# Attention Models in Graphs: A Survey

JOHN BOAZ LEE, WPI, USA
RYAN A. ROSSI, Adobe Research, USA
SUNGCHUL KIM, Adobe Research, USA
NESREEN K. AHMED, Intel Labs, USA
EUNYEE KOH, Adobe Research, USA

Graph-structured data arise naturally in many different application domains. By representing data as graphs, we can capture entities (*i.e.*, nodes) as well as their relationships (*i.e.*, edges) with each other. Many useful insights can be derived from graph-structured data as demonstrated by an ever-growing body of work focused on graph mining. However, in the real-world, graphs can be both large – with many complex patterns – and noisy which can pose a problem for effective graph mining. An effective way to deal with this issue is to incorporate "attention" into graph mining solutions. An attention mechanism allows a method to focus on task-relevant parts of the graph, helping it to make better decisions. In this work, we conduct a comprehensive and focused survey of the literature on the emerging field of graph attention models. We introduce three intuitive taxonomies to group existing work. These are based on problem setting (type of input and output), the type of attention mechanism used, and the task (*e.g.*, graph classification, link prediction, *etc.*). We motivate our taxonomies through detailed examples and use each to survey competing approaches from a unique standpoint. Finally, we highlight several challenges in the area and discuss promising directions for future work.

CCS Concepts: • **Computing methodologies** → **Artificial intelligence**; **Machine learning**; **Logical and relational learning**; • **Mathematics of computing** → **Graph algorithms**; *Combinatorics*; *Graph theory*; • **Information systems** → **Data mining**; • **Theory of computation** → **Graph algorithms analysis**; **Streaming, sublinear and near linear time algorithms**;

Additional Key Words and Phrases: Attention mechanism, graph attention, deep learning, graph attention survey

## 1 INTRODUCTION

Data which can be naturally modeled as graphs are found in a wide variety of domains including the world-wide web [Albert et al. 1999], bioinformatics [Pei et al. 2005], neuroscience [Lee et al. 2017], chemoinformatics [Duvenaud et al. 2015], social networks [Backstrom and Leskovec 2011], scientific citation and collaboration [Liu et al. 2017], urban computing [Zheng et al. 2014], recommender

Authors' addresses: John Boaz Lee, WPI, MA, USA, jtlee@wpi.edu; Ryan A. Rossi, Adobe Research, San Jose, CA, USA, rrossi@adobe.com; Sungchul Kim, Adobe Research, San Jose, CA, USA, sukim@adobe.com; Nesreen K. Ahmed, Intel Labs, Santa Clara, CA, USA, nesreen.k.ahmed@intel.com; Eunyee Koh, Adobe Research, San Jose, CA, USA, eunyee@adobe.com.
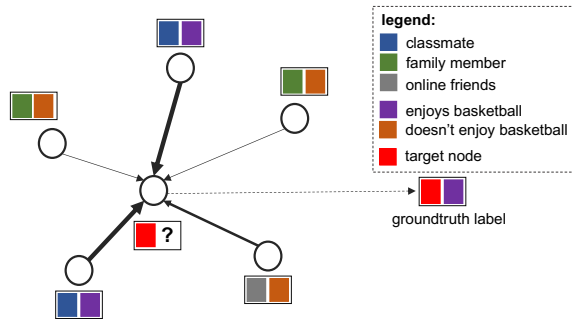
Fig. 1. Attention is used to assign importance to each type of neighbor. The link size denotes how much attention we want to apply to each neighbor. In this example, we see that by using attention to focus on a node's classmates, we can better predict the kind of activity the target is interested in. Best viewed with color.

systems [Deng et al. 2017], sensor networks [Aggarwal et al. 2017], epidemiology [Moslonka-Lefebvre et al. 2011], anomaly and fraud analysis [Akoglu et al. 2015], and ecology [Allesina et al. 2005]. For instance, interactions between users on a social network can be captured using a graph where the nodes represent users and links denote user interaction and/or friendship [Backstrom and Leskovec 2011]. On the other hand, in chemoinformatics, we can build molecular graphs by treating the atoms as nodes and the bonds between atoms as edges [Duvenaud et al. 2015].

A large body of work – which we broadly categorize as the field of graph mining [Aggarwal and Wang 2010] – has emerged that focuses on gaining insights from graph data. Many interesting and important problems have been studied in this area. These include graph classification [Duvenaud et al. 2015], link prediction [Sun et al. 2011], community detection [Girvan and Newman 2002], functional brain network discovery [Bai et al. 2017], node classification and clustering [Perozzi et al. 2014], and influence maximization [He and Kempe 2014] with new problems constantly being proposed.

In the real-world, however, graphs can be both structurally large and complex [Ahmed et al. 2014; Leskovec and Faloutsos 2006] as well as noisy [He and Kempe 2014]. These pose a significant challenge to graph mining techniques, particularly in terms of performance [Ahmed et al. 2017].

Various techniques have been proposed to address this issue. For instance, the method proposed by [Wu et al. 2015] utilizes multiple views of the same graph to improve classification performance while [Zhang et al. 2016a] leverages auxiliary non-graph data as side views under similar motivation. Another popular technique involves the identification of task-relevant or discriminative subgraphs [Shi et al. 2012; Zhu et al. 2012].

Recently, a new approach has emerged to address the above-mentioned problem and this is by incorporating *attention* into graph mining solutions. An attention mechanism aids a model by allowing it to "focus on the most relevant parts of the input to make decisions" [Velickovic et al. 2018]. Attention was first introduced in the deep learning community to help models attend to important parts of the data [Bahdanau et al. 2015; Mnih et al. 2014]. The mechanism has been successfully adopted by models solving a variety of tasks. Attention was used by [Mnih et al. 2014] to take glimpses of relevant parts of an input image for image classification; on the other hand, [Xu et al. 2015] used attention to focus on important parts of an image for the image captioning task. Meanwhile [Bahdanau et al. 2015] utilized attention for the machine translation task by assigning weights which reflected the importance of different words in the input sentence when generating corresponding words in the output sentence. Finally, attention has also been used for both the

image [Yang et al. 2016] as well as the natural language [Kumar et al. 2016] question answering tasks. However, most of the work involving attention has been done in the computer vision or natural language processing domains.

More recently, there has been a growing interest in attention models for graphs and various techniques have been proposed [Choi et al. 2017; Feng et al. 2016; Han et al. 2018; Lee et al. 2018; Ma et al. 2017; Velickovic et al. 2018]. Although attention is defined in slightly different ways in all these papers, the competing approaches do share common ground in that attention is used to allow the model to focus on task-relevant parts of the graph. We discuss and define, more precisely, the main strategies used by these methods to apply attention to graphs in Section 6. Figure 1 shows a motivating example of when attention can be useful in a graph setting. In particular, the main advantages of using attention on graphs can be summarized as follows:

(1) Attention allows the model to avoid or ignore noisy parts of the graph, thus improving the signal-to-noise (SNR) ratio [Lee et al. 2018; Mnih et al. 2014].
(2) Attention allows the model to assign a relevance score to elements in the graph (for instance, the different neighbors in a node's neighborhood) to highlight elements with the most task-relevant information, further improving SNR [Velickovic et al. 2018].
(3) Attention also provides a way for us to make a model's results more interpretable [Choi et al. 2017; Velickovic et al. 2018]. For example, by analyzing a model's attention over different components in a medical ontology graph we can identify the main factors that contribute to a particular medical condition [Choi et al. 2017].

In this paper, we conduct a comprehensive and focused review of the literature on graph attention models. To the best of our knowledge, this is the first work on this topic. We introduce three different taxonomies to group existing work into intuitive categories. We then motivate our taxonomies through detailed examples and use each to survey competing approaches from a particular standpoint. In particular, we group the existing work by problem setting (defined by the type of input graph and the primary problem output), by the kind of attention that is used, and by the task (*e.g.*, node classification, graph classification, *etc.*).

In previous work, different kinds of graphs (*e.g.*, homogeneous graphs, heterogeneous graphs, directed acyclic graphs, *etc.*) have been studied with different properties (*e.g.*, attributed, weighted or unweighted, directed or undirected) and different outputs (*e.g.*, node embedding, link embedding, graph embedding). The first taxonomy allows us to survey the field from this perspective. The second taxonomy tackles the main strategies that have been proposed for applying attention in graphs. We then introduce a final taxonomy that groups the methods by application area; this shows the reader what problems have already been tackled while, perhaps more importantly, revealing important graph-based problems where attention models have yet to be applied. Figure 2 shows the proposed taxonomies.

Additionally, we also summarize the challenges that have yet to be addressed in the area of graph attention and provide promising directions for future work.

## 1.1 Main contributions

The main contributions of this work are as follows:

(1) We introduce three intuitive taxonomies for categorizing various graph attention models and survey existing methods using these taxonomies. To the best of our knowledge, this is the first survey on the important field of graph attention.
(2) We motivate each taxonomy by discussing and comparing different graph attention models from the taxonomy's perspective.
(3) We highlight the challenges that have yet to be addressed in the area of graph attention.
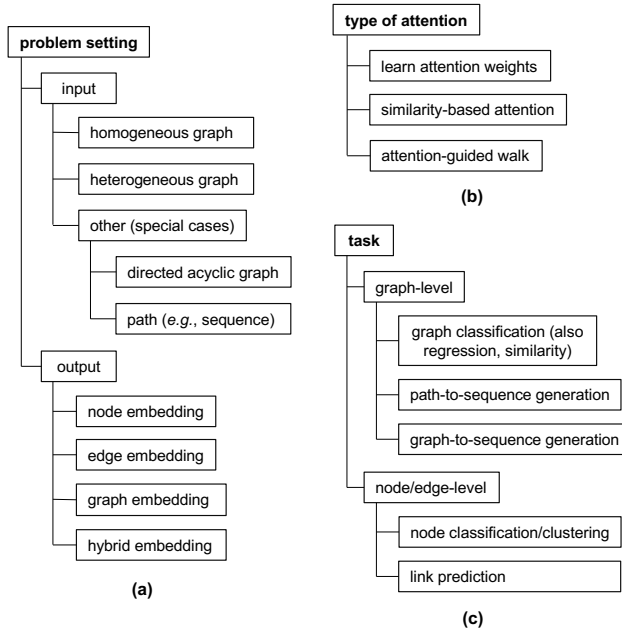
Fig. 2. Proposed taxonomies to group graph attention models based on (a) problem setting, (b) type of attention used, and (c) task or problem.

(4) We also provide suggestions on potential areas for future work in this emerging field.

## 1.2 Scope of this article

In this article, we focus on examining and categorizing various techniques that apply attention to graphs (we give a general definition of attention in Sec. 2). Most of these methods take graphs as input and solve some graph-based problem such as link prediction [Sun et al. 2011], graph classification/regression [Duvenaud et al. 2015], or node classification [Kipf and Welling 2017]. However, we also consider methods that apply attention to graphs although the graph is only one of several types of input to the problem.

We do *not* attempt to survey the vast field of general graph-based methods that do not explicitly apply attention, multiple work have already been done on this with each having a particular focus [Cai et al. 2018; Getoor and Diehl 2015; Jiang et al. 2013].

## 1.3 Organization of the survey

The rest of this survey is organized as follows. We start by introducing useful notations and definitions in Section 2. We then use Sections 3 through 5 to discuss related work using our main taxonomy (Fig. 2a). We organized Sections 3-5 such that the methods in existing work are grouped by the main type of embedding they calculate (*e.g.*, node embedding, edge embedding, graph embedding, or hybrid embedding); these methods are then further divided by the type of graphs they support. In Sections 6 and 7, we switch to a different perspective and use the remaining taxonomies (Fig. 2b and Fig. 2c) to guide the discussion. We then discuss challenges as well as interesting opportunities for future work in Section 8. Finally, we conclude the survey in Section 9.

## 2 PROBLEM FORMULATION

In this section, we define the different types of graphs that appear in the discussion; we also introduce related notations.

**DEFINITION 1 (HOMOGENEOUS GRAPH).** *Let $G = (V, E)$ be a graph where $V$ is the set of nodes (vertices) and $E$ is the set of edges between the nodes in $V$.*

Further, let $\mathbf{A} = \begin{bmatrix} A_{ij} \end{bmatrix}$ be the $n \times n$, for $n = |V|$, adjacency matrix of $G$ where $A_{ij} = 1$ if there exists $(v_i, v_j) \in E$ and $A_{ij} = 0$ otherwise. When $\mathbf{A}$ is a binary matrix (*i.e.*, $A_{ij} \in \{0, 1\}$), we consider the graph to be *unweighted*. Note that $\mathbf{A}$ may also encode edge weights; given a weight $w \in \mathbb{R}$ for $(v_i, v_j) \in E$, $A_{ij} = w$. In this case, the graph is said to be *weighted*. Also, if $A_{ij} = A_{ji}$, for any $1 \le i, j \le n$, then the graph is *undirected*, otherwise it is *directed*.

Given an adjacency matrix $\mathbf{A}$, we can construct a right stochastic matrix $\mathbf{T}$ – also called the transition matrix – which is simply $\mathbf{A}$ with rows normalized to 1. Also, given a vertex $v \in V$, let $\Gamma_v$ be the set of vertices in node $v$'s neighborhood (for instance, a popular definition for neighborhood is simply the one-hop neighborhood of $v$, or $\Gamma_v = \{w | (v, w) \in E\}$).

**DEFINITION 2 (HETEROGENEOUS GRAPH).** *A heterogeneous graph is defined as $G = (V, E)$ consisting of a set of node objects $V$ and a set of edges $E$ connecting the nodes in $V$. A heterogeneous graph also has a node type mapping function $\theta : V \to \mathcal{T}_V$ and an edge type mapping function defined as $\xi : E \to \mathcal{T}_E$ where $\mathcal{T}_V$ and $\mathcal{T}_E$ denote the set of node types and edge types, respectively. The type of node $i$ is denoted as $\theta(i)$ (which may be an author, a paper, or a conference in a heterogeneous bibliographic network) whereas the type of edge $e = (i, j) \in E$ is denoted as $\xi(i, j) = \xi(e)$ (e.g., a co-authorship relationship, or a "published in" relationship).*

Note a heterogeneous graph is sometimes referred to as a typed network. Furthermore, a bipartite graph is a simple heterogeneous graph with two node types and a single edge type. A homogeneous graph is simply a special case of a heterogeneous graph where $|\mathcal{T}_V| = |\mathcal{T}_E| = 1$.

**DEFINITION 3 (ATTRIBUTED GRAPH).** *Let $G = (V, E, \mathbf{X})$ be an attributed graph where $V$ is a set of nodes, $E$ is a set of edges, and $\mathbf{X}$ is an $n \times d$ matrix of* node input attributes *where each $\bar{\mathbf{x}}_i$ (or $\mathbf{X}_{i:}$) is a d-dimensional (row) vector of attribute values for node $v_i \in V$ and $\mathbf{x}_j$ (or $\mathbf{X}_{:j}$) is an n-dimensional vector corresponding to the jth attribute (column) of $\mathbf{X}$. Alternatively, $\mathbf{X}$ may also be a $m \times d$ matrix of* edge input attributes. *These may be real-valued or not.*

**DEFINITION 4 (DIRECTED ACYCLIC GRAPH (DAG)).** *A directed graph with no cycles.*

The different "classes" of graphs defined in Definitions 1-3 can be directed or undirected as well as weighted or unweighted. Other classes of graphs arise from composing the distinct "graph classes" (Definition 1-3). For instance, it is straightforward to define an attributed heterogeneous network $G = (V, E, \mathcal{T}_V, \mathcal{T}_E, \mathbf{X})$.

**DEFINITION 5 (PATH).** *A path $P$ of length $L$ is defined as a sequence of* unique *indices $i_1, i_2, \ldots, i_{L+1}$ such that $(v_{i_t}, v_{i_{t+1}}) \in E$ for all $1 \le t \le L$. The length of a path is defined as the number of edges it contains.*

Unlike walks that allow loops, paths do not and therefore a walk is a path iff it has no loops (all nodes are unique). Note that a path is a special kind of graph with a very rigid structure where all the nodes have at most 2 neighbors. We now give a general, yet formal, definition of the notion of *graph attention* as follows:

**DEFINITION 6 (GRAPH ATTENTION).** *Given a target graph object (e.g., node, edge, graph, etc), $v_0$ and a set of graph objects in $v_0$'s "neighborhood" $\{v_1, \cdots, v_{|\Gamma_{v_0}|}\} \in \Gamma_{v_0}$ (the neighborhood is task-specific*

Table 1. Summary of notation. Matrices are bold upright roman letters; vectors are bold lowercase letters.

| | |
|---|---|
| $G$ | (un)directed (attributed) graph |
| $\mathbf{A}$ | adjacency matrix of the graph $G = (V, E)$ |
| $\mathbf{T}$ | stochastic transition matrix for $G$ constructed by normalizing over the rows in $\mathbf{A}$ |
| $n, m$ | number of nodes and edges in the graph |
| $k$ | number of embedding dimensions |
| $d$ | node attribute dimension |
| $\Gamma_i$ | neighbors of $i$ defined using some neighborhood |
| $\Gamma_i^{(\ell)}$ | $\ell$-neighborhood $\Gamma_i^{(\ell)} = \{j \in V \mid \text{dist}(i, j) \leq \ell\}$ |
| $\text{dist}(i, j)$ | shortest distance between $i$ and $j$ |
| $\mathbf{X}$ | $n \times d$ input attribute matrix |
| $\mathbf{x}$ | a $d$-dimensional feature vector |
| $x_i$ | the $i$-th element of $\mathbf{x}$ |
| $\mathbf{Z}$ | output node embedding matrix (or vector $\mathbf{z}$ in the case of graph embeddings) |

and can mean the $\ell$-hop neighborhood of a node or something more general). Attention is defined as a function $f' : \{v_0\} \times \Gamma_{v_0} \to [0, 1]$ that maps each of the object in $\Gamma_{v_0}$ to a relevance score which tells us how much attention we want to give to a particular neighbor object. Furthermore, it is usually expected that $\sum_{i=1}^{|\Gamma_{v_0}|} f'(v_0, v_i) = 1$.

A summary of the notation used throughout the manuscript is provided in Table 1.

## 3 ATTENTION-BASED NODE/EDGE EMBEDDING

In this section – as well as the two succeeding sections – we introduce various graph attention models and categorize them by problem setting. For easy reference, we show all of the surveyed methods in Table 2, taking care to highlight where they belong under each of the proposed taxonomies. We now begin by defining the traditional node embedding problem.

DEFINITION 7 (NODE EMBEDDING). *Given a graph $G = (V, E)$ with $V$ as the node set and $E$ as the edge set, the objective of node embedding is to learn a function $f : V \to \mathbb{R}^k$ such that each node $i \in V$ is mapped to a $k$-dimensional vector $\mathbf{z}_i$ where $k \ll |V|$. The node embedding matrix is denoted as $\mathbf{Z}$.*

The learned node embeddings given as output can subsequently be used as input to mining and machine learning algorithms for a variety of tasks such as link prediction [Sun et al. 2011], classification [Velickovic et al. 2018], community detection [Girvan and Newman 2002], and role discovery [Rossi and Ahmed 2015]. We now define the problem of attention-based node embedding as follows:

DEFINITION 8 (ATTENTION-BASED NODE EMBEDDING). *Given the same inputs as above, we learn a function $f : V \to \mathbb{R}^k$ that maps each node $i \in V$ to an embedding vector $\mathbf{z}_i$. Additionally, we learn a second function $f' : V \times V \to [0, 1]$ to assign "attention weights" to the elements in a target node $i$'s neighborhood $\Gamma_i$. The function $f'$ defines each neighbor $j$'s, for $j \in \Gamma_i$, importance relative to the target node $i$. Typically, $f'$ is constrained to have $\sum_{j \in \Gamma_i} f'(i, j) = 1$ for all $i$ with $f'(i, k) = 0$ for all $k \notin \Gamma_i$. The goal of attention-based node embedding is to assign a similar embedding to node $i$ and to its more similar or important neighbors, i.e., $\delta(\mathbf{z}_i, \mathbf{z}_j) > \delta(\mathbf{z}_i, \mathbf{z}_k)$ iff $f'(i, j) > f'(i, k)$, where $\delta$ is some similarity measure.*

Above, we defined attention-based node embedding; attention-based edge embedding can also be defined similarly.

Table 2. Qualitative and quantitative comparison of graph-based attention models.

| Method | INPUT | | | | OUTPUT | | | | MECHANISM | | | | TASK | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Homogeneous graph | Heterogeneous graph | Directed acyclic graph | Path | Node embedding | Edge embedding | Graph embedding | Hybrid embedding | Learn attention weights | Similarity-based attention | Attention-guided walk | Other | Node/Link classification | Link prediction | Graph classification | Graph regression | Seq-to-seq generation | Graph-to-seq generation | Other |
| **AttentionWalks** [Abu-El-Haija et al. 2017] | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **GAKE** [Feng et al. 2016] | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **GAT** [Velickovic et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **AGNN** [Thekumparampil et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **PRML** [Zhao et al. 2017] | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| **EAGCN** [Shang et al. 2018] | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| **Modified − GAT** [Ryu et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ |
| **graph2seq** [Xu et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ |
| **GAM** [Lee et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **RNNSearch** [Bahdanau et al. 2015] | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| **Att − NMT** [Luong et al. 2015] | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| **Dipole** [Ma et al. 2017] | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ |
| **GRAM** [Choi et al. 2017] | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| **CCM** [Zhou et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✗ |
| **JointD/E + SATT** [Han et al. 2018] | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ | ✗ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

Here, we use a single section to discuss both node and edge embeddings since there has not been a lot of work on attention-based edge embeddings and also because the two problems are quite similar [Cai et al. 2018]. We now categorize the different methods that calculate attention-based node/edge embeddings based on the type of graph they support.

## 3.1 Homogeneous graph

Most of the work that calculate attention-based node embeddings focus on homogeneous graphs [Abu-El-Haija et al. 2017; Thekumparampil et al. 2018; Velickovic et al. 2018; Zhao et al. 2017]. Also, all the methods assume the graph is attributed although [Abu-El-Haija et al. 2017] only needs node labels (in this case, attribute size $d = 1$).

The method proposed by [Abu-El-Haija et al. 2017], called ATTENTIONWALKS, is most reminiscent of popular node embedding approaches such as DeepWalk and node2vec [Grover and Leskovec 2016; Perozzi et al. 2014] in that a random walk is used to calculate node contexts. Given a graph $G$ with its corresponding transition matrix $\mathbf{T}$, and a context window size $c$, we can calculate the expected number of times walks started at each node visits other nodes via:

$$\mathbb{E}[\mathbf{D}|a_1, \cdots, a_c] = \mathbf{I}_n \sum_{i=1}^{c} a_i(\mathbf{T})^i.$$

Here $\mathbf{I}_n$ is the size-$n$ identity matrix, $a_i$, for $1 \le i \le c$, are learnable attention weights, and $\mathbf{D}$ is the walk distribution matrix where $D_{uv}$ encodes the number of times node $u$ is expected to visit node $v$

given that a walk is started from each node. In this scenario, attention is thus used to steer the walk towards a broader neighborhood or to restrict it within a narrower neighborhood (*e.g.*, when the weights are top-heavy). This solves the problem of having to do a grid-search to identify the best hyper-parameter $c$ as studies have shown that this has a noticeable impact on performance [Perozzi et al. 2014] – note that for DeepWalk the weights are fixed at $a_i = 1 - \frac{i-1}{c}$.

Another attention-based method that is very similar in spirit to methods like DeepWalk, node2vec, and LINE [Grover and Leskovec 2016; Perozzi et al. 2014; Tang et al. 2015] in that it uses co-occurrence information to learn node embeddings is GAKE [Feng et al. 2016]. It is important to point out that GAKE builds a graph, commonly referred to as a knowledge graph, from knowledge triplets $(h, t, r)$ where $h$ and $t$ are terms connected by a relationship $r$. Given three triplets (Jose, Tagalog, speaks), (Sato, Nihongo, speaks), and (Jose, Sinigang, eats) we can construct a simple heterogeneous graph with three types of nodes, namely {person, language, food}, and two types of relationships, namely {speaks, eats}. However we categorize GAKE as a homogeneous graph method as it doesn't seem to make a distinction between different kinds of relationships or node types (see metapath2vec [Dong et al. 2017] for a node embedding method that explicitly models the different kinds of relationships and nodes in a heterogeneous graph). Instead, the method takes a set of knowledge triplets and builds a directed graph by taking each triplet $(h, t, r)$ and adding $h$ and $t$ as the head and tail nodes, respectively, while adding an edge from the head to the tail (they also add a reverse link). The main difference between GAKE and methods such as DeepWalk or node2vec is that they include edges when calculating a node's context. They define three different contexts to get related subjects (nodes or edges), formally they maximize the log-likelihood:

$$\sum_{s \in V \cup E} \sum_{c \in \Gamma_s} log \ \Pr(s \mid c)$$

where $\Gamma_s$ is a set of nodes and/or edges defining the neighborhood context of $s$. To get the final node embedding for a given subject $s$ in the graph, they use attention to obtain the final embedding $\mathbf{z}_s = \sum_{c \in \Gamma'_s} \alpha(c) \mathbf{z}'_c$ where $\Gamma'_s$ is some neighborhood context for $s$, $\alpha(c)$ defines the attention weights for context object $c$ and $\mathbf{z}'_c$ is the learned embedding for $c$.

On the other hand, methods like GAT [Velickovic et al. 2018] and AGNN [Thekumparampil et al. 2018] extend graph convolutional networks (GCN) [Kipf and Welling 2017] by incorporating an explicit attention mechanism. Recall that a GCN is able to propagate information via an input graph's structure using the propagation rule:

$$\mathbf{H}^{(l+1)} = \sigma(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}^{(l)})$$

where $\mathbf{H}^{(l)}$ indicates the learned embedding matrix at layer $l$ of the GCN with $\mathbf{H}^{(0)} = \mathbf{X}$. Also, $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_n$ is the adjacency matrix of an undirected graph $\mathbf{A}$ with added self loop. The matrix $\tilde{\mathbf{D}}$, on the other hand, is defined as the diagonal degree matrix of $\tilde{\mathbf{A}}$ so, in other words, $\tilde{D}_{i,i} = \sum_j \tilde{A}_{i,j}$. Hence, the term $\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$ computes a symmetric normalization (similar to the normalized graph Laplacian) for the graph defined by $\tilde{\mathbf{A}}$. Finally, $\mathbf{W}^{(l)}$ is the trainable weight-matrix for level $l$ and $\sigma(\cdot)$ is a nonlinearity like ReLU, Sigmoid, or tanh.

A GCN works like an end-to-end differentiable version of the Weisfeiler-Lehman algorithm [Shervashidze et al. 2011] where each additional layer allows us to expand and integrate information from a larger neighborhood. However, because we use the term $\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{H}$ in the propagation, we are essentially applying a weighted sum of the features of neighboring nodes normalized by their degrees. GAT and AGNN essentially introduce attention weights $a_{uv}$ to determine how much attention we want to give to a neighboring node $v$ from node $u$'s perspective. The two methods differ primarily in the way attention is defined (we expound on this in Sec. 6). Furthermore, [Velickovic

et al. 2018] introduce the concept of "multi-attention" which basically defines multiple attention heads (*e.g.*, weights) between a pair of objects in the graph.

PRML [Zhao et al. 2017] is another approach that learns edge embeddings but they use a different strategy to define attention. In [Zhao et al. 2017], a path-length threshold $L$ is defined and a recurrent neural network [Gers et al. 2000] learns a path-based embedding for paths of length $1, \cdots, L$ between pairs of nodes $u$ and $v$. Attention is defined in two ways. First, attention is used to identify important nodes along paths and this helps in calculating the intermediary path embeddings. Second, attention then assigns priority to paths that are more indicative of link formation and this is used to highlight important or task-relevant path embeddings when these are integrated to form the final edge embedding. We can think of this approach as an attention-based model that calculates the Katz betweenness score for pairs of nodes [Liben-Nowell and Kleinberg 2003].

## 3.2 Heterogeneous graph

The only work, to the best of our knowledge, that has been proposed to calculate attention-guided node embeddings for heterogeneous graphs is EAGCN [Shang et al. 2018]. Very similar to both GAT and AGNN, [Shang et al. 2018] proposes to apply attention to a GCN [Duvenaud et al. 2015; Kipf and Welling 2017]. However, they handle the case where there can be multiple types of links connecting nodes in a graph. Thus they propose to use "multi-attention," like [Velickovic et al. 2018], and each of the attention mechanisms considers neighbors defined only by a particular link type. Although the authors validate EAGCN using the graph regression task, the method can be used without much adjustments for node-level tasks since the graph embeddings in EAGCN are simply concatenations or sums of the learned attention-guided node embeddings.

However, the above-mentioned work assumes the given heterogeneous network only has one type of node, *i.e.*, $|\mathcal{T}_V| = 1$.

## 3.3 Other special cases

Unlike the areas of attention-based graph embedding, there does not seem to be work on calculating attention-guided node-embeddings for special types of graphs which appear in certain domains (*e.g.*, medical ontologies represented as DAGs, or certain Heterogeneous networks represented as star graphs).

Since the above-mentioned methods work for general graphs, they should be suitable for more specific types of graphs and one should be able to apply them directly to these cases.

## 4 ATTENTION-BASED GRAPH EMBEDDING

Similarly, we begin the discussion by defining the traditional graph embedding problem.

DEFINITION 9 (GRAPH EMBEDDING). *Given a set of graphs, the objective of graph embedding is to learn a function $f : \mathcal{G} \to \mathbb{R}^k$ that maps an input graph $G \in \mathcal{G}$ to a low dimensional embedding vector $\mathbf{z}$ of length $k$; here $\mathcal{G}$ is the input space of graphs. Typically, we want to learn embeddings that group similar graphs (e.g., graphs belonging to the same class) together.*

The learned graph embeddings can then be fed as input to machine learning/data mining algorithms to solve a variety of graph-level tasks such as graph classification [Lee et al. 2018], graph regression [Duvenaud et al. 2015], and graph-to-sequence generation [Xu et al. 2018]. We now define the problem of attention-based graph embedding as follows:

DEFINITION 10 (ATTENTION-BASED GRAPH EMBEDDING). *Given the same inputs as above, we learn a function $f : \mathcal{G} \to \mathbb{R}^k$ that maps each input graph $G \in \mathcal{G}$ to an embedding vector $\mathbf{z}$ of length $k$.*

*Additionally, we learn a second function f′ that assigns "attention weights" to different subgraphs of a given graph to allow the model to prioritize more important parts (i.e., subgraphs) of the graph when calculating its embedding.*

### 4.1  Homogeneous graph

Several methods have been proposed for learning graph embeddings on homogeneous graphs [Lee et al. 2018; Ryu et al. 2018; Xu et al. 2018].

[Ryu et al. 2018] simply use the formulation of GAT [Velickovic et al. 2018] and make a few adjustments to the way attention is calculated. However, aside from the change in how attention weights are calculated, there isn't much difference between GAT and the method proposed in [Ryu et al. 2018]. To obtain the final graph embedding for the task of graph regression on molecular graphs, the proposed method adds fully connected layers after the final GAT layer to flatten the per-node outputs.

On the other hand, [Xu et al. 2018] propose GRAPH2SEQ which solves the natural language question answering task. Given a set of facts, in the form of sentences their approach is quite unique in that they convert the input into an attributed directed homogeneous graph. To obtain the graph embedding that is used for the sequence generation task, they first learn node embeddings for each node in the graph. Node embeddings are formed by concatenating a forward and a backward representation for each node which are representations derived by aggregating information from each node's forward (neighbors traversed using forward links) and backward (similarly, neighbors traversed using reverse links) neighborhoods, respectively. The final graph embedding is obtained by pooling the individual node embeddings or by simply aggregating them. Attention, in the case of GRAPH2SEQ however, is applied by also attending to the individual node embeddings during sequence generation. Since each node embedding $\mathbf{z}_i$ captures information in the region around node $i$ (we can think of this as a subgraph focused around $i$), attention allows use to prioritize a particular part of the graph when generating a word in the output sentence (sequence). This captures the intuition that different parts of the graph can be associated, primarily, with different concepts or words.

Finally, in a previous work [Lee et al. 2018], we proposed GAM which uses two types of attention to learn a graph embedding. The main idea is to use an attention-guided walk to sample relevant parts of the graph to form a subgraph embedding. In GAM, we took a walk of length $L$. Let $1 \leq i \leq L$ be the $i$-th node discovered in the walk and $\mathbf{x}_i$ the corresponding node attribute vector, an RNN is used to integrate the node attributes $(\mathbf{x}_1, \cdots, \mathbf{x}_L)$ to form a subgraph embedding $\mathbf{s}$ (the subgraph or region covered during the walk). During each step, an attention mechanism is used to determine which neighboring node is more relevant to allow the walk to cover more task-relevant parts of the graph. To get the final graph embedding, we deploy $z$ "agents" to sample from various parts of the graph yielding embeddings $\{\mathbf{s}_1, \cdots, \mathbf{s}_z\}$. A second attention mechanism is then used to determine the relevance of the various subgraph embeddings before these are combined to form a graph embedding. In other words an attention mechanism is defined as a function $\alpha : \mathbb{R}^k \rightarrow [0, 1]$ which maps a given subgraph embedding to a relevance score. The graph embedding is thus defined as $\sum_{i=1}^{z} \alpha(\mathbf{s}_i)\, \mathbf{s}_i$.

### 4.2  Heterogeneous graph

Recall from our previous discussion that EAGCN [Shang et al. 2018] was used to study the task of graph regression. EAGCN used an approach similar to GAT to generate a graph embedding. Since the method uses a similar attention mechanism as GAT, attention is focused on determining important neighbors to attend to when calculating node embeddings for a graph. The final graph embedding is then a concatenation/sum of the attention-guided node embeddings. It shouldn't be

difficult to apply a second attention mechanism, similar to that of [Lee et al. 2018], to weight the importance of the different node embeddings.

## 4.3 Other special cases

Attention was originally studied in the Computer Vision and NLP domains and various RNN models using attention on sequence-based tasks were proposed [Bahdanau et al. 2015; Luong et al. 2015]. Since sequences are technically no different from paths, we introduce some notable attention models under this setting.

A sequence (*e.g.*, a sentence) of length $L$ can be represented as a directed attributed graph of length $L$ – the $i$-th attribute vector $\mathbf{x}_i$ of node $i$ is then a representation of the $i$-th component in the sequence (a one-hot word embedding or a word2vec embedding, for instance). In the proposed methods [Bahdanau et al. 2015; Luong et al. 2015; Ma et al. 2017], the node attribute vectors $\{\mathbf{x}_i, \cdots, \mathbf{x}_L\}$ are fed one after the other into an RNN. Recall that in the simplest case, a recurrent neural network calculates a hidden embedding $\mathbf{h}_i$ for each input $i$ via the rule

$$\mathbf{h}_i = \sigma(\mathbf{W}_h \mathbf{x}_i + \mathbf{U}_h \mathbf{h}_{i-1}) + \mathbf{b}_h,$$

where $\mathbf{W}_h$ and $\mathbf{U}_h$ are trainable weight matrices for the input $\mathbf{x}_i$ and the previous hidden embedding $\mathbf{h}_{i-1}$, respectively; $\mathbf{b}_h$ is a bias vector and $\mathbf{h}_0$ is usually initialized to the zero vector.

In this case, we can think of $\mathbf{h}_i$ as node $i$'s node embedding which integrated information from the sub-path encompassing nodes $1, \cdots, i$. Attention can then be applied on node embeddings to generate an attention-guided graph embedding. In particular, the method proposed in [Bahdanau et al. 2015] called RNNSEARCH defined the graph embedding as $\mathbf{z} = \sum_{i=1}^{L} \alpha_i \mathbf{h}_i$ where attention is assigned to each hidden embedding $\mathbf{h}_i$ depending on its relevance to the task.

When the length of the path $L$ is large, however, attending over all the hidden embeddings as in RNNSEARCH may not yield the best results. [Luong et al. 2015] proposed to use two attention mechanisms, the first is similar in spirit to that of [Bahdanau et al. 2015]. However, the second attention allows the model to select a local point within the input and focus attention there. More concretely, depending on the needs of the task, the second attention mechanism outputs a position $1 \leq p \leq L$ and the graph embedding is constructed from the hidden node embeddings $\mathbf{h}_{p-D}, \cdots, \mathbf{h}_{p+D}$ where $D$ is an empirically selected attention window size.

Finally, DIPOLE [Ma et al. 2017] applied an attention model similar to that of [Bahdanau et al. 2015] to sequential medical data. They used it to diagnose or predict medical conditions from the sequence. In practice, both [Ma et al. 2017] and [Bahdanau et al. 2015] used a bi-directional model which processes an input sequence using two RNNs, one taking the input sequence in its original order and another which takes the input sequence in reverse order.

## 5 ATTENTION-BASED HYBRID EMBEDDING

In this section we discuss methods that apply attention on graph data. However, the calculated embeddings are "hybrid" since the methods here also take data of other modalities (*e.g.*, text) and the learned embedding is a combination of all inputs.

## 5.1 Homogeneous graph

Like GAKE [Feng et al. 2016], CCM [Zhou et al. 2018] deals with knowledge graphs. However, in [Zhou et al. 2018], the sequence-to-sequence generation problem is studied. CCM uses an encoder-decoder model [Sutskever et al. 2014] to encode a sentence and output a sentence (*e.g.*, for sentence translation, or question and answering). However, the setting used in CCM assumes that the model refers to a knowledge graph for the question-and-answer task. Their knowledge graph

is made up of multiple connected components, each of which is made up of multiple knowledge triplets. A subgraph embedding $\mathbf{z}_i$ is learned for each component $i$:

$$\mathbf{z}_i = \sum_{n=1}^{|T_i|} \alpha_n^{(i)} [\mathbf{h}_n^{(i)}; \mathbf{t}_n^{(i)}]$$

where $T_i = \{(\mathbf{h}_1^{(i)}, \mathbf{t}_1^{(i)}, \mathbf{r}_1^{(i)}), \cdots, (\mathbf{h}_{n_i}^{(i)}, \mathbf{t}_{n_i}^{(i)}, \mathbf{r}_{n_i}^{(i)})\}$ consists of the embeddings for the corresponding knowledge triplets for $i$ and $\alpha_n^{(i)}$ is used to assign importance to different triplets (represented by the concatenation of the term embeddings) taking into account the terms as well as their relationship.

While the model is in the process of decoding (generating words for the output), it selectively identifies important subgraphs by attending on their embeddings $\mathbf{z}_i$. Furthermore, an attention mechanism is also used to identify the triplets within a selected knowledge subgraph to identify important triplets for word generation. Like GAKE, we consider CCM a homogeneous graph method since it doesn't seem to make an explicit distinction between different types of nodes and relationships.

JointD/E + SATT [Han et al. 2018] is another work that applies attention on a knowledge graph, similar to [Feng et al. 2016; Zhou et al. 2018]. However, they also use a large text corpus that may contain references to the different terms in the knowledge graph implying certain relationships. They introduce a method that uses attention to learn a joint embedding from graph and text data for the task of knowledge graph link prediction.

Under certain settings (brain network construction, for instance) the datasets tend to be quite small and noisy [Zhang et al. 2016b]. Methods such as that proposed by [Zhang et al. 2016b] proposes to use side information to regularize the graph construction process to highlight more discriminative patterns – which is useful when the output graphs are used for graph classification. Exploring the possibility of adding attention in this setting is interesting.

## 5.2 Heterogeneous graph

While there has been work proposed like [Han et al. 2018; Zhou et al. 2018] that applies attention on knowledge graphs – which are considered heterogeneous graphs – these models do not distinguish between the different types of links and nodes explicitly. To the best of our knowledge, there currently does not exist any work that has considered this setting. One possibility is to extend the idea proposed by EAGCN [Shang et al. 2018] to more general heterogeneous graphs.

## 5.3 Other special cases

GRAM [Choi et al. 2017] is a graph-based attention model for doing classification/regression on clinical data. Clinical records can usually be described by clinical codes $c_1, \cdots, c_{|C|} \in C$ in a vocabulary $C$. GRAM constructs a DAG whose leaves are the codes $c_1, \cdots, c_{|C|}$ while the ancestor nodes are more general medical concepts. [Choi et al. 2017] uses an attention mechanism to learn a $k$-dimensional final embedding of each leaf node $i$ (medical concept) via:

$$\mathbf{g}_i = \sum_{j \in \mathcal{A}(i)} \alpha_{i,j} \, \mathbf{e}_j.$$

where $\mathcal{A}(i)$ denotes the set comprised of $i$ and all its ancestors and $\mathbf{e}_j$ is the $k$-dimensional basic embedding of the node $j$. The use of attention allows the model to refer to more informative or relevant general concepts when a concept in $C$ is less helpful for medical diagnosis (*e.g.*, it is a rare concept). Since a patient's clinical visit record is represented as a binary vector $\mathbf{x} \in \{0, 1\}^{|C|}$ indicating which clinical codes were present for a particular visit the learned embeddings
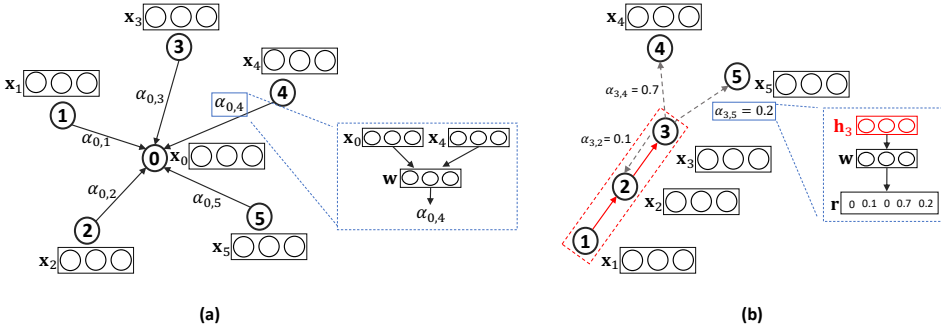
Fig. 3. (a) Given a target object $v_0$, we assign importance weights $\alpha_{0,i}$ to the objects $i$ in our neighborhood. This can be done by learning a function that assigns importance by examining (possibly) the hidden embeddings of $v_0$ and $v_i$, $\mathbf{x}_0$ and $\mathbf{x}_i$. (b) The hidden embedding $\mathbf{h}_3$ represents information $(\mathbf{x}_1, \cdots, \mathbf{x}_3)$ we've integrated after a walk of length $L = 3$, we input this into a ranking function that determines the importance of various neighbor nodes and this is used to bias the next step. For both examples, we use $\mathbf{w}$ to represent the trainable parameters of the attention function.

$\mathbf{g}_1, \cdots, \mathbf{g}_{|C|}$ can be stacked to form a $k \times |C|$ embedding matrix to embed $\mathbf{x}$; this can then be inputted into a predictive model for medical diagnosis. Note that the problem of medical diagnosis is a classical supervised learning task which takes clinical code feature vectors but GRAM applies attention on a medical concept DAG for the purpose of embedding the given feature vectors.

## 6 TYPES OF GRAPH ATTENTION MECHANISM

We now describe the three main types of attention that have been applied to graph data. While all three types of attention share the same purpose or intent, they differ in how attention is defined or implemented. In this section we provide examples from the literature to illustrate how each type is implemented.

Recall from our general definition of attention in Def. 6 that we are given a target graph object (*e.g.*, node, edge, graph, *etc*) $v_0$ and a set of graph objects in $v_0$'s "neighborhood" $\{v_1, \cdots, v_{|\Gamma_{v_0}|}\} \in \Gamma_{v_0}$. Attention is defined as a function $f' : \{v_0\} \times \Gamma_{v_0} \to [0, 1]$ that maps each of the objects in $\Gamma_{v_0}$ to a relevance score. In practice this is usually done in one of three ways which we introduce below.

### 6.1 Learn attention weights

Given the corresponding attributes/features $\mathbf{x}_0, \mathbf{x}_1, \cdots, \mathbf{x}_{|\Gamma_{o^*}|}$ for $v_0, v_1, \cdots, v_{|\Gamma_{v_0}|}$, attention weights can be learned via:

$$\alpha_{0,j} = \frac{e_{0,j}}{\sum_{k \in \Gamma_{v_0}} e_{0,k}}$$

where $e_{0,j}$ is node $v_j$'s relevance to $v_0$. In practice, this is typically implemented using softmax with a trainable function learning $v_j$'s relevance to $v_0$ by considering their attributes. The implementation in GAT [Velickovic et al. 2018] illustrates this:

$$\alpha_{0,j} = \frac{\exp\Big(\text{LeakyReLU}\Big(\mathbf{a}[\mathbf{W}\mathbf{x}_0||\mathbf{W}\mathbf{x}_j]\Big)\Big)}{\sum_{k \in \Gamma_{v_0}} \exp\Big(\text{LeakyReLU}\Big(\mathbf{a}[\mathbf{W}\mathbf{x}_0||\mathbf{W}\mathbf{x}_k]\Big)\Big)}$$

where $\mathbf{a}$ is a trainable attention vector, $\mathbf{W}$ is a trainable weight matrix mapping the input features to the hidden space, and $||$ represents concatenation. An illustration of this is shown in Fig. 3a.

## 6.2 Similarity-based attention

Again, given the corresponding attributes or features, the second type of attention can be learned similarly as above except for a key difference. We call this approach similarity-based attention as more attention is given to object's that share more similar hidden representations or features, this is also often referred to in the literature as alignment [Bahdanau et al. 2015]. To illustrate this, we use the definition given in AGNN [Thekumparampil et al. 2018]:

$$\alpha_{0,j} = \frac{\exp\left(\beta \cdot \cos\left(\mathbf{W}\mathbf{x}_0, \mathbf{W}\mathbf{x}_j\right)\right)}{\sum_{k \in \Gamma_{v_0}} \exp\left(\beta \cdot \cos\left(\mathbf{W}\mathbf{x}_0, \mathbf{W}\mathbf{x}_k\right)\right)}$$

where $\beta$ is a trainable bias and "cos" represents cosine-similarity; like before, $\mathbf{W}$ is a trainable weight matrix to map inputs to the hidden space. Note that this is very similar to the above definition. The difference is that the model explicitly learns similar hidden embeddings for objects that are relevant to each other since attention is based on similarity or alignment.

## 6.3 Attention-guided walk

While the first two types of attention focuses on choosing relevant information to integrate into a target object's hidden representation, the third type of attention has a slightly different purpose. We use GAM [Lee et al. 2018] to illustrate this idea. GAM takes a series of steps on an input graph and encodes information from visited nodes using an RNN to construct a subgraph embedding. The RNN's hidden state at a time $t$, $\mathbf{h}_t \in \mathbb{R}^h$ encodes information from the nodes that were previously visited by the walk from steps $1, \cdots, t$. Attention is then defined as a function $f' : \mathbb{R}^h \to \mathbb{R}^k$ that maps an input hidden vector $f'(\mathbf{h}_t) = \mathbf{r}_{t+1}$ to a $k$-dimensional rank vector that tells us which of the $k$ types of nodes we should prioritize for our next step. The model will then prioritize neighboring nodes of a particular type for the next step. We illustrate this in Fig. 3b.

## 7 GRAPH ATTENTION TASKS

The different attention-based graph methods can be divided broadly by the kind of problem they solve: node-level or graph-level.

### 7.1 Node-level

A number of work have been proposed to study graph attention for node-level tasks, the most notable of which are node classification and link prediction [Abu-El-Haija et al. 2017; Feng et al. 2016; Han et al. 2018; Thekumparampil et al. 2018; Velickovic et al. 2018; Zhao et al. 2017]. Although each method differs in their approach and assumptions, they share a common technique which is to learn an attention-guided node or edge embedding which can then be used to train a classifier for classification or link prediction. It isn't hard to see these methods implemented for the related tasks of node clustering. Some notable node-level tasks for which attention-based graph methods have not been proposed include node/edge role discovery [Rossi and Ahmed 2015], and node ranking [Sun et al. 2009a].

### 7.2 Graph-level

Multiple works have also studied graph-level tasks such as graph classification and graph regression. In this setting, an attention-guided graph embedding is constructed by attending to relevant parts of the graph. Methods like EAGCN, Modified-GAT, GAM, and Dipole [Lee et al. 2018; Ma et al. 2017; Ryu et al. 2018; Shang et al. 2018] learn a graph embedding for the more standard graph-based

tasks of classification and regression. It is not hard to see how these methods can be applied to the related problem of graph similarity search [Zheng et al. 2013].

On the other had, work like [Bahdanau et al. 2015; Luong et al. 2015; Xu et al. 2018; Zhou et al. 2018] generate sequences from input graph data. Notably, GRAPH2SEQ proposes a method that outputs a sequence given an input graph instead of the more methods that do sequence-to-sequence generation.

Finally, GRAM [Choi et al. 2017] applied attention to a medical ontology graph to help learn attention-based embeddings for medical codes. While the problem they studied was the problem of classifying a patient record (described by certain medical codes), the novelty of their work was in applying attention on the ontology graph to improve model performance.

## 8 DISCUSSION AND CHALLENGES

In this section we discuss additional issues and highlight important challenges for future work.

### 8.1 Attention-based methods for heterogeneous graphs

The study of heterogeneous graphs, also called heterogeneous information networks, has become quite popular in recent years with many papers published in the area [Dong et al. 2017; Lee and Adorna 2012; Sun et al. 2011, 2009a,b].

While methods like GAKE, CCM, JOINTD/E+SATT all consider knowledge graphs which is a type of heterogeneous graph, they do not differentiate between the different kinds of links and nodes. While methods such as EAGCN deal with multiple types of links they do not consider the general case where there can also be multiple kinds of nodes. GAM is another method that distinguishes, in a way, between different kinds of nodes since the attention-guided walk prioritizes the node to visit based on node type.

There is a need for attention-based methods that study how attention can be applied to general heterogeneous networks, especially taking into consideration how different kinds of meta-paths [Dong et al. 2017] can affect the learned embeddings. This is important as methods based on heterogeneous graphs have been shown to outperform methods that make the assumption that graphs only have a single type of link/edge. One can refer to [Shi et al. 2017] for a survey of heterogeneous graph-based methods.

### 8.2 Scalability of graph attention models

The majority of methods [Bahdanau et al. 2015; Ma et al. 2017; Zhou et al. 2018] that calculate an attention-based graph embedding work for relatively small graphs and may have trouble scaling effectively to larger graphs.

Methods like [Lee et al. 2018; Ryu et al. 2018; Shang et al. 2018] may be able to work on larger graphs but they still have their shortcomings. For instance, GAM [Lee et al. 2018] uses a walk to sample a graph. For large graphs, natural questions that arise include: (1) Do we need longer walks and can an RNN handle these? (2) Can a walk effectively capture all relevant and useful structures, especially if they are more complex?

Methods like [Ryu et al. 2018; Shang et al. 2018] that apply attention to a GCN architecture seem like a step in the right direction as this can be described as an attention-based end-to-end differentiable version of the Weisfeiler-Lehman algorithm [Duvenaud et al. 2015; Shervashidze et al. 2011] especially if we train stochastically [Chen et al. 2018]. However, there is a need to evaluate graph attention models on a variety of large real-world graphs to test their effectiveness. It would also be useful to explore other ways to apply attention to improve performance.

## 8.3 Inductive graph learning

Recently, there has been an interest in exploring inductive learning of graph-based problems [Guo et al. 2018; Hamilton et al. 2017] which is different from transductive learning. The former is more useful than the latter as it can generalize to yet unseen data.

This setting is important as it allows graph-based methods to handle many real-world cases where nodes and edges in graphs can appear dynamically. It also allows us to learn general patterns in one graph dataset that can be useful in another dataset much like how transfer learning is used to train a model on one text corpora for application on another text dataset [Dai et al. 2007] or to train a model to recognize general shapes in a large image dataset to be applied to a sparser dataset [Oquab et al. 2014]. Another interesting example of transfer learning is shown by [Zhu et al. 2011] where the transfer is done across data of different domains (*e.g.*, text to images).

An interesting direction for future study is looking at attention-based inductive learning techniques that can be used to identify relevant graph patterns that are generalizable to other graph datasets. Looking further, we can also explore attention-based techniques that do cross-domain or heterogeneous transfer learning [Zhu et al. 2011].

While the authors of methods like GAT [Velickovic et al. 2018] have conducted an initial study of inductive learning on a small graph dataset, we believe more focused experiments should be done on graph-based inductive learning taking into account a large set of datasets and settings.

## 8.4 Attention-guided attributed walk

Recently, Ahmed *et al.* [Ahmed et al. 2018] proposed the notion of *attributed walk* and showed that it can be used to generalize graph-based deep learning and embedding methods making them more powerful and able to learn more appropriate embeddings. This is achieved by replacing the notion of random walk (based on node ids) used in graph-based deep learning and embedding methods with the more appropriate and powerful notion of *attributed walk*. More formally,

DEFINITION 11 (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*

$$\phi(\mathbf{x}_{i_1}), \phi(\mathbf{x}_{i_2}), \ldots, \phi(\mathbf{x}_{i_{L+1}}) \tag{1}$$

*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* $\phi : \mathbf{x} \to y$ *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 the node types that occur during a walk (as opposed to the node ids). It was shown that the original deep graph learning models that use traditional random walks are recovered as a special case of the attributed walk framework when the number of unique types $t \to n$ [Ahmed et al. 2018]. It should be noted that the node types here do not necessarily refer to node types in heterogeneous graphs but can be calculated for nodes in a homogeneous graph from their local structure.

Attention-based methods that leverage random walks (based on node ids) may also benefit from the notion of *attributed walks* (typed walks) proposed by Ahmed *et al.*[Ahmed et al. 2018].

## 9 CONCLUSION

In this work, we conducted a comprehensive and focused survey of the literature on the important field of graph attention models. To the best of our knowledge, this is the first work of this kind. We introduced three intuitive taxonomies to group existing work. These are based on problem setting, the type of attention mechanism used, and the task. We motivated our taxonomies through detailed examples and used each to survey competing approaches from the taxonomy's unique standpoint.

We also highlighted several challenges in the area and provided discussion on possible directions for future work.

## REFERENCES

Sami Abu-El-Haija, Bryan Perozzi, Rami Al-Rfou, and Alex Alemi. 2017. Watch Your Step: Learning Graph Embeddings Through Attention. In *arXiv:1710.09599v1*.

Charu C. Aggarwal, Amotz Bar-Noy, and Simon Shamoun. 2017. On sensor selection in linked information networks. *Computer Networks* 126, C (2017), 100–113.

Charu C. Aggarwal and Haixun Wang. 2010. *Graph Data Management and Mining: A Survey of Algorithms and Applications*. Advances in Database Systems, Vol. 40. Springer.

Nesreen K. Ahmed, Nick Duffield, Theodore L. Willke, and Ryan A. Rossi. 2017. On Sampling from Massive Graph Streams. In *Proc. of VLDB*. 1430–1441.

Nesreen K. Ahmed, Jennifer Neville, and Ramana Kompella. 2014. Network Sampling: From Static to Streaming Graphs. *ACM TKDD* 8, 2 (2014), 1–56.

Nesreen K Ahmed, Ryan Rossi, John Boaz Lee, Xiangnan Kong, Theodore L Willke, Rong Zhou, and Hoda Eldardiry. 2018. Learning Role-based Graph Embeddings. In *arXiv:1802.02896*.

Leman Akoglu, Hanghang Tong, and Danai Koutra. 2015. Graph based anomaly detection and description: a survey. *DMKD* 29, 3 (2015), 626–688.

Reka Albert, Hawoong Jeong, and Albert-Laszlo Barabasi. 1999. Internet: Diameter of the World-Wide Web. *Nature* 401, 1 (1999), 130–131.

Stefano Allesina, Antonio Bodini, and Cristina Bondavalli. 2005. Ecological subsystems via graph theory: the role of strongly connected components. *Oikos* 110, 1 (2005), 164–176.

Lars Backstrom and Jure Leskovec. 2011. Supervised random walks: predicting and recommending links in social networks. In *Proc. of WSDM*. 635–644.

Dzmitry Bahdanau, KyungHyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *Proc. of ICLR*. 1–15.

Zilong Bai, Peter B. Walker, Anna E. Tschiffely, Fei Wang, and Ian Davidson. 2017. Unsupervised Network Discovery for Brain Imaging Data. In *Proc. of KDD*. 55–64.

HongYun Cai, Vincent W. Zheng, and Kevin Chen-Chuan Chang. 2018. A Comprehensive Survey of Graph Embedding: Problems, Techniques and Applications. *ACM TKDE* (2018), 1–20.

Jianfei Chen, Jun Zhu, and Le Song. 2018. Stochastic Training of Graph Convolutional Networks with Variance Reduction. In *arXiv:1710.10568v3*.

Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F. Stewart, and Jimeng Sun. 2017. GRAM: Graph-based Attention Model for Healthcare Representation Learning. In *Proc. of KDD*. 787–795.

Wenyuan Dai, Gui-Rong Xue, Qiang Yang, and Yong Yu. 2007. Transferring Naive Bayes Classifiers for Text Classification. In *Proc. of AAAI*. 540–545.

Shuiguang Deng, Longtao Huang, Guandong Xu, Xindong Wu, and Zhaohui Wu. 2017. On Deep Learning for Trust-Aware Recommendations in Social Networks. *IEEE TNNLS* 28, 5 (2017), 1164–1177.

Yuxiao Dong, Nitesh V Chawla, and Ananthram Swami. 2017. metapath2vec: Scalable Representation Learning for Heterogeneous Networks. In *Proc. of KDD*. 135–144.

David Duvenaud, Dougal Maclaurin, Jorge Aguilera-Iparraguirre, Rafael Gomez-Bombarelli, Timothy Hirzel, Alan Aspuru-Guzik, and Ryan P. Adams. 2015. Convolutional Networks on Graphs for Learning Molecular Fingerprints. In *Proc. of NIPS*. 2224–2232.

Jun Feng, Minlie Huang, Yang Yang, and Xiaoyan Zhu. 2016. GAKE: Graph Aware Knowledge Embedding. In *Proc. of COLING*. 641–651.

Felix A. Gers, Jurgen Schmidhuber, and Fred A. Cummins. 2000. Learning to Forget: Continual Prediction with LSTM. *Neural Computation* 12, 10 (2000), 1–20.

Lise Getoor and Christopher P. Diehl. 2015. Link Mining: A Survey. *SIGKDD Explor. Newsl.* 7, 2 (2015), 3–12.

M. Girvan and M. E. J. Newman. 2002. Community structure in social and biological networks. *PNAS* 99, 12 (2002), 7821–7826.

Aditya Grover and Jure Leskovec. 2016. node2vec: Scalable feature learning for networks. In *Proc. of KDD*. 855–864.

Junliang Guo, Linli Xu, and Enhong Cheng. 2018. SPINE: Structural Identity Preserved Inductive Network Embedding. In *arXiv:1802.03984v1*.

Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Proc. of NIPS*. 1–11.

Xu Han, Zhiyuan Liu, and Maosong Sun. 2018. Neural Knowledge Acquisition via Mutual Attention Between Knowledge Graph and Text. In *Proc. of AAAI*. 1–8.

Xinran He and David Kempe. 2014. Stability of influence maximization. In *Proc. of KDD*. 1256–1265.

Chuntao Jiang, Frans Coenen, and Michele Zito. 2013. A survey of frequent subgraph mining algorithms. *The Know. Eng. Review* 28, 1 (2013), 75–105.

Thomas N Kipf and Max Welling. 2017. Semi-supervised classification with graph convolutional networks. In *Proc. of ICLR*. 1–14.

Ankit Kumar, Ozan Irsoy, Peter Ondruska, Mohit Iyyer, James Bradbury, Ishaan Gulrajani, Victor Zhong, Romain Paulus, and Richard Socher. 2016. Ask Me Anything: Dynamic Memory Networks for Natural Language Processing. In *Proc. of ICML*. 2397–2406.

John Boaz Lee and Henry Adorna. 2012. Link Prediction in a Modified Heterogeneous Bibliographic Network. In *Proc. of ASONAM*. 442–449.

John Boaz Lee, Xiangnan Kong, Yihan Bao, and Constance Moore. 2017. Identifying Deep Contrasting Networks from Time Series Data: Application to Brain Network Analysis. In *Proc. of SDM*. 543–551.

John Boaz Lee, Ryan Rossi, and Xiangnan Kong. 2018. Graph Classification using Structural Attention. In *Proc. of KDD*. 1–9.

Jure Leskovec and Christos Faloutsos. 2006. Sampling from large graphs. In *Proc. of KDD*. 631–636.

David Liben-Nowell and Jon Kleinberg. 2003. The link prediction problem for social networks. In *Proc. of CIKM*. 556–559.

Qi Liu, Biao Xiang, Nicholas Jing Yuan, Enhong Chen, Hui Xiong, Yi Zheng, and Yu Yang. 2017. An Influence Propagation View of PageRank. *ACM TKDD* 11, 3 (2017), 30:1–30:30.

Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. Effective Approaches to Attention-based Neural Machine Translation. In *Proc. of EMNLP*. 1412–1421.

Fenglong Ma, Radha Chitta, Jing Zhou, Quanzeng You, Tong Sun, and Jing Gao. 2017. Dipole: Diagnosis Prediction in Healthcare via Attention-based Bidirectional Recurrent Neural Networks. In *Proc. of KDD*. 1903–1911.

Volodymyr Mnih, Nicolas Heess, Alex Graves, and Koray Kavukcuoglu. 2014. Recurrent Models of Visual Attention. In *Proc. of NIPS*. 2204–2212.

Mathieu Moslonka-Lefebvre, Ann Finley, Ilaria Dorigatti, Katharina Dehnen-Schmutz, Tom Harwood, Michael J. Jeger, Xiangming Xu, Ottmar Holdenrieder, and Marco Pautasso. 2011. Networks in Plant Epidemiology: From Genes to Landscapes, Countries, and Continents. *Phytopathology* 101, 4 (2011), 392–403.

Maxime Oquab, Leon Bottou, Ivan Laptev, and Josef Sivic. 2014. Learning and Transferring Mid-Level Image Representations using Convolutional Neural Networks. In *Proc. of CVPR*. 1717–1724.

Jian Pei, Daxin Jiang, and Aidong Zhang. 2005. On Mining Cross-Graph Quasi-Cliques. In *Proc. of KDD*. 228–238.

Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. DeepWalk: online learning of social representations. In *Proc. of KDD*. 701–710.

Ryan A. Rossi and Nesreen K. Ahmed. 2015. Role Discovery in Networks. *ACM TKDE* 27, 4 (2015), 1112–1131.

Seongok Ryu, Jaechang Lim, and Woo Youn Kim. 2018. Deeply learning molecular structure-property relationships using graph attention neural network. In *arXiv:1805.10988v2*.

Chao Shang, Qinqing Liu, Ko-Shin Chen, Jiangwen Sun, Jin Lu, Jinfeng Yi, and Jinbo Bi. 2018. Edge Attention-based Multi-Relational Graph Convolutional Networks. In *arXiv:1802.04944v1*.

Nino Shervashidze, Pascal Schweitzer, Erik Jan van Leeuwen, Kurt Mehlhorn, and Karsten M. Borgwardt. 2011. Weisfeiler-Lehman Graph Kernels. *JMLR* (2011), 2539–2561.

Chuan Shi, Yitong Li, Jiawei Zhang, Yizhou Sun, and Philip S. Yu. 2017. A Survey of Heterogeneous Information Network Analysis. *IEEE TKDE* 29, 1 (2017), 17–37.

Xiaoxiao Shi, Xiangnan Kong, and Philip S. Yu. 2012. Transfer Significant Subgraphs across Graph Databases. In *Proc. of SDM*. 552–563.

Yizhou Sun, Rick Barber, Manish Gupta, Charu C. Aggarwal, and Jiawei Han. 2011. Co-author Relationship Prediction in Heterogeneous Bibliographic Networks. In *Proc. of ASONAM*. 121–128.

Yizhou Sun, Jiawei Han, Peixiang Zhao, Zhijun Yin, Hong Cheng, and Tianyi Wu. 2009a. RankClus: Integrating Clustering with Ranking for Heterogeneous Information Network Analysis. In *Proc. of EDBT*. 565–576.

Yizhou Sun, Yintao Yu, and Jiawei Han. 2009b. Ranking-based Clustering of Heterogeneous Information Networks with Star Network Schema. In *Proc. of KDD*. 797–806.

Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks. In *Proc. of NIPS*. 3104–3112.

Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. 2015. LINE: Large-scale Information Network Embedding. In *Proc. of WWW*. 1067–1077.

Kiran K. Thekumparampil, Chong Wang, Sewoong Oh, and Li-Jia Li. 2018. Attention-based Graph Neural Network for Semi-supervised Learning. In *arXiv:1803.03735v1*.

Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2018. Graph Attention Networks. In *Proc. of ICLR*. 1–12.

Jia Wu, Zhibin Hong, Shirui Pan, Xingquan Zhu, Zhihua Cai, and Chengqi Zhang. 2015. Multi-graph-view subgraph mining for graph classification. *KAIS* 48, 1 (2015), 29–54.

Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron C. Courville, Ruslan Salakhutdinov, Richard S. Zemel, and Yoshua Bengio. 2015. Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. In *Proc. of ICML*. 2048–2057.

Kun Xu, Lingfei Wu, Zhiguo Wang, Yansong Feng, Michael Witbrock, and Vadim Sheinin. 2018. Graph2Seq: Graph to Sequence Learning with Attention-based Neural Networks. In *arXiv:1804.00823v3*.

Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. 2016. Stacked Attention Networks for Image Question Answering. In *Proc. of CVPR*. 21–29.

Jingyuan Zhang, Bokai Cao, Sihong Xie, Chun-Ta Lu, Philip S. Yu, and Ann B. Ragin. 2016a. Identifying Connectivity Patterns for Brain Diseases via Multi-side-view Guided Deep Architectures. In *Proc. of SDM*. 36–44.

Jingyuan Zhang, Bokai Cao, Sihong Xie, Chun-Ta Lu, Philip S. Yu, and Ann B. Ragin. 2016b. Identifying Connectivity Patterns for Brain Diseases via Multi-side-view Guided Deep Architectures. In *Proc. of SDM*. 36–44.

Zhou Zhao, Ben Gao, Vicent W. Zheng, Deng Cai, Xiaofei He, and Yueting Zhuang. 2017. Link Prediction via Ranking Metric Dual-level Attention Network Learning. In *Proc. of IJCAI*. 3525–3531.

Weiguo Zheng, Lei Zou, Xiang Lian, Dong Wang, and Dongyan Zhao. 2013. Graph similarity search with edit distance constraint in large graph databases. In *Proc. of CIKM*. 1595–1600.

Yu Zheng, Licia Capra, Ouri Wolfson, and Hai Yang. 2014. Introduction to the Special Section on Urban Computing. *ACM TIST* 5, 3 (2014), 38:1–38:55.

Hao Zhou, Tom Yang, Minlie Huang, Haizhou Zhao, Jingfang Xu, and Xiaoyan Zhu. 2018. Commonsense Knowledge Aware Conversation Generation with Graph Attention. In *Proc. of IJCAI-ECAI*. 1–7.

Yin Zhu, Yuqiang Chen, Zhongqi Lu, Sinno Jialin Pan, Gui-Rong Xue, Yong Yu, and Qiang Yang. 2011. Heterogeneous Transfer Learning for Image Classification. In *Proc. of AAAI*. 1304–1309.

Yuanyuan Zhu, Jeffrey Xu Yu, Hong Cheng, and Lu Qin. 2012. Graph Classification: A Diversified Discriminative Feature Selection Approach. In *Proc. of CIKM*. 205–214.