# A Framework for Generalizing Graph-based Representation Learning Methods

**Nesreen K. Ahmed**
Intel Labs

**Ryan A. Rossi**
PARC

**Rong Zhou**
Google

**John Boaz Lee**
WPI

**Xiangnan Kong**
WPI

**Theodore L. Willke**
Intel Labs

**Hoda Eldardiry**
PARC

## Abstract

Random walks are at the heart of many existing deep learning algorithms for graph data. However, such algorithms have many limitations that arise from the use of random walks, *e.g.*, the features resulting from these methods are unable to transfer to new nodes and graphs as they are tied to node identity. In this work, we introduce the notion of *attributed random walks* which serves as a basis for generalizing existing methods such as DeepWalk, node2vec, and many others that leverage random walks. Our proposed framework enables these methods to be more widely applicable for both transductive and inductive learning as well as for use on graphs with attributes (if available). This is achieved by learning functions that generalize to new nodes and graphs. We show that our proposed framework is effective with an average AUC improvement of 16.1% while requiring on average 853 times less space than existing methods on a variety of graphs from several domains.

**Keywords** random walk · representation learning · inductive learning · deep learning · attributed graphs

## 1 Introduction

Graphs (networks) are ubiquitous [5, 9] and allow us to model entities (nodes) and the dependencies (edges) between them. Graph data is often observed directly in the natural world (*e.g.*, biological or social networks) [53] or constructed from non-relational data by deriving a metric space between entities and retaining only the most significant edges [21, 47, 58]. Learning a useful feature representation from graph data lies at the heart and success of many machine learning tasks such as node classification [33, 35], anomaly detection [5], link prediction [6], dynamic network analysis [37], community detection [36, 44], role discovery [46], visualization and sensemaking [1, 41], graph classification [25, 29], and network alignment [26].

Many existing techniques use *random walks* as a basis for learning features or estimating the parameters of a graph model for a downstream prediction task. Examples include recent node embedding methods such as DeepWalk [40],

node2vec [19], as well as graph-based deep learning algorithms. However, the simple random walk used by these methods is fundamentally tied to the *identity* of the node. This has three main disadvantages. First, these approaches are inherently transductive and do not generalize to unseen nodes and other graphs. Furthermore, they are unable to be used for graph-based transfer learning tasks such as across-network classification [15, 27], graph similarity [17, 57], and matching [39, 43]. Second, they are not space-efficient as a feature vector is learned for each node which is impractical for large graphs. Third, most of these approaches lack support for *attributed graphs*. This includes graphs with intrinsic (node) attributes such as age or gender as well as structural features such as higher-order subgraph counts, *e.g.*, number of 4-cliques each node participates.

To make these methods more generally applicable, we introduce a reinterpretation of the notion of random walk that is not tied to node identity and is instead based on a function $\Phi : \mathbf{x} \to w$ that maps a node attribute vector to a type. Using this reinterpretation as a basis we propose a framework that naturally generalizes many existing methods. This framework provides a number of important advantages. First, the learned features generalize to new nodes and across graphs and therefore can be used for transfer learning tasks such as across-network link prediction and classification. Our proposed approach supports both transductive and inductive learning. Second, the generalized approach is inherently space-efficient since embeddings are learned for types (as opposed to nodes) and therefore requires significantly less space than existing methods. Third, the generalized approach supports learning from attributed graphs. However, we stress that the approach does not require graphs with input attributes since these can be derived from the graph structure. Furthermore, our approach is shown to be effective with an average improvement of 16.1% in AUC while requiring on average 853x less space than existing methods on a variety of graphs.

**Contributions**: This paper proposes a framework for generalizing many existing algorithms making them more broadly applicable for attributed graphs and inductive learning. It has the following key properties:

- **Space-efficient**: It requires on average 853x less space

than a number of existing methods.

- **Accurate**: It is accurate with an average improvement of 16.1% across a variety of graphs from several domains.
- **Inductive**: It is an inductive learning approach that is able to learn embeddings for new nodes and graphs.
- **Attributed**: It naturally supports graphs with attributes (if available) and serves as a foundation for generalizing existing methods for use on attributed graphs.

## 2 Framework

In this section, we formally introduce the attributed random walk framework which can serve as a basis for generalizing many existing deep graph models and embedding methods. The framework consists of two general components:

**C1**. (**Function Mapping Nodes to Types**): A function $\Phi$ that maps nodes to types based on a $N \times K$ matrix $\mathbf{X}$ of attributes. The function $\Phi$ can be defined or learned automatically. Note $\mathbf{X}$ may be given as input or computed based on the structure of the graph. For more details, see Section 2.1.

**C2**. (**Attributed Random Walks**): Generate *attributed random walks* based on the function $\Phi$ mapping nodes to types (Section 2.2). Informally, an attributed walk is simply the node types that occur during a walk.

A summary of the notation is provided in Table 1. To avoid abuse of notation, this paper presents the framework for nodes. However, the approach is applicable to both nodes and edges.

### 2.1 Function Mapping Nodes to Types

Given an (un)directed graph $G = (V, E)$ and an $N \times K$ matrix $\mathbf{X}$ of attributes (if available), the first component of the framework maps the $N$ nodes to a set $W = \{w_1, ..., w_M\}$ of $M$ types where $1 \leq M \leq N$ and $M$ is often much smaller than $N$, *i.e.*, $M \ll N$. More formally,

$$\Phi : \mathbf{x} \rightarrow w \qquad (1)$$

Table 1: Summary of notation.

| | |
|---|---|
| $G$ | (un)directed (attributed) graph |
| $N$ | number of nodes $|V|$ |
| $M$ | number of unique types |
| $K$ | number of attributes used for deriving types |
| $D$ | number of dimensions to use for embedding |
| $R$ | number of walks for each node |
| $L$ | length of a random walk |
| $\Gamma_i$ | neighbors (adjacent nodes) of $v_i$ |
| $\mathcal{S}$ | set of (attributed) random walks |
| $\Phi(\cdot)$ | a function $\Phi : \mathbf{x} \rightarrow w$ that maps an attribute vector $\mathbf{x} \in \mathbb{R}^K$ to a corresponding type |
| $W$ | set of types where $w \in W$ is a type assigned to one or more nodes |
| $\mathbf{X}$ | an $N \times K$ attribute matrix where the rows represent nodes and the columns represent attributes |
| $\mathbf{x}_i$ | a $K$-dimensional attribute vector for node $v_i$ |
| $\bar{\mathbf{x}}_j$ | a $N$-dimensional vector for attribute $j$ (column of $\mathbf{X}$) |
| $\alpha$ | transformation hyperparameter |
| $\mathbf{Z}$ | an $M \times D$ matrix of type embeddings |
| $\mathbf{z}_k$ | a $D$-dimensional embedding for type $w_k$ |

where $\Phi$ is a function mapping nodes to *types* based on the $N \times K$ attribute matrix $\mathbf{X}$. The function $\Phi$ can be learned automatically or defined manually by the user. The framework is general and flexible for use with any arbitrary function $\Phi$ that maps nodes to types based on an $N \times K$ attribute matrix $\mathbf{X}$. For graphs that are not attributed, we simply derive $\mathbf{X}$ by extracting a set of graph features based on the structure of the graph (*e.g.*, $\mathbf{X}$ may represent the graphlet features shown in Figure 2).[1] The set of types $W$ is defined as:

$$W = \Phi(\mathbf{x}_1) \cup \Phi(\mathbf{x}_2) \cup \cdots \cup \Phi(\mathbf{x}_N) \qquad (2)$$

where $\Phi$ is a function that maps the set of $N = |V|$ nodes to $M = |W|$ types such that $1 \leq M \leq N$. In node2vec and other embedding and deep graph algorithms, notice that each node is mapped to a unique identifier that identifies the specific node which is then used in the random walk. However, in the attributed random walk framework two or more nodes may map to the same type.

There are two general classes of functions for mapping nodes to types. The first class of functions are simple functions taking the form:

$$\Phi(\mathbf{x}) = x_1 \circ x_2 \circ \cdots \circ x_K \qquad (3)$$

where $\mathbf{x}$ is an attribute vector:

$$\mathbf{x} = \begin{bmatrix} x_1 & x_2 & \cdots & x_K \end{bmatrix} \qquad (4)$$

and $\circ$ is a binary operator such as concatenation, sum, among others. The second class of functions are learned by solving an objective function. This includes functions based on a low-rank factorization of the $N \times K$ matrix $\mathbf{X}$ having the form $\mathbf{X} \approx f\langle \mathbf{U}\mathbf{V}^T \rangle$ with factor matrices $\mathbf{U} \in \mathbb{R}^{N \times F}$ and $\mathbf{V} \in \mathbb{R}^{K \times F}$ where $F > 0$ is the rank and $f$ is a linear or non-linear function. More formally,

$$\arg \min_{\mathbf{U}, \mathbf{V} \in \mathcal{C}} \left[ \mathcal{D}(\mathbf{X}, f\langle \mathbf{U}\mathbf{V}^T \rangle) + \mathcal{R}(\mathbf{U}, \mathbf{V}) \right] \qquad (5)$$

where $\mathcal{D}$ is the loss, $\mathcal{C}$ is constraints (*e.g.*, non-negativity constraints $\mathbf{U} \geq 0, \mathbf{V} \geq 0$), and $\mathcal{R}(\mathbf{U}, \mathbf{V})$ is a regularization penalty. Let $V_k$ denote the set of nodes mapped to type $w_k \in W$. We partition $\mathbf{U} \in \mathbb{R}^{N \times F}$ into $M$ disjoint sets of nodes (for each of the $M$ types) $V_1, \ldots, V_M$ such that $V = V_1 \cup \ldots \cup V_M$ by solving the k-means objective:

$$\min_{\{V_j\}_{j=1}^M} \sum_{j=1}^M \sum_{\mathbf{u}_i \in V_j} \|\mathbf{u}_i - \mathbf{c}_j\|^2, \text{ where } \mathbf{c}_j = \frac{\sum_{\mathbf{u}_i \in V_j} \mathbf{u}_i}{|V_j|} \quad (6)$$

Alternatively, the nodes can be mapped to types directly by setting $F = M$ and using $\mathbf{U}$ as follows:

$$\Phi(\mathbf{x}_i) = \arg \max_{w_k \in W} U_{ik}, \quad \forall i = 1, \ldots, N \qquad (7)$$

For other possibilities, we refer the readers to a recent survey by Xu *et al.* [54]. Notice that the approaches in Eq. (5)-(7) can be used for inductive learning as follows: Given a new

---

[1]Graphlets (and orbits) can be derived using exact algorithms [4] or accurately estimated using fast parallel approximation methods with provable error bounds [3, 49].
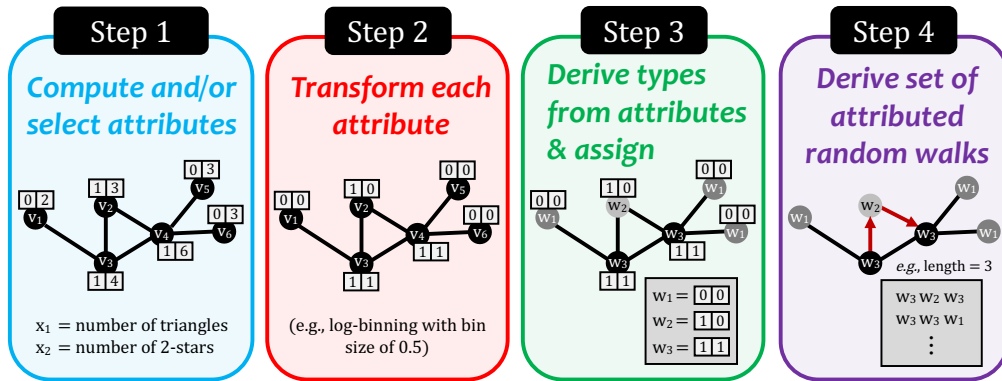
Figure 1: A simple illustrative toy example. This toy example shows only one potential instantiation of the general framework. In particular, Step 1 computes the number of triangles and 2-stars that each node participates whereas Step 2 transforms each individual attribute using logarithmic binning. Next, Step 3 derives types using a simple function $\Phi$ representing a concatenation of each node's attribute values resulting in $M = 3$ types $\{w_1, w_2, w_3\}$ among the $N = 6$ nodes. Finally, Step 4 derives a set of attributed random walks (type sequences) which are then used to learn embeddings. See Section 2.3 for further discussion.

graph $G' = (V', E')$ (or node $v_{N+1}$), we compute the attribute matrix $\mathbf{X}'$ (or attribute vector $\mathbf{x}_{N+1}$) using the same $K$ attributes and then use $\mathbf{X}'$ and $\mathbf{V}$ from Eq. (5) to estimate $\mathbf{U}'$ directly. More formally, given $\mathbf{X}'$ and $\mathbf{V}$, we find $\mathbf{U}'$ by solving:

$$\underset{\mathbf{U}' \in \mathcal{C}}{\operatorname{argmin}} \left[ \mathcal{D}\big(\mathbf{X}', f\langle \mathbf{U}'\mathbf{V}^T \rangle\big) + \mathcal{R}(\mathbf{U}') \right] \qquad (8)$$

Observe that $\mathbf{V}$ gives the mapping from the latent feature space to the input attribute space and represents how each of the $K$ attributes map to the latent features and therefore $\mathbf{V}$ is not tied to particular nodes but more generally to the (structural) attributes. Many of these methods avoid the issue of selecting an appropriate subset of attributes to use apriori and are more robust to the selection of different sets of attributes compared to simple functions such as $\Phi(\mathbf{x}) = x_1 \circ x_2 \circ \cdots \circ x_K$. As an aside, it is straightforward to formulate a function learning problem to automate the choice of function $\Phi$.

The framework naturally supports both attributed and non-attributed graphs as the attribute matrix $\mathbf{X}$ can be derived by extracting a set of graph features based on the structure of the graph (*e.g.*, $\mathbf{X}$ may represent the graphlet features shown in Figure 2) or learned automatically using an inductive feature learning approach such as DeepGL [48].

## 2.2 Attributed Random Walks

Most existing deep learning models and embedding algorithms for graphs use simple random walks based on node identity [19, 40]. The features learned from these methods are fundamentally tied to the identity of a node and therefore are unable to generalize to new nodes as well as for graph-based transfer learning tasks or across-network learning, among other limitations. In this section, we introduce the notion of attributed random walks and demonstrate how they serve as a basis for generalizing many existing deep graph models and embedding algorithms.

Previous work uses random walks based on the following traditional definition of a walk:

**Definition 1 (Walk)** *A walk $S$ of length $L$ is defined as a sequence of indices $i_1, i_2, \ldots, i_{L+1}$ such that $(v_{i_t}, v_{i_{t+1}}) \in E$ for all $1 \le t \le L$.*

Note that nodes in a walk $S$ are not necessarily distinct. We replace the traditional notion of walk (Definition 1) with the more appropriate and useful notion of *attributed walk* defined as:

**Definition 2 (Attributed walk)** *Let $\mathbf{x}_i$ be a $K$-dimensional vector for node $v_i$. An attributed walk $S$ of length $L$ is defined as a sequence of adjacent node types*

$$\Phi(\mathbf{x}_{i_1}), \Phi(\mathbf{x}_{i_2}), \ldots, \Phi(\mathbf{x}_{i_{L+1}}) \qquad (9)$$

*associated with a sequence of indices $i_1, i_2, \ldots, i_{L+1}$ such that $(v_{i_t}, v_{i_{t+1}}) \in E$ for all $1 \le t \le L$ and $\Phi : \mathbf{x} \to w$ is a function that maps the input vector $\mathbf{x}$ of a node to a corresponding type $\Phi(\mathbf{x})$.*

The type sequence $\Phi(\mathbf{x}_{i_1}), \Phi(\mathbf{x}_{i_2}), \ldots, \Phi(\mathbf{x}_{i_{L+1}})$ is simply the node types that occur during a walk. The framework allows both uniform and non-uniform attributed random walks. It is also straightforward to bias the attributed random walk in an arbitrary fashion. The notion of attributed walk can be extended in various ways (*e.g.*, for edge types).

Given an arbitrary embedding method or deep graph model $\mathcal{A}$ that uses simple random walks based on node identity (*e.g.*, node2vec, DeepWalk, among others), we generalize $\mathcal{A}$ by replacing the traditional notion of random walk with the proposed notion of *attributed random walks*. We will show that existing methods are actually a special case of the proposed framework when $\Phi(\mathbf{x}_i)$ uniquely identifies node $v_i$. Suppose the graph $G$ has $N$ nodes mapped to $M$ types using an arbitrary function $\Phi$ (Section 2.1). If we select $\Phi$ such that $M \to N$ then we recover the original method $\mathcal{A}$ as a special case of the framework. More intuitively, each node $v_i \in V$ must be assigned to a unique type $\Phi(\mathbf{x}_i) = w_i$ that uniquely identifies it from any other node $v_j \in V$ since $|W| = |V|$ where $W = \Phi(\mathbf{x}_1) \cup \cdots \cup \Phi(\mathbf{x}_N)$. Observe that $M \to N$ is actually an extreme case of the framework and corresponds exactly to the original method $\mathcal{A}$.

There are three key advantages that arise when existing methods are generalized by the proposed framework. First, the generalized method naturally supports attributed graphs. Second, the features learned generalize to new nodes and across graphs and therefore are naturally inductive and able to be used for transfer learning tasks. Finally, the learned embeddings are significantly more space-efficient with a space complexity of only $\mathcal{O}(MD)$ compared to previous methods that require $\mathcal{O}(ND)$ space where $M \ll N$ ($M$ is much smaller than $N$) and $D$ is the embedding dimensionality.

## 2.3 Illustrative Example

We emphasize that the proposed framework is general and flexible; and the example discussed in this section represents only one such simple instantiation of the framework. An overview of the key steps involved in this specific instantiation are shown in Figure 1 and summarized succinctly below as follows:

**Step 1.** **Compute and/or Select Attributes**: Given an (un)directed graph $G$, the first step is to compute a set of attributes using the graph structure (*e.g.*, graphlets).

**Step 2.** **Transform Attributes**: Next, we transform each attribute vector. The goal is to map the values of each individual attribute vector $\bar{\mathbf{x}}_j$ to a smaller set of values (via quantization or discretization algorithm). In Figure 1, each attribute $\bar{\mathbf{x}}_j$ is transformed using logarithmic binning.[2] For convenience, each initial $\bar{\mathbf{x}}_j$ is replaced by the transformed attribute.

**Step 3.** **Derive Types from Attributes & Assign**: Now, the set of *types* are derived as $W = \Phi(\mathbf{x}_1) \cup \Phi(\mathbf{x}_2) \cup \cdots \cup \Phi(\mathbf{x}_N)$ where $|W| \le N$ and $\Phi : \mathbf{x} \to w$ such that $\Phi(\mathbf{x}) = x_1 \circ x_2 \circ \cdots \circ x_K$ and $\circ$ is a binary operator such as concatenation. We also assign each node to its corresponding type. See Figure 1 for intuition.

**Step 4.** **Attributed Random Walks**: Finally, we generate a set of attributed random walks using the node types (as opposed to the unique node identifiers used traditionally). The attribute random walks represent sequences of node types which are then used to learn embeddings. For instance, the set of attributed random walks are used in Section 3 to learn the $M \times D$ type embedding matrix $\mathbf{Z}$.

From the perspective of the general framework, Steps 1-3 in Figure 1 correspond to the first main component that maps nodes to types via the function $\Phi$ while Step 4 corresponds to the last component that generates *attributed random walks* representing sequences of node types.

## 3 Node2Vec Generalization

We use the proposed framework in Section 2 to generalize node2vec [19].[3] This gives rise to a new generalized

---

[2]Logarithmic binning assigns the first $\alpha N$ nodes with smallest attribute value to 0 (where $0 < \alpha < 1$), then assigns the $\alpha$ fraction of remaining unassigned nodes with smallest value to 1, and so on.

[3]However, the framework could have been used to generalize any node embedding or deep graph model that leverages traditional random walks.

---

**Algorithm 1** Generalized Node2Vec

1: **procedure** GENERALIZEDNODE2VEC( $G = (V, \text{E})$ and possibly attributes $\mathbf{X}$, embedding dimensions $D$, walks per node $R$, walk length $L$, context (window) size $\omega$, return $p$, and in-out $q$, function $\Phi$ )
2:      Initialize set of *attributed walks* $\mathcal{S}$ to $\emptyset$
3:      Extract (graphlet) features if needed and append to $\mathbf{X}$
4:      Transform each attribute in $\mathbf{X}$ (*e.g.*, using logarithmic binning)
5:      Precompute transition probabilities $\pi$ using $G$, $p$, and $q$
6:      $G' = (V, E, \pi)$
7:      **parallel for** $j = 1, 2, ..., R$ **do**      ▷ walks per node
8:          Set $\Pi$ to be a random permutation of the nodes $V$
9:          **for each** $v \in \Pi$ in order **do**
10:              $S = \text{ATTRIBUTEDWALK}(G', \mathbf{X}, v, \Phi, L)$
11:              Add the *attributed walk* $S$ to $\mathcal{S}$
12:          **end for**
13:      **end parallel**
14:      $\mathbf{Z} = \text{STOCHASTICGRADIENTDESCENT}(\omega, D, \mathcal{S})$      ▷ **parallel**
15:      **return** the learned *type* embeddings $\mathbf{Z}$

---

**Algorithm 2** Attributed Random Walk

1: **procedure** ATTRIBUTEDWALK($G'$, $\mathbf{X}$, start node $s$, function $\Phi$, $L$)
2:      Initialize attributed walk $S$ to $\left[\Phi(\mathbf{x}_s)\right]$
3:      Set $i = s$      ▷ current node
4:      **for** $\ell = 1$ **to** $L$ **do**
5:          $\Gamma_i = $ Set of the neighbors for node $i$
6:          $j = \text{ALIASSAMPLE}(\Gamma_i, \pi)$      ▷ select node $j \in \Gamma_i$
7:          Append $\Phi(\mathbf{x}_j)$ to $S$
8:          Set $i$ to be the current node $j$
9:      **end for**
10:      **return** attributed walk $S$ of length $L$ rooted at node $s$

---

node2vec algorithm with the following advantages: (a) naturally supports attributed graphs, (b) learns features that generalize for new nodes as well as for extraction on another arbitrary graph and thus able to be used for transfer learning tasks, and (c) space-efficient by learning an embedding for each type as opposed to each node. The generalized node2vec algorithm is shown in Alg. 1. In particular, we replace the notion of random walk in node2vec with the more appropriate notion of *attributed random walk* (Line 10-11). Nodes are mapped to types using an arbitrary function $\Phi(\mathbf{x})$ which can be defined or learned; see Section 2.1 for further details. As an aside, Alg. 1 assumes the function $\Phi$ is defined or learned apriori. Note that if no attribute matrix $\mathbf{X}$ is given as input then we derive features based on the graph structure (Line 3). However, even if (intrinsic/self- or structural) attributes are given as input we may also choose to derive additional structural features such as graphlets [2] and append these to $\mathbf{X}$.

The attributed random walk routine for the generalized node2vec algorithm is shown in Alg. 2. Given a start node $s$, we first obtain the *type* of the start node $s$ given by the function $\Phi(\mathbf{x}_s)$ and initialize an attributed walk $S$ (stored as

a list[4]) to be $S = \big[ \Phi(\mathbf{x}_s) \big]$ (Line 2). Next, Line 3 sets the current node $i$ to be the start node $s$. The attributed walk $S$ of length $L$ rooted at the start node $s$ is obtained in Lines 4-9. In particular, Line 5 gets the neighbors of the current node $i$ which is used in Line 6 to sample a new node $j \in \Gamma_i$. The type $\Phi(\mathbf{x}_j)$ of node $j$ is appended to $S$ (*e.g.*, if $S$ is stored using an ordered list then we append $\Phi(\mathbf{x}_j)$ to the back of $S$ to preserve order). Afterwards, Line 8 sets the current node $i$ to be $j$ and Lines 4-9 are repeated for $L$ steps. Finally, an attributed random walk $S$ of length $L$ starting at node $s$ is returned in Line 10. One can also leverage Alg. 2 (or an adapted variant of it) for generalizing other deep graph models and node embedding methods. It is straightforward to see that Alg. 1 is a generalization of node2vec since if we allow $M \to N$ then we recover the original node2vec algorithm as a special case of the framework. Furthermore, DeepWalk [40] is also a special case where $M \to N$, $p = 1$, and $q = 1$.

**Implementation details:** In our implementation, we evaluate $\Phi(\mathbf{x}_i) = x_1 \circ x_2 \circ \cdots \circ x_K$ only once for each node $i \in V$ and store each $\Phi(\mathbf{x}_i)$ in a hash table giving $o(1)$ time lookup. This allows us to derive $\Phi(\mathbf{x}_s)$ and $\Phi(\mathbf{x}_j)$ in Line 2 and 7 of Alg. 2 in $o(1)$ constant time. To construct the hash table it takes only $\mathcal{O}(NK)$ time to obtain $\Phi(\mathbf{x}_i), \forall i = 1, ..., N$ given $\mathbf{X}$ and only $\mathcal{O}(N)$ space to store them efficiently using a hash table.

# 4 Experiments

In this section, we investigate the effectiveness of the proposed framework using a variety of graphs from several domains. In these experiments, we use the node2vec generalization given in Section 3.
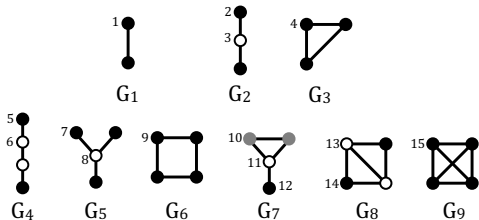
Figure 2: Summary of the 9 node graphlets and 15 orbits (graphlet automorphisms) with 2-4 nodes.

## 4.1 Experimental Setup

Unless otherwise mentioned, all experiments use logarithmic binning (defined in Section 2.3) with $\alpha = 0.5$. In these experiments, we use a simple function $\Phi(\mathbf{x})$ that represents a concatenation of the attribute values in the node attribute vector $\mathbf{x}$. In particular, $\Phi(\mathbf{x})$ is defined as $\Phi(\mathbf{x}) = x_1 \circ \cdots \circ x_K$ where $\circ$ is the concatenation operator. For these experiments, we searched over 10 subsets of the 9 graphlet attributes shown in Figure 2.[5] However, the approach trivially

---

[4]A list or other efficient data structure may be used for storing a sequence of items (nodes/edges) with potential duplicates

[5]Note we also investigated the 15 graphlet orbits (as opposed to the 9 graphlet attributes) and found similar results.

handles other categories of attributes including structural attributes, intrinsic/self-attributes (such as age or other non-relational attributes), and relational features derived using the graph $G$ along with existing structural or self-attributes. We evaluate the *generalized node2vec* approach presented in Section 3 that leverages the attributed random walk framework (Section 2) against a number of popular methods including: node2vec[6], DeepWalk [40], and LINE [52]. For our approach and node2vec, we use the same hyperparameters ($D = 128$, $R = 10$, $L = 80$) and grid search over $p, q \in \{0.25, 0.50, 1, 2, 4\}$ as mentioned in [19]. We use logistic regression (LR) with an L2 penalty. The model is selected using 10-fold cross-validation on $10\%$ of the labeled data. Experiments are repeated for 10 random seed initializations. All results are statistically significant with p-value $< 0.01$. Unless otherwise mentioned, we use AUC to evaluate the models. Data has been made available at NetworkRepository [45].[7]
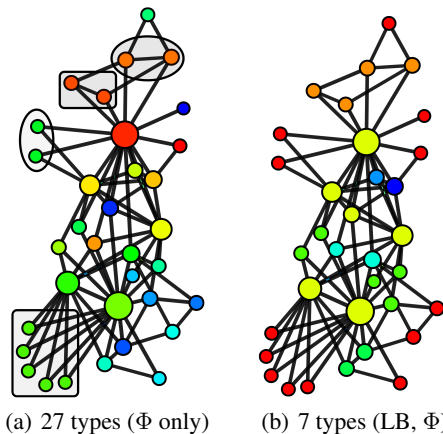
(a) 27 types (Φ only)          (b) 7 types (LB, Φ)

Figure 3: Using only the node types are useful for revealing interesting insights into the graph data including nodes that are structurally similar. Node types alone reveal nodes that are structurally similar (karate network). Node color encodes type. Only the number of triangles and wedges are used by Φ. In (a) nodes with identical types are grouped; (b) uses log binning (LB) with $\alpha = 0.5$. See text for further discussion.

## 4.2 Node Mapping Experiments

This section investigates the intermediate representation defined by an arbitrary function Φ that maps nodes to a set of types. Recall that for simplicity Φ is defined as a concatenation of the attribute values for a given node and the attributes correspond to the graphlet counts in Figure 2. Strikingly, we observe in Figure 3 that even these simple mappings alone are useful and effective for understanding the important structural and behavioral characteristics of nodes, that is before even generating attributed random walks and learning an embedding based on them. In particular, nodes assigned to the same type in Figure 3(a) obey a strict notion

---

[6]https://github.com/aditya-grover/node2vec
[7]http://networkrepository.com/

of node equivalence on a feature representation:

$$[\forall j, 1 \le j \le K : f_j(u) = f_j(v)] \Rightarrow u \equiv v \qquad (10)$$

where $f_j$ is a feature. Eq. 10 is strict since two nodes belong to the same type iff they have identical feature vectors. However, Figure 3(b) captures a more relaxed notion of structural equivalance called structural similarity [46]. This result is surprising since no embedding has been learned yet, only the mapping into the intermediate representation (types). Hence, it validates the intermediate representation of mapping nodes to types while also demonstrating the effectiveness of it for grouping structurally similar nodes.

## 4.3 Comparison

This section compares the proposed approach to other embedding methods for link prediction. Given a partially observed graph $G$ with a fraction of missing edges, the link prediction task is to predict these missing edges. We generate a labeled dataset of edges as done in [19]. Positive examples are obtained by removing $50\%$ of edges randomly, whereas *negative examples* are generated by randomly sampling an equal number of node pairs that are not connected with an edge, *i.e.*, each node pair $(i, j) \notin E$. For each method, we learn features using the remaining graph that consists of only positive examples. Using the feature representations from each method, we then learn a model to predict whether a given edge in the test set exists in $E$ or not. Notice that node embedding methods such as DeepWalk and node2vec require that each node in $G$ appear in at least one edge in the training graph (*i.e.*, the graph remains connected), otherwise these methods are unable to derive features for such nodes. This is a significant limitation that prohibits their use in many real-world applications.

For comparison, we use the same set of binary operators [19] to construct features for the edges *indirectly* using the learned node representations. Moreover, the AUC scores
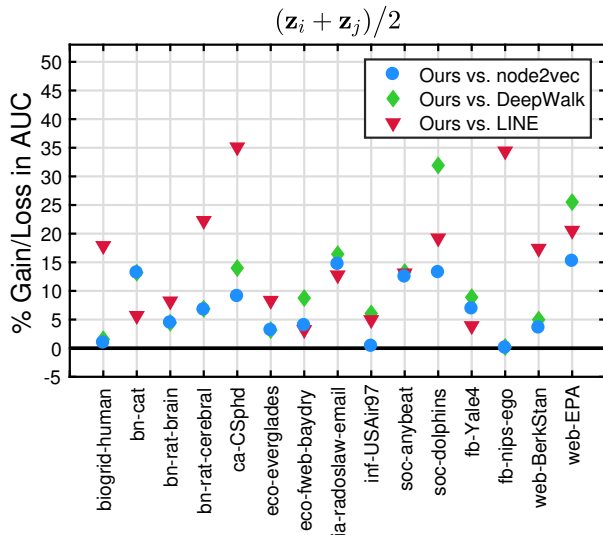


Figure 4: AUC gain over other methods for link prediction bootstrapped using $(\mathbf{z}_i + \mathbf{z}_j)/2$. The average gain over all methods and across all graphs is $10.5\%$.
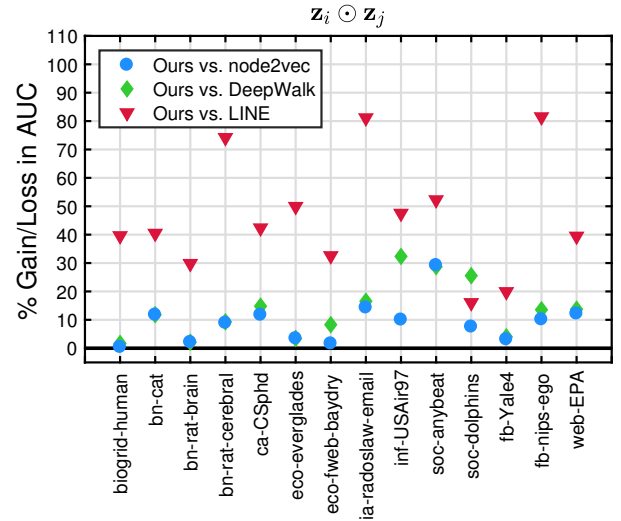


Figure 5: AUC gain of our approach over the other methods for link prediction bootstrapped using Hadamard. The average gain over all methods and across all graphs is $21.7\%$.

from our method are all significantly better than the other methods at $p < 0.01$. We summarize the gain/loss in predictive performance over the other methods in Figure 4 and 5. In all cases, our method achieves better predictive performance over the other methods across a wide variety of graphs with different characteristics. Overall, the mean and product binary operators give the best results with an average gain in predictive performance (over all graphs) of $10.5\%$ and $21.7\%$, respectively.

## 4.4 Space-efficient Embeddings

We now investigate the space-efficiency of the learned embeddings from the proposed framework and intermediate representation. Observe that any embedding method that implements the proposed attributed random walk framework (and intermediate representation) learns an embedding for each distinct node type $w \in W$. In the worst case, an embedding is learned for each of the $N$ nodes in the graph and we recover the original method as a special case. In general, the best embedding most often lies between such extremes and therefore the embedding learned from a method implementing the framework is often orders of magnitude smaller in size since $M \ll N$ where $M = |W|$ and $N = |V|$.

Given an attribute vector $\mathbf{x}$ of graphlet counts (Figure 2) for an arbitrary node in $G$, we derive embeddings using each of the following:

$$\Phi(\mathbf{x}_i = [\,x_2\ x_3\,]), \quad \text{for } i = 1, ..., N \qquad (11)$$

$$\Phi(\mathbf{x}_i = [\,x_2\ x_3\ x_4\ x_6\ x_9\,]), \quad \text{for } i = 1, ..., N \qquad (12)$$

$$\Phi(\mathbf{x}_i = [\,x_2\ x_3\ x_4\ x_5\ \cdots\ x_9\,]), \quad \text{for } i = 1, ..., N \qquad (13)$$

$$\Phi(\mathbf{x}_i = [\,x_1\ x_2\ x_3\ x_4\ x_5\ \cdots\ x_9\,]), \quad \text{for } i = 1, ..., N \qquad (14)$$

where $\Phi(\cdot)$ is a function that maps $\mathbf{x}_i$ to a type $w \in W$. In these experiments, we use logarithmic binning (applied to each $N$-dimensional graphlet feature) with $\alpha = 0.5$ and use $\Phi$ defined as the concatenation of the logarithmically binned attribute values. Embeddings are learned using the different
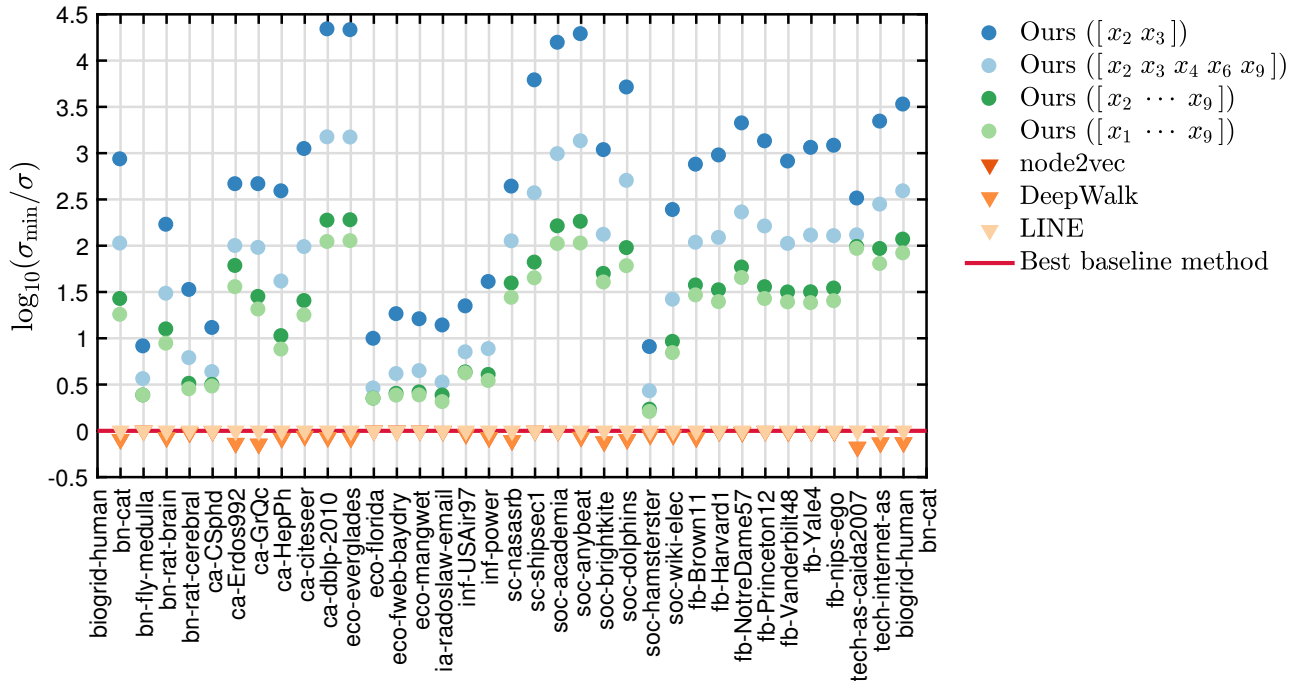
Figure 6: The space requirements of the proposed approach are orders of magnitude less than existing methods. Note $\log_{10}(\sigma_{\min}/\sigma)$ (y-axis) is the relative "gain/loss" (in space) over the best baseline method where $\sigma$ is the size (bytes) of an embedding and $\sigma_{\min}$ is the embedding with the smallest size among the baseline methods (node2vec, DeepWalk, LINE); and $[\mathrm{x}_2\ \mathrm{x}_3], \cdots, [\mathrm{x}_1\ \cdots\ \mathrm{x}_9]$ is the attribute sets in Eq. (11)-(14) used as input to $\Phi$ for mapping nodes to types. These represent variations of our method. The baseline method with the minimum space ($\sigma_{\min}$) is shown at 0 and thus methods that use more space lie below the horizontal line.

subsets of attributes in Eq. (11)-(14). For instance, Eq. (11) indicates that node types are derived using the (logarithmic binned) number of 2-stars $\mathrm{x}_2$ and triangles $\mathrm{x}_3$ that a given node participates (Figure 2). To evaluate the space-efficiency of the various methods, we measure the space (in bytes) required to store the embedding learned by each method. In Figure 6, we summarize the reduction in space from our approach compared to the other methods. In all cases, the embeddings learned from our approach require significantly less space and thus more space-efficient. Specifically, the embeddings from our approach require on average 853 times less space than the best method averaged across all graphs. In addition, Table 2 provides the number of types derived when using the various attribute combinations.

## 5 Related Work

Recent embedding techniques for graphs have largely been based on the popular skip-gram model [11, 34] originally introduced for learning vector representations of words in the natural language processing (NLP) domain. In particular, DeepWalk [40] used this approach to embed the nodes such that the co-occurrence frequencies of pairs in short random walks are preserved. More recently, node2vec [19] introduced hyperparameters to DeepWalk that tune the depth and breadth of the random walks. These approaches are becoming increasingly popular and have been shown to outperform a number of existing methods. These methods are all

Table 2: Comparing the number of unique types.

| Graph | $N$ | $M$ | | | |
| --- | --- | --- | --- | --- | --- |
| | | Eq. (11) | Eq. (12) | Eq. (13) | Eq. (14) |
| biogrid–human | 9527 | 9 | 73 | 290 | 429 |
| ca–citeseer | 227320 | 9 | 132 | 1049 | 1783 |
| ca–dblp–2010 | 226413 | 9 | 130 | 1018 | 1706 |
| fb–Harvard1 | 15126 | 7 | 64 | 253 | 328 |
| fb–NotreDame57 | 12155 | 9 | 75 | 340 | 454 |
| sc–shipsec1 | 140385 | 9 | 144 | 866 | 1346 |
| soc–academia | 200169 | 9 | 129 | 954 | 1634 |
| soc–brightkite | 56739 | 9 | 92 | 489 | 768 |
| tech–internet–as | 40164 | 9 | 78 | 259 | 365 |

based on simple random walks and thus are well-suited for generalization using the *attributed random walk framework*.

On the other hand, methods such as SkipGraph [29] make use of simple random walks to learn embeddings for entire graphs (as opposed to individual nodes). These methods can be used for graph-level tasks such as graph classification and clustering. Since these methods are still based on simple random walks, it is straightforward to generalize them using our proposed framework.

While most work has focused on transductive (within-network) learning, there has been some recent work on graph-based inductive approaches [20, 48, 56]. Yang *et al.* proposed an inductive approach called Planetoid [56]. How-

ever, Planetoid is an embedding-based approach for semi-supervised learning and does not use any structural features. Rossi *et al.* proposed an inductive approach for (attributed) networks called DeepGL that learns (inductive) relational functions representing compositions of one or more operators applied to an initial set of graph features [48]. More recently, Hamilton *et al.* [20] proposed a similar approach that also aggregates features from node neighborhoods. However, these approaches are not based on random-walks.

Heterogeneous networks [51] have also been recently considered [10, 14] as well as attributed networks [22, 23]. In particular, Huang *et al.* proposed an approach for attributed networks with labels [23] whereas Yang *et al.* used text features to learn node representations [55]. Liang *et al.* proposed a semi-supervised approach for networks with outliers [30]. Bojchevski *et al.* proposed an unsupervised rank-based approach [8]. More recently, Coley *et al.* introduced a convolutional approach for attributed molecular graphs that learns graph embeddings [13] as opposed to node embeddings. However, most of these approaches are not inductive nor space-efficient.

Our work is also related to uniform and non-uniform random walks on graphs [12, 32]. Random walks are at the heart of many important applications such as ranking [38], community detection [36, 42], recommendation [7], link prediction [31], influence modeling [24], search engines [28], image segmentation [18], routing in wireless sensor networks [50], and time-series forecasting [16]. These applications and techniques may also benefit from the notion of attributed random walks.

## 6 Conclusion

To make existing methods more generally applicable, this work proposed a flexible framework based on the notion of attributed random walks. The framework serves as a basis for generalizing existing techniques (that are based on random walks) for use with attributed graphs, unseen nodes, graph-based transfer learning tasks, and allowing significantly larger graphs due to the inherent space-efficiency of the approach. Instead of learning individual embeddings for each node, our approach learns embeddings for each type based on functions that map feature vectors to types. This allows for both inductive and transductive learning. Furthermore, the framework was shown to have the following desired properties: space-efficient, accurate, inductive, and able to support graphs with attributes (if available). Finally, the approach is guaranteed to perform at least as well as the original methods since they are recovered as a special case in the framework.

## References

[1] N. K. Ahmed and R. A. Rossi. Interactive visual graph analytics on the web. In *ICWSM*, 2015.

[2] N. K. Ahmed, J. Neville, R. A. Rossi, and N. Duffield. Efficient graphlet counting for large networks. In *ICDM*, page 10, 2015.

[3] N. K. Ahmed, T. L. Willke, and R. A. Rossi. Estimation of Local Subgraph Counts. In *BigData*, pages 586–595, 2016.

[4] N. K. Ahmed, J. Neville, R. A. Rossi, N. G. Duffield, and T. L. Willke. Graphlet Decomposition: Framework, Algorithms, and Applications. *Knowledge and Information Systems*, 50 (3):689–722, Mar 2017.

[5] L. Akoglu, H. Tong, and D. Koutra. Graph based anomaly detection and description: a survey. *DMKD*, 29(3):626–688, 2015.

[6] M. Al Hasan and M. J. Zaki. A survey of link prediction in social networks. In *Social Network Data Analytics*, pages 243–275. 2011.

[7] T. Bogers. Movie recommendation using random walks over the contextual graph. In *Context-Aware Recommender Systems*, 2010.

[8] A. Bojchevski and S. Günnemann. Deep gaussian embedding of attributed graphs: Unsupervised inductive learning via ranking. *arXiv:1707.03815*, 2017.

[9] S. Borgatti and M. Everett. Notions of position in social network analysis. *Sociological methodology*, 22(1):1–35, 1992.

[10] S. Chang, W. Han, J. Tang, G.-J. Qi, C. C. Aggarwal, and T. S. Huang. Heterogeneous network embedding via deep architectures. In *SIGKDD*, pages 119–128, 2015.

[11] W. Cheng, C. Greaves, and M. Warren. From n-gram to skip-gram to concgram. *Int. J. of Corp. Linguistics*, 11(4):411–433, 2006.

[12] F. Chung. Random walks and local cuts in graphs. *Linear Algebra and its applications*, 423(1):22–32, 2007.

[13] C. W. Coley, R. Barzilay, W. H. Green, T. S. Jaakkola, and K. F. Jensen. Convolutional embedding of attributed molecular graphs for physical property prediction. *J. Chem. Info. & Mod.*, 2017.

[14] Y. Dong, N. V. Chawla, and A. Swami. metapath2vec: Scalable representation learning for heterogeneous networks. In *SIGKDD*, pages 135–144, 2017.

[15] L. Getoor and B. Taskar, editors. *Intro. to SRL*. MIT Press, 2007.

[16] D. F. Gleich and R. A. Rossi. A dynamical system for pagerank with time-dependent teleportation. *Inter. Math.*, pages 188–217, 2014.

[17] T. E. Goldsmith and D. M. Davenport. Assessing structural similarity of graphs. 1990.

[18] L. Grady. Random walks for image segmentation. *TPAMI*, 28(11):1768–1783, 2006.

[19] A. Grover and J. Leskovec. node2vec: Scalable feature learning for networks. In *SIGKDD*, pages 855–864, 2016.

[20] W. L. Hamilton, R. Ying, and J. Leskovec. Inductive representation learning on large graphs. *arXiv:1706.02216*, 2017.

[21] M. Henaff, J. Bruna, and Y. LeCun. Deep convolutional networks on graph-structured data. *arXiv:1506.05163*, 2015.

[22] X. Huang, J. Li, and X. Hu. Accelerated attributed network embedding. In *SDM*, pages 633–641, 2017.

[23] X. Huang, J. Li, and X. Hu. Label informed attributed network embedding. In *WSDM*, pages 731–739, 2017.

[24] A. Java, P. Kolari, T. Finin, and T. Oates. Modeling the spread of influence on the blogosphere. In *WWW*, pages 22–26, 2006.

[25] S. Kearnes, K. McCloskey, M. Berndl, V. Pande, and P. Riley. Molecular graph convolutions: moving beyond fingerprints.

*Journal of computer-aided molecular design*, 30(8):595–608, 2016.

[26] M. Koyutürk, Y. Kim, U. Topkara, S. Subramaniam, W. Szpankowski, and A. Grama. Pairwise alignment of protein interaction networks. *JCB*, 13(2):182–199, 2006.

[27] A. Kuwadekar and J. Neville. Relational active learning for joint collective classification models. In *ICML*, pages 385–392, 2011.

[28] J.-L. Lassez, R. Rossi, and K. Jeev. Ranking links on the web: Search and surf engines. In *IEA/AIE*, pages 199–208, 2008.

[29] J. B. Lee and X. Kong. Skip-Graph: Learning graph embeddings with an encoder-decoder model. In *ICLR OpenReview*, 2017.

[30] J. Liang, P. Jacobs, and S. Parthasarathy. Seano: Semi-supervised embedding in attributed networks with outliers. *arXiv:1703.08100*, 2017.

[31] W. Liu and L. Lü. Link prediction based on local random walk. *Europhysics Letters*, 89:58007, 2010.

[32] L. Lovász. Random walks on graphs. *Combinatorics*, 2:1–46, 1993.

[33] L. K. McDowell, K. M. Gupta, and D. W. Aha. Cautious collective classification. *JMLR*, 10:2777–2836, 2009.

[34] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv:1301.3781*, 2013.

[35] J. Neville and D. Jensen. Iterative classification in relational data. In *AAAI SRL Workshop*, pages 13–20, 2000.

[36] A. Y. Ng, M. I. Jordan, and Y. Weiss. On spectral clustering: Analysis and an algorithm. In *NIPS*, pages 849–856, 2002.

[37] V. Nicosia, J. Tang, C. Mascolo, M. Musolesi, G. Russo, and V. Latora. Graph metrics for temporal networks. In *Temporal Networks*, pages 15–40. Springer, 2013.

[38] L. Page, S. Brin, R. Motwani, and T. Winograd. PageRank citation ranking: Bringing order to the web. *Stanford Tech. Report*, 1998.

[39] Y. Park, D. Reeves, V. Mulukutla, and B. Sundaravel. Fast malware classification by automated behavioral graph matching. In *Workshop on Cyber Security and Info. Intel. Res.*, page 45, 2010.

[40] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online learning of social representations. In *SIGKDD*, pages 701–710, 2014.

[41] R. Pienta, J. Abello, M. Kahng, and D. H. Chau. Scalable graph exploration and visualization: Sensemaking challenges and opportunities. In *BigComp*, pages 271–278, 2015.

[42] P. Pons and M. Latapy. Computing communities in large networks using random walks. *J. Graph Alg. Appl.*, 10(2):191–218, 2006.

[43] N. Pržulj. Biological network comparison using graphlet degree distribution. *Bioinformatics*, 23(2):e177–e183, 2007.

[44] F. Radicchi, C. Castellano, F. Cecconi, V. Loreto, and D. Parisi. Defining and identifying communities in networks. *Proceedings of the National Academy of Sciences*, 101(9): 2658–2663, 2004.

[45] R. A. Rossi and N. K. Ahmed. The network data repository with interactive graph analytics and visualization. In *AAAI*, pages 4292–4293, 2015. URL http://networkrepository.com.

[46] R. A. Rossi and N. K. Ahmed. Role discovery in networks. *TKDE*, 27(4):1112–1131, 2015.

[47] R. A. Rossi, L. K. McDowell, D. W. Aha, and J. Neville. Transforming graph data for statistical relational learning. *JAIR*, 45:363–441, 2012.

[48] R. A. Rossi, R. Zhou, and N. K. Ahmed. Deep feature learning for graphs. In *arXiv:1704.08829*, 2017.

[49] R. A. Rossi, R. Zhou, and N. K. Ahmed. Estimation of Graphlet Statistics. In *arXiv preprint*, pages 1–14, 2017.

[50] S. D. Servetto and G. Barrenechea. Constrained random walks on random graphs: routing algorithms for large scale wireless sensor networks. In *Wireless Sensor Networks & App.*, pages 12–21, 2002.

[51] C. Shi, X. Kong, Y. Huang, S. Y. Philip, and B. Wu. HeteSim: A General Framework for Relevance Measure in Heterogeneous Networks. *TKDE*, 26(10):2479–2492, 2014.

[52] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei. LINE: Large-scale Information Network Embedding. In *WWW*, pages 1067–1077, 2015.

[53] B. Viswanath, A. Mislove, M. Cha, and K. P. Gummadi. On the evolution of user interaction in Facebook. In *OSN*, pages 37–42, 2009.

[54] R. Xu and D. Wunsch. Survey of clustering algorithms. *IEEE Transactions on Neural Networks*, 16(3):645–678, 2005.

[55] C. Yang, Z. Liu, D. Zhao, M. Sun, and E. Y. Chang. Network representation learning with rich text information. In *IJCAI*, pages 2111–2117, 2015.

[56] Z. Yang, W. W. Cohen, and R. Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. *arXiv:1603.08861*, 2016.

[57] L. A. Zager and G. C. Verghese. Graph similarity scoring and matching. *Applied mathematics letters*, 21(1):86–94, 2008.

[58] X. Zhu, Z. Ghahramani, and J. D. Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*, pages 912–919, 2003.