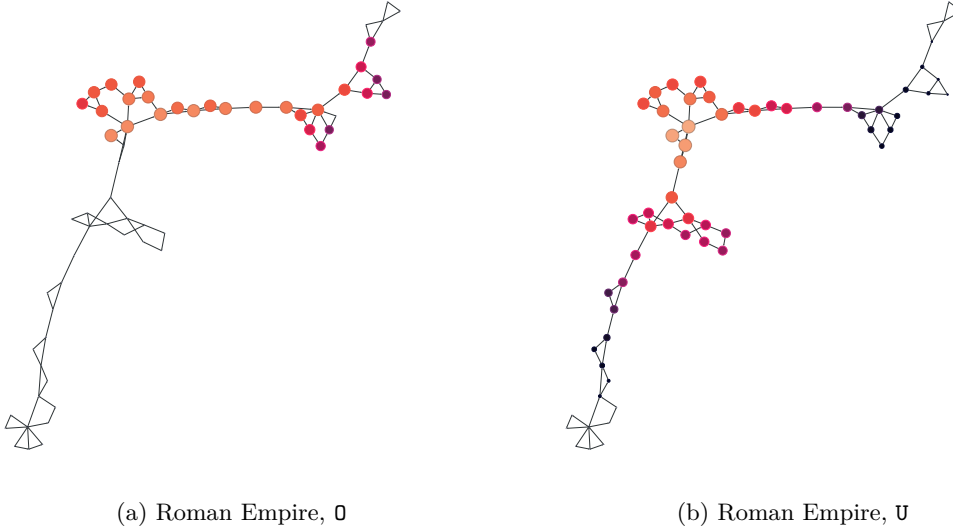


Figure 14: Both figures show the same subgraph of the Roman Empire graph. The node colours reflect the reachability vector for the same node under the uniform walk length distribution. (a) uses default directed edges ($\sigma = 0$) while (b) uses undirected edge ($\sigma = U$).

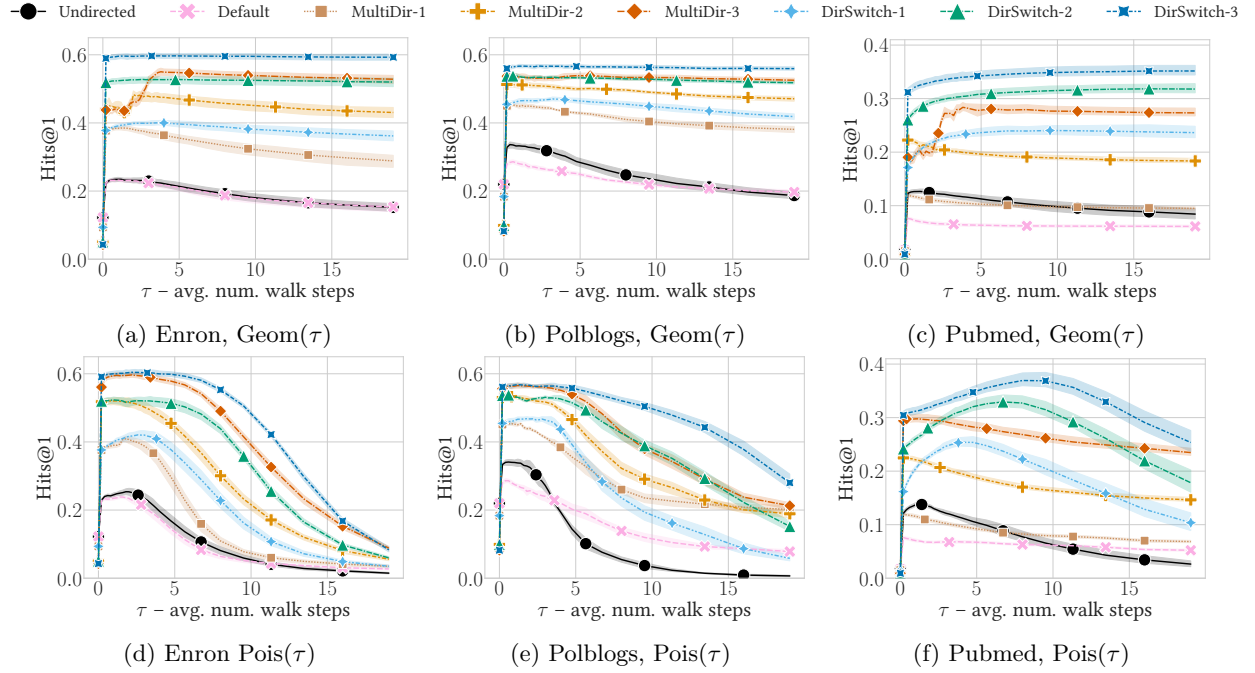


Figure 15: Evaluation of edge direction expressivity for four graphs using the geometric, Geom(τ), and Poisson, Pois(τ), walk length distributions. The y-axes represent graph alignment accuracy under 15% edge removal, while the x-axes correspond to τ , the average walk length. The curve colours and styles denote different sets of edge direction specifiers, σ .

G Additional DirSwitch embedding quality evaluation

This section provides further analysis and discussion of the results presented in Section 6.3, focusing on the benefits of multi-scale embeddings, the limitations introduced by fixed embedding dimensionality and the presence of 2-step homophily in the Pokec dataset.

Multi-scale embeddings and diminishing returns. The results in Tables 4 and 6 offer clear evidence of the benefits of multi-scale embeddings. Across all datasets, we observe at least moderate improvements when transitioning from single-scale to multi-scale walk length distributions, with several datasets exhibiting substantial accuracy gains. For instance, on Cora-ML, the classification accuracy for undirected embeddings increases from 31% to 76%; on Citeseer, from 76% to 89%; and on Cora, from 59% to 69% when switching from **Geom** to **Geom-4**. These results highlight the importance of embedding expressivity in unsupervised learning.

Similar trends are observed for proximity embeddings in Table 6. For example, EU-Email improves from 30% to 72%, and Polblogs from 64% to 85% under the same change in walk length distributions.

However, the benefits of multi-scale embeddings diminish as the number of multi-directional neighbourhoods increases, particularly for **MultiDir- r** and **DirSwitch- r** . In Cora-ML, for example, **DirSwitch-1** improves from 57% to 82% when switching from **Geom** to **Geom-4**, whereas **DirSwitch-3** shows a more modest improvement from 79% to 86%.

This trend is even more pronounced in the results for $p = 512$ embedding dimensions (Tables 16 and 18). In these cases, increasing both the number of multi-directional neighbourhoods and the number of scales can lead to accuracy degradation. For example, in Table 16, the accuracy for Cora declines as r increases and when using **Binom-3** or **Geom-4**. A similar pattern is observed for EU-Email in Table 18.

This behaviour arises from the diminishing returns of combining multi-directional and multi-scale embeddings via concatenation under fixed dimensionality constraints. As described in Section 4, the total embedding dimension p must be divided among all combinations of direction specifiers and walk length distributions. When p is fixed, increasing the number of concatenated components reduces the dimensionality available for each, limiting their expressivity. As a result, performance gains plateau at $p = 1024$ and begin to decline at $p = 512$.

While increasing p can help alleviate this issue by allowing more capacity for embedding components, this approach has practical limitations. Larger embeddings require more memory and computational resources, and may introduce challenges such as overfitting or the curse of dimensionality in downstream tasks. Future work is needed to explore more efficient ways to compress and integrate multi-directional and multi-scale information into compact embeddings without sacrificing expressivity.

Homophily in the Pokec dataset. As noted in the main paper, the classification results for the Pokec dataset in Table 4 exhibit behaviour that diverges from other heterophilic datasets. In particular, we observe a notable accuracy increase when using **Geom-4**, especially for undirected edges: accuracy rises from 62% with **Binom-3** to 72% with **Geom-4**. While Pokec typically follows the heterophilic trend—where undirected embeddings underperform compared to **DirSwitch- r** and **MultiDir- r** —this pattern is reversed for **Geom-4**.

To understand this, we need to examine the structural properties of the Pokec dataset. Pokec is an online Slovak social network, with nodes representing users and labels corresponding to reported gender. As reported by Lim et al. (2021), who curated this benchmark, the graph is heterophilic at the 1-step level due to the predominance of heterosexual connections. This is reflected in a low 1-step homophily score of 0.43.

However, the homophily increases significantly at the 2-step level, rising to 0.61. This indicates that while direct connections tend to be between users of different genders, the extended 2-step neighbourhoods are more gender-homogeneous. In such cases, undirected smoothing becomes more effective, especially when the embedding method emphasizes these mid-range neighbourhoods.

This is precisely what the **Geom-4** distribution achieves. As seen in Figure 11d, one of its components is **Geom($k - 2; \tau = 3$)**, which emphasizes 2-step neighbourhoods more than the other multi-scale distributions. This focus on 2-step structure explains the performance gain observed with **Geom-4** and undirected edges.

Table 15: Further investigation into the Pokec dataset. The walk length distributions are single-scale indicator distributions that place the entire probability mass on a single k -value.

(a) $p = 512$					(b) $p = 1024$				
EDGE DIRECTIONS	$\mathbf{1}_{k=0}$	$\mathbf{1}_{k=1}$	$\mathbf{1}_{k=2}$	$\mathbf{1}_{k=3}$	EDGE DIRECTIONS	$\mathbf{1}_{k=0}$	$\mathbf{1}_{k=1}$	$\mathbf{1}_{k=2}$	$\mathbf{1}_{k=3}$
DEFAULT	61.4±0.1	63.0±0.1	63.2±0.1	62.5±0.1	DEFAULT	61.4±0.1	63.0±0.1	63.1±0.1	62.5±0.1
UNDIRECTED	61.4±0.1	63.0±0.1	69.9±0.1	66.2±0.1	UNDIRECTED	61.4±0.1	63.0±0.1	69.9±0.0	66.2±0.1
MULTIDIR-1	61.4±0.1	63.4±0.1	65.3±0.1	63.3±0.1	MULTIDIR-1	61.4±0.1	63.4±0.1	65.3±0.1	63.3±0.1
MULTIDIR-2	61.4±0.1	63.4±0.1	68.0±0.1	63.9±0.1	MULTIDIR-2	61.4±0.1	63.4±0.1	68.0±0.1	63.9±0.1
DIRSWITCH-1	61.4±0.1	63.4±0.1	68.0±0.1	66.3±0.1	DIRSWITCH-1	61.4±0.1	63.4±0.1	68.0±0.1	66.3±0.1
DIRSWITCH-2	61.4±0.1	63.4±0.1	68.0±0.1	66.1±0.1	DIRSWITCH-2	61.4±0.1	63.4±0.1	68.0±0.1	66.1±0.1
DIRSWITCH-3	61.4±0.1	60.7±0.1	67.1±0.1	63.5±0.1	DIRSWITCH-3	61.4±0.1	63.4±0.1	68.0±0.1	65.9±0.1

To further validate this interpretation, we conduct an additional experiment using indicator walk length distributions that place all probability mass on a single step length:

$$\mathbf{P}_w = \mathbf{1}_{k=\tau} = \begin{cases} 1 & \text{if } k = \tau \\ 0 & \text{otherwise} \end{cases}. \quad (35)$$

This isolates the effect of each specific neighbourhood scale. The results, shown in Table 15, confirm that 2-step neighbourhoods ($\mathbf{1}_{k=2}$) yield the highest accuracy for undirected edges. In contrast, no similar performance gain is observed for directed edges (Default), reinforcing the notion that the combination of undirected smoothing and mid-range reachability is key to improved performance on Pokec.

Beyond explaining this anomaly, the analysis highlights a broader point: homophily and heterophily are not binary properties but exist on a spectrum. The optimal embedding strategy may vary across this spectrum, and models tailored exclusively for one end may perform poorly on datasets lying in the middle. This underscores the utility of DirSwitch, which we have shown to perform consistently well across both homophilic and heterophilic settings.

Unsupervised community detection. Our evaluation of DirSwitch in Section 6.3 focused on node classification, which offers a convenient and objective metric: classification accuracy.

In contrast, evaluating unsupervised downstream tasks is more challenging, as these tasks often lack a single "correct" solution. For example, many valid groupings of nodes can exist in a graph, yet only one such grouping is typically labelled for evaluation. Similar issues arise in other unsupervised tasks such as anomaly detection. As a result, these evaluations often require more nuanced analysis.

Nevertheless, since DirSwitch is intended for unsupervised embedding learning, it is important to assess its performance on unsupervised tasks as well. We therefore include a community detection benchmark. For this, we use the same proximity node classification datasets from Section 6.3, treating the node labels as ground-truth communities. We apply k -means clustering (Arthur & Vassilvitskii, 2007) to the embeddings, setting the number of clusters equal to the number of labelled communities. Performance is reported using normalized mutual information (NMI), averaged over 10 runs with different random seeds to compute mean and standard deviation.

Tables 20 and 21 present the community detection results for DirSwitch- r with $p = 512$ and $p = 1024$, respectively. The findings mirror those from the node classification experiments: DirSwitch- r generally performs well across datasets, while the sparse citation graphs (CoCite, Pubmed, and Cora (subelj)) tend to benefit from fully undirected edge orientation.

Results on the Polblogs dataset are particularly noteworthy due to their large standard deviations. This reflects the challenge of evaluating unsupervised tasks: we interpret these fluctuations as evidence that multiple valid community structures exist, with k -means sometimes recovering the one aligned with the labels, and sometimes converging to an alternative.

Table 16: Node classification accuracy for message-passing ReachNEs with $p = 512$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. Each entry reports the mean accuracy and standard deviation. Bold blue highlights the highest accuracy in each column, with light blue indicating results within one standard deviation of the best. Similarly, bold orange denotes the lowest accuracy, and light orange represents values within one standard deviation of the worst.

EDGE DIRECTIONS	ARXIV					ARXIV YEAR				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	55.0±0.2	42.4±0.4	55.4±0.2	54.2±0.2	55.4±0.2	36.5±0.3	35.2±0.5	36.7±0.2	37.8±0.3	38.2±0.3
UNDIRECTED	61.6±0.2	64.9±0.2	67.0±0.2	69.8±0.2	69.1±0.2	36.5±0.2	36.5±0.2	37.1±0.2	37.6±0.3	38.3±0.2
MULTIDIR-1	59.8±0.3	59.4±0.3	60.1±0.3	60.2±0.3	60.0±0.3	37.0±0.3	35.1±0.7	38.1±0.2	39.1±0.2	39.4±0.3
MULTIDIR-2	61.2±0.3	65.3±0.3	65.1±0.3	64.4±0.3	64.6±0.3	38.6±0.3	38.0±0.3	40.5±0.3	41.1±0.2	41.0±0.3
MULTIDIR-3	60.9±0.3	65.5±0.3	65.3±0.3	64.7±0.3	64.1±0.2	39.7±0.3	40.2±0.3	41.4±0.3	42.1±0.3	41.3±0.2
DIRSWITCH-1	61.5±0.3	65.9±0.2	67.3±0.3	69.5±0.2	68.6±0.3	38.2±0.2	37.8±0.2	39.9±0.3	39.8±0.3	41.2±0.3
DIRSWITCH-2	61.5±0.3	65.7±0.2	66.7±0.3	68.9±0.3	67.8±0.3	39.6±0.3	39.5±0.3	40.9±0.3	41.1±0.3	41.5±0.3
DIRSWITCH-3	60.8±0.3	65.6±0.3	65.9±0.3	67.7±0.3	66.3±0.3	39.6±0.3	40.5±0.3	41.4±0.2	42.1±0.2	42.2±0.3

EDGE DIRECTIONS	CORA-ML					CITSEER					CORA				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	50.5±3.1	51.9±2.1	64.2±2.6	70.5±1.3	73.5±1.7	70.7±1.8	58.4±2.2	74.1±2.6	75.8±1.9	78.2±1.6	57.4±0.7	54.1±0.8	60.7±0.8	60.9±0.8	61.1±0.9
UNDIRECTED	60.4±2.5	68.1±2.2	73.8±1.8	82.2±1.3	83.7±1.4	80.6±1.4	86.6±1.5	87.0±1.5	93.9±0.9	92.8±1.1	64.9±0.7	68.1±0.7	69.2±0.6	70.1±0.9	70.0±0.8
MULTIDIR-1	68.5±2.4	67.5±2.4	75.2±2.6	77.2±2.2	79.4±2.3	76.1±1.9	75.4±2.1	80.7±1.9	81.9±1.5	82.8±1.5	65.5±0.7	64.8±0.8	65.3±1.1	63.5±1.0	62.7±0.9
MULTIDIR-2	79.1±1.7	78.8±1.6	82.7±1.6	83.1±1.4	84.3±1.2	83.9±1.7	85.1±1.7	86.0±1.4	85.8±1.6	86.6±1.2	66.7±1.1	68.1±0.9	66.2±0.9	64.5±1.0	63.6±1.1
MULTIDIR-3	83.1±1.5	83.6±1.7	85.2±1.5	84.8±1.5	84.7±1.6	85.2±1.5	88.0±1.4	87.0±1.3	87.1±1.2	87.2±0.9	65.0±0.8	67.2±1.0	64.4±1.0	62.0±1.1	61.1±0.9
DIRSWITCH-1	72.9±1.9	73.7±2.1	81.3±1.9	85.0±1.5	85.2±1.2	82.8±1.0	88.4±1.2	89.2±1.3	91.5±0.9	90.7±1.2	67.5±0.9	69.3±0.7	69.1±0.9	69.2±1.0	68.4±0.7
DIRSWITCH-2	79.5±2.0	80.8±1.9	84.1±1.6	85.7±1.9	85.7±1.5	85.3±1.1	88.8±0.9	88.4±1.1	90.0±1.2	88.8±1.3	66.7±0.9	68.6±0.8	67.4±0.8	67.4±0.6	65.8±0.7
DIRSWITCH-3	82.8±1.3	83.5±1.4	85.3±1.0	86.1±1.3	85.7±1.4	85.4±1.3	88.2±1.1	87.1±0.9	88.4±0.8	87.5±1.1	64.9±0.9	67.1±0.9	64.5±0.9	64.7±0.5	62.2±0.6

EDGE DIRECTIONS	ROMAN EMPIRE					POKEC					SNAP-PATENTS				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	69.3±0.9	32.2±0.6	67.0±0.5	49.9±0.6	69.6±0.6	60.3±0.1	59.0±0.8	61.6±0.1	60.7±0.2	65.5±0.1	35.6±0.1	41.7±0.1	39.1±0.0	40.8±0.1	40.9±0.1
UNDIRECTED	67.4±0.9	47.0±0.9	69.3±0.7	63.7±0.6	70.0±0.7	60.3±0.1	59.4±0.4	62.0±0.1	62.3±0.1	71.7±0.1	30.8±0.1	29.9±0.0	31.7±0.1	32.7±0.1	33.9±0.4
MULTIDIR-1	70.1±0.8	40.7±0.8	68.8±0.5	59.6±0.8	71.5±0.6	61.2±0.1	60.2±0.0	62.3±0.1	61.7±0.1	67.3±0.1	40.9±0.0	43.4±0.3	43.9±0.2	46.7±0.1	45.2±0.1
MULTIDIR-2	71.6±0.6	63.2±0.7	72.0±0.8	66.7±0.5	75.5±0.7	61.6±0.1	61.2±0.2	63.2±0.1	62.7±0.1	67.9±0.1	42.9±0.0	44.0±0.1	45.3±0.2	47.4±0.2	45.0±0.1
MULTIDIR-3	72.0±0.5	69.1±0.5	73.0±0.8	67.4±0.7	72.5±0.4	61.8±0.1	63.5±0.1	63.1±0.1	62.3±0.1	62.4±0.2	44.2±0.1	47.2±0.1	45.2±0.0	46.4±0.1	44.8±0.1
DIRSWITCH-1	71.6±0.5	57.1±0.9	73.8±0.5	69.2±0.7	75.3±0.7	61.2±0.1	60.4±0.3	62.8±0.1	62.8±0.1	69.9±0.0	40.7±0.0	40.8±0.1	40.4±0.1	41.6±0.1	42.7±0.4
DIRSWITCH-2	71.9±0.5	65.5±0.7	74.7±0.4	70.6±0.7	76.1±0.6	61.5±0.1	60.9±0.2	63.5±0.1	63.6±0.1	69.9±0.1	44.2±0.2	46.0±0.1	44.7±0.1	44.9±0.2	45.2±0.3
DIRSWITCH-3	71.8±0.5	69.1±0.5	74.3±0.5	70.1±0.4	72.7±0.5	61.6±0.1	62.7±0.1	62.9±0.1	62.1±0.1	62.0±0.1	44.4±0.1	48.8±0.1	46.3±0.1	46.9±0.2	46.9±0.2

Lastly, we compare DirSwitch- r against the state-of-the-art proximity embedding methods from Section 6.4. We reuse the same hyperparameters as in the node classification setting, without performing additional tuning (see Section D.2).

The results are shown in Table 22. DirSwitch and APP (Zhou et al., 2017) achieve the best overall performance. APP obtains higher NMI scores on EU-Email, CoCite, and Pubmed, while DirSwitch performs best on Fly Larva. The methods are tied on Cora (subelj). As before, we do not consider Polblogs due to its high variance.

As discussed in Section 5, APP is a sampled random walk method that asymptotically performs the same matrix decomposition as proximity DirSwitch with $\sigma = 0$ and $P_w = \text{Geom}(\tau)$. It is therefore noteworthy that APP outperforms DirSwitch in some cases.

One possible explanation is that the hyperparameters—tuned for node classification—may happen to favour APP in the community detection setting. Another plausible explanation is that APP benefits from random walk sampling itself. Similar effects have been observed in recommendation systems, where sampling-based matrix factorization often outperforms direct SVD due to reduced overfitting (Koren et al., 2009). Further research is needed to better understand the relationship between the sampled factorization in APP and the explicit decomposition used in proximity ReachNEs.

Table 17: Relative improvements in node classification accuracy for message-passing ReachNEs with $p = 512$. The values reflect the maximum accuracy for per row and dataset in Table 16. The top row displays absolute accuracies for the default edge directions, with standard deviations expressed as percentages. The subsequent rows present the relative improvements compared to the top row. The table is structured with the four homophilic datasets on the left and the four heterophilic datasets on the right.

EDGE DIRECTIONS	CORA-ML	CITeseer	CORA	ARXIV	ROMAN EMPIRE	ARXIV YEAR	POKEC	SNAP-PATENTS
DEFAULT	73.5±2.9%	78.2±2.6%	61.1±1.3%	55.4±0.5%	69.6±0.9%	38.2±0.8%	65.5±0.4%	41.7±0.2%
UNDIRECTED	+13.9%	+20.1%	+14.6%	+25.9%	+0.7%	+0.3%	+9.5%	-18.8%
MULTIDIR-1	+8.1%	+6.0%	+7.1%	+8.7%	+2.8%	+3.1%	+2.8%	+12.0%
MULTIDIR-2	+14.8%	+10.8%	+11.4%	+17.8%	+8.5%	+7.8%	+3.8%	+13.7%
MULTIDIR-3	+16.0%	+12.6%	+9.9%	+18.2%	+5.0%	+10.4%	-3.1%	+13.1%
DIRSWITCH-1	+16.0%	+17.0%	+13.3%	+25.4%	+8.3%	+7.9%	+6.8%	+2.4%
DIRSWITCH-2	+16.7%	+15.2%	+12.2%	+24.3%	+9.4%	+8.7%	+6.8%	+10.4%
DIRSWITCH-3	+17.3%	+13.1%	+9.8%	+22.1%	+6.8%	+10.5%	-3.9%	+17.1%

Table 18: Node classification accuracy for proximity ReachNEs with $p = 512$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. The values denote average accuracies with standard deviations. Bold blue highlights the best results in each column, with light blue indicating results within one standard deviation. Similarly, bold orange marks the worst results, with light orange showing values within one standard deviation of the lowest performance.

EDGE DIRECTIONS	FLYLARVA					EU-EMAIL					POLBLOGS				
	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4
DEFAULT	45.3±1.8	46.1±2.4	48.8±1.9	49.2±2.0	49.8±2.0	46.1±3.1	56.6±3.5	64.2±3.5	68.8±3.4	70.4±3.6	71.7±2.7	74.7±2.2	82.6±2.2	85.0±2.1	85.6±1.9
UNDIRECTED	47.8±2.0	50.4±1.9	53.6±2.0	55.1±1.7	54.9±2.0	50.3±3.8	58.6±3.3	70.1±3.2	75.3±2.7	75.8±2.6	77.0±2.6	79.2±2.7	85.6±1.8	87.4±2.1	87.0±1.9
MULTIDIR-1	48.5±1.9	48.5±2.1	51.4±2.3	52.5±1.9	52.0±1.8	65.4±3.0	67.5±3.1	73.7±3.1	73.8±3.0	73.7±2.7	81.7±1.9	83.1±2.1	86.4±2.0	87.7±1.7	88.4±1.7
MULTIDIR-2	51.8±2.3	51.4±1.7	53.2±1.9	52.7±2.3	52.7±2.1	74.0±2.7	74.5±2.7	75.2±3.3	72.5±2.9	70.4±2.7	87.9±1.8	87.6±1.7	89.1±1.7	89.2±1.6	89.6±1.6
MULTIDIR-3	53.4±1.9	53.5±2.1	53.7±1.9	54.0±1.6	54.7±1.5	73.9±2.3	75.3±3.0	70.0±2.5	66.4±3.4	63.6±2.7	89.0±1.4	89.5±1.3	89.8±1.8	89.5±2.0	89.2±1.4
DIRSWITCH-1	50.0±2.3	51.5±1.7	54.2±1.7	56.1±2.1	55.4±1.5	66.4±3.3	69.5±3.4	74.9±2.9	75.7±2.3	75.4±2.6	83.7±2.2	84.1±1.8	88.3±2.1	87.5±2.0	89.7±1.5
DIRSWITCH-2	52.4±2.1	52.8±2.0	55.2±2.0	55.9±2.1	55.5±1.9	73.6±2.8	74.0±3.1	74.9±2.7	72.6±2.5	70.1±2.7	87.6±2.1	88.3±2.0	89.5±1.8	87.8±1.8	90.0±1.6
DIRSWITCH-3	54.2±1.9	54.6±1.6	54.9±1.9	56.2±1.7	56.0±2.2	74.8±3.0	74.6±2.7	70.5±3.0	66.2±3.0	63.9±3.3	88.8±1.7	89.4±1.6	89.7±1.6	89.7±1.4	90.0±1.9

EDGE DIRECTIONS	CoCITE					PUBMED					CORA (SUBELJ)				
	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4
DEFAULT	37.6±0.5	38.2±0.6	38.9±0.6	38.9±0.5	38.9±0.5	71.2±0.7	71.0±0.7	71.3±0.6	70.7±0.7	69.9±0.7	55.8±0.9	55.8±0.8	56.4±0.8	56.5±0.8	56.3±0.6
UNDIRECTED	44.7±0.5	45.7±0.4	46.4±0.5	47.7±0.6	47.3±0.5	81.5±0.6	81.8±0.6	81.6±0.6	82.1±0.6	81.5±0.7	65.1±0.7	65.7±0.8	66.1±0.8	66.4±0.8	66.3±0.7
MULTIDIR-1	38.8±0.6	39.3±0.6	40.0±0.5	40.1±0.5	40.3±0.5	71.5±0.6	71.0±0.6	70.2±0.7	68.9±0.7	68.4±0.8	55.5±0.7	55.6±0.9	55.5±0.7	55.3±0.7	55.3±0.9
MULTIDIR-2	41.5±0.6	42.3±0.5	42.9±0.7	44.0±0.5	44.0±0.5	76.2±0.6	76.0±0.7	75.1±0.7	74.4±0.8	74.1±0.6	58.1±0.7	58.6±0.6	58.6±0.8	58.7±0.7	58.9±0.7
MULTIDIR-3	42.3±0.5	43.2±0.5	44.3±0.5	45.9±0.4	45.6±0.3	78.3±0.6	78.0±0.7	77.0±0.7	76.7±0.8	75.9±0.7	58.8±0.9	59.3±0.9	59.0±0.7	58.8±1.0	59.7±0.8
DIRSWITCH-1	44.3±0.5	45.4±0.4	46.2±0.5	47.8±0.5	47.3±0.6	81.1±0.7	81.3±0.6	81.2±0.6	81.5±0.6	80.6±0.7	63.4±0.8	64.2±0.7	64.6±0.8	65.3±0.7	65.0±0.7
DIRSWITCH-2	43.1±0.6	44.5±0.5	45.4±0.4	48.0±0.4	46.8±0.5	80.8±0.6	81.0±0.7	80.5±0.6	80.9±0.7	80.0±0.6	61.1±0.8	62.4±0.7	62.6±0.8	64.0±0.8	63.4±0.8
DIRSWITCH-3	41.8±0.5	43.5±0.5	44.5±0.5	47.1±0.5	46.3±0.5	80.1±0.7	80.4±0.7	80.0±0.7	79.8±0.6	78.6±0.7	58.0±0.7	59.6±0.7	60.1±0.8	61.0±0.9	61.2±0.8

Table 19: Relative improvements in node classification accuracy for proximity ReachNEs with $p = 512$. The values reflect the maximum accuracy for per row and dataset in Table 18. The top row displays absolute accuracies for the default edge directions, with standard deviations expressed as percentages. The subsequent rows present the relative improvements compared to the top row. The table is structured with the three denser graphs on the left and the three sparser graphs on the right.

EDGE DIRECTIONS	FLYLARVA	EU-EMAIL	POLBLOGS	CoCITE	PUBMED	CORA (SUBELJ)
DEFAULT	49.8±4.1%	70.4±4.8%	85.6±2.6%	38.9±1.4%	71.3±0.9%	56.5±1.4%
UNDIRECTED	+10.7%	+7.8%	+2.0%	+22.4%	+15.1%	+17.7%
MULTIDIR-1	+5.4%	+4.9%	+3.2%	+3.5%	+0.3%	-1.5%
MULTIDIR-2	+6.8%	+6.8%	+4.6%	+13.1%	+6.9%	+4.4%
MULTIDIR-3	+9.8%	+7.1%	+4.8%	+17.8%	+9.8%	+5.8%
DIRSWITCH-1	+12.8%	+7.6%	+4.8%	+22.9%	+14.3%	+15.6%
DIRSWITCH-2	+12.4%	+6.5%	+5.1%	+23.2%	+13.5%	+13.4%
DIRSWITCH-3	+12.8%	+6.4%	+5.1%	+20.9%	+12.7%	+8.3%

Table 20: Community detection normalized mutual information (NMI) scores for proximity ReachNEs with $p = 512$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. The values denote average NMI with standard deviations. Bold blue highlights the best results in each column, with light blue indicating results within one standard deviation. Similarly, bold orange marks the worst results, with light orange showing values within one standard deviation of the lowest performance.

EDGE DIRECTIONS	FLY LARVA						EU-EMAIL						POLBLOGS					
	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4		GEOM	POIS	GEOM-U	BINOM-3	GEOM-4		GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	
DEFAULT	54.2±0.5	55.0±0.4	55.7±0.4	56.1±0.4	56.2±0.3		51.6±3.0	60.7±1.1	59.7±1.1	61.2±1.0	61.3±0.9		5.8±11.0	11.3±15.3	15.1±17.9	13.3±18.4	18.1±18.7	
UNDIRECTED	54.9±0.4	56.1±0.3	56.5±0.3	57.7±0.3	57.7±0.4		56.4±4.9	64.4±1.5	65.8±1.2	64.5±1.0	65.3±0.8		16.9±17.6	21.3±12.7	25.7±15.8	28.2±16.6	30.1±17.9	
MULTIDIR-1	55.3±0.5	56.0±0.4	56.6±0.5	56.7±0.5	57.2±0.3		57.9±2.0	60.4±1.9	61.9±1.0	63.0±1.1	63.1±0.8		25.8±14.9	35.0±13.3	36.8±10.5	34.6±12.2	38.7±12.4	
MULTIDIR-2	56.4±0.4	57.1±0.3	57.6±0.4	57.3±0.4	58.0±0.3		61.8±1.2	65.0±1.4	65.1±0.8	65.1±0.6	64.2±0.4		40.2±17.5	47.1±10.8	44.7±15.3	41.8±14.5	45.8±10.3	
MULTIDIR-3	57.0±0.5	57.7±0.4	58.4±0.4	57.6±0.3	58.3±0.2		64.0±0.8	66.5±0.8	64.7±0.5	60.0±0.6	58.6±0.6		27.0±16.6	41.9±17.1	43.3±16.2	46.7±10.4	45.1±14.5	
DIRSWITCH-1	55.5±0.3	56.6±0.5	57.2±0.3	58.2±0.4	58.2±0.4		59.8±2.3	64.0±1.4	65.0±1.1	64.9±0.9	65.6±0.6		17.5±13.3	23.8±15.5	31.2±17.9	25.7±16.1	26.8±16.8	
DIRSWITCH-2	56.4±0.4	57.2±0.4	58.0±0.3	58.9±0.4	59.0±0.3		62.8±1.3	66.0±1.0	65.8±0.8	65.3±0.8	64.7±0.5		23.0±15.6	30.2±16.3	32.8±18.1	29.5±17.8	29.8±17.3	
DIRSWITCH-3	57.1±0.4	57.8±0.4	58.9±0.3	59.2±0.3	59.8±0.4		64.0±0.9	66.7±0.6	65.0±0.5	60.4±0.5	59.3±0.8		27.3±17.3	45.4±16.7	32.9±16.6	29.3±16.3	30.5±15.5	
EDGE DIRECTIONS	CoCITE						PUBMED						CORA (SUBELJ)					
	GEOM	POIS	GEOM-U	BINOM-3	GEOM-4		GEOM	POIS	GEOM-U	BINOM-3	GEOM-4		GEOM	POIS	GEOM-U	BINOM-3	GEOM-4	
DEFAULT	14.9±2.3	16.7±1.7	17.0±1.5	17.1±1.8	17.2±2.2		2.7±0.9	2.7±1.0	2.7±1.0	2.4±0.6	2.6±1.3		31.4±1.1	31.8±0.8	32.2±0.7	33.7±0.7	33.3±0.6	
UNDIRECTED	20.4±1.1	22.4±1.0	23.9±1.0	28.0±0.5	26.7±0.9		11.9±4.1	15.4±2.9	17.7±2.4	21.6±4.0	20.6±3.0		43.2±0.4	45.9±0.4	46.8±0.4	49.8±0.2	48.6±0.3	
MULTIDIR-1	16.7±0.6	17.3±0.5	17.9±0.7	19.3±0.6	18.7±0.6		4.5±1.2	4.6±1.4	4.9±1.4	5.1±1.5	5.8±0.8		31.5±0.8	31.7±0.9	31.9±0.6	33.0±0.7	32.4±0.5	
MULTIDIR-2	18.1±0.7	19.2±0.6	20.1±0.7	22.6±0.6	22.4±0.4		5.7±1.4	6.2±1.3	6.6±1.1	7.7±1.4	7.5±1.0		35.8±0.8	36.4±0.8	35.6±0.7	36.2±0.5	35.7±0.7	
MULTIDIR-3	18.9±0.9	20.2±0.9	21.6±0.6	23.9±0.7	23.5±0.4		6.9±2.1	8.0±2.1	8.2±1.5	9.6±1.7	8.1±0.6		36.8±0.7	38.5±0.5	37.7±0.5	37.8±0.5	38.1±0.6	
DIRSWITCH-1	19.8±1.2	21.8±1.0	23.4±1.2	28.8±0.9	27.8±1.0		12.0±3.9	14.4±2.8	15.1±2.0	21.3±3.8	18.7±3.0		43.0±0.4	45.3±0.5	46.0±0.4	49.6±0.3	48.1±0.3	
DIRSWITCH-2	19.7±1.2	21.3±1.0	24.0±1.3	29.6±0.7	27.7±0.8		8.8±3.2	12.1±2.3	11.8±2.4	19.0±4.1	15.3±3.3		41.9±0.6	44.1±0.4	44.7±0.6	47.9±0.3	46.5±0.3	
DIRSWITCH-3	18.5±0.9	21.1±1.0	23.9±1.1	29.3±0.7	28.0±0.5		8.9±3.1	11.0±2.5	10.5±3.1	17.3±2.9	13.6±3.3		38.5±0.5	41.7±0.5	41.7±0.4	44.6±0.3	43.6±0.3	

Table 21: Community detection normalized mutual information (NMI) scores for proximity ReachNEs with $p = 1024$ embedding dimensions. Columns correspond to different datasets and multi-scale walk length distributions, while rows represent various edge direction specifiers. The values denote average NMI with standard deviations. Bold blue highlights the best results in each column, with light blue indicating results within one standard deviation. Similarly, bold orange marks the worst results, with light orange showing values within one standard deviation of the lowest performance.

EDGE DIRECTIONS	FLY LARVA						EU-EMAIL						POLBLOGS				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4		GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4		GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	52.9±0.5	54.3±0.5	55.0±0.4	55.6±0.4	55.7±0.4		54.0±1.8	59.1±1.6	60.0±1.3	62.1±1.3	60.5±1.0		11.1±15.3	13.8±17.9	16.4±18.4	12.4±18.3	13.3±18.0
UNDIRECTED	54.0±0.6	55.7±0.4	55.9±0.3	57.5±0.4	57.5±0.3		55.2±2.6	63.5±1.1	63.8±1.3	63.8±0.9	64.6±0.9		11.8±15.9	25.5±19.7	24.0±16.0	22.3±17.7	28.4±18.2
MULTIDIR-1	54.4±0.5	55.4±0.5	55.9±0.4	56.5±0.4	56.8±0.4		51.9±4.1	61.5±1.7	62.1±1.1	62.5±1.0	63.0±1.0		20.3±18.5	29.6±13.7	28.0±15.1	36.1±15.6	38.2±14.3
MULTIDIR-2	55.7±0.4	56.6±0.3	57.0±0.4	57.1±0.4	57.6±0.4		59.3±1.7	63.1±1.5	64.3±1.2	65.8±0.9	65.9±0.8		33.8±21.7	41.8±17.2	36.0±19.3	43.6±15.3	45.9±13.5
MULTIDIR-3	56.3±0.5	57.1±0.3	57.6±0.3	57.7±0.4	58.3±0.4		61.9±1.5	65.8±1.0	65.7±0.9	65.5±0.5	64.3±0.6		28.3±17.9	43.6±17.8	41.7±19.4	48.5±11.1	45.0±15.4
DIRSWITCH-1	54.7±0.4	56.2±0.4	56.7±0.3	57.8±0.3	57.9±0.3		53.6±5.3	64.5±1.3	65.0±1.3	64.7±1.1	65.0±1.3		11.2±14.1	22.0±13.7	23.9±15.2	29.5±18.0	25.6±15.9
DIRSWITCH-2	55.6±0.4	56.6±0.5	57.3±0.4	58.5±0.4	58.5±0.3		59.5±1.8	63.3±1.3	64.7±1.1	65.0±0.8	65.7±0.7		16.2±14.5	23.2±13.3	31.7±17.8	27.9±16.9	28.8±17.2
DIRSWITCH-3	56.5±0.3	57.2±0.3	58.0±0.3	58.9±0.3	59.2±0.3		62.4±1.7	65.7±1.1	65.8±0.8	65.2±0.5	64.6±0.7		23.3±15.4	35.5±18.0	31.9±18.7	38.4±17.4	27.7±17.6
EDGE DIRECTIONS	CoCITE						PUBMED						CORA (SUBELJ)				
	GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4		GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4		GEOM	POIS	GEOM- \mathcal{U}	BINOM-3	GEOM-4
DEFAULT	14.7±3.0	17.1±1.5	15.7±2.7	17.4±1.5	16.7±2.7		2.7±0.9	2.5±1.0	2.8±1.1	2.4±0.9	2.3±0.9		30.5±1.6	30.9±1.4	31.7±0.8	33.8±0.7	33.3±0.9
UNDIRECTED	20.0±1.3	22.7±0.8	23.4±1.1	27.6±0.8	26.2±1.0		10.8±2.9	15.0±3.2	16.5±3.8	21.1±4.1	20.6±2.7		42.9±0.6	45.6±0.5	46.4±0.4	49.8±0.3	49.0±0.4
MULTIDIR-1	17.0±0.6	17.2±0.7	17.7±0.6	19.6±0.5	19.2±0.5		3.8±1.9	4.3±1.6	5.1±1.5	5.6±1.8	5.4±1.8		31.6±0.9	31.5±1.1	31.6±0.9	33.0±0.8	32.6±0.7
MULTIDIR-2	18.2±0.5	19.3±0.6	20.1±0.7	22.8±0.5	22.1±0.7		5.1±1.7	6.8±1.2	6.4±1.1	7.2±2.3	7.1±1.6		36.7±0.8	37.0±0.8	36.2±0.7	37.1±0.7	36.4±0.6
MULTIDIR-3	19.2±0.9	20.0±0.8	21.0±0.6	24.1±0.6	23.4±0.4		6.6±1.6	7.2±1.4	8.1±1.8	10.7±1.7	8.8±1.5		38.8±0.6	39.8±0.5	39.3±0.5	39.4±0.5	38.9±0.4
DIRSWITCH-1	19.3±1.3	21.6±0.9	22.8±1.0	27.6±1.0	26.8±1.1		11.1±4.0	14.3±3.0	15.5±3.6	21.7±3.2	19.6±3.3		42.7±0.5	45.1±0.4	46.4±0.4	49.7±0.3	48.4±0.3
DIRSWITCH-2	19.3±1.0	21.1±1.0	23.6±0.8	28.6±1.2	27.5±1.2		10.3±3.8	13.2±3.2	13.9±2.5	21.2±4.1	17.5±3.9		42.4±0.4	44.6±0.4	45.5±0.3	49.3±0.2	47.8±0.2
DIRSWITCH-3	19.3±1.2	21.0±1.0	23.6±1.3	29.4±0.6	27.5±0.6		8.7±3.3	10.8±2.8	11.3±2.6	19.2±3.5	14.3±3.3		40.9±0.5	43.3±0.4	44.1±0.5	47.8±0.3	46.2±0.3

Table 22: Community detection normalized mutual information (NMI) scores for proximity node embedding models. Bold indicate the top NMI scores, and results within one standard deviation, for each dataset. Average and standard deviations are calculated over 10 different k -means initialisations, and 5 different random seeds for the embedding models.

MODEL	FLY LARVA	EU-EMAIL	POLBLOGS	CoCITE	PUBMED	CORA (SUBELJ)
HOPE ¹	46.3 ± 0.8	19.9 ± 0.5	9.4 ± 1.3	2.3 ± 0.4	0.4 ± 0.2	12.4 ± 1.1
APP ²	54.9 ± 0.4	68.7 ± 1.0	41.0 ± 12.9	31.7 ± 0.6	23.9 ± 0.2	50.0 ± 0.3
NERD ³	50.8 ± 0.4	50.0 ± 1.5	15.1 ± 16.4	3.3 ± 0.0	4.9 ± 4.6	22.0 ± 0.6
DGGAN ⁴	35.2 ± 0.7	21.1 ± 0.4	1.4 ± 0.0	3.6 ± 0.0	0.2 ± 0.1	4.2 ± 0.0
BLADE ⁵	54.3 ± 0.8	42.5 ± 1.2	19.6 ± 9.6	4.1 ± 0.6	0.6 ± 0.4	9.3 ± 0.6
DIRSWITCH	59.6 ± 0.7	65.0 ± 0.7	26.8 ± 16.4	27.8 ± 0.9	21.6 ± 3.8	49.7 ± 0.3

¹OU ET AL. (2016) ²ZHOU ET AL. (2017) ³KHOSLA ET AL. (2020) ⁴ZHU ET AL. (2021B) ⁵VIRINCHI & SALADI (2023)