

# gl2vec: Learning Feature Representation Using Graphlets for Directed Networks

Kun Tu<sup>1</sup>, Jian Li<sup>1</sup>, Don Towsley<sup>1</sup>, Dave Braines<sup>2</sup>, and Liam D. Turner<sup>3</sup>

<sup>1</sup>University of Massachusetts Amherst, <sup>2</sup>IBM UK, <sup>3</sup>Cardiff University

<sup>1</sup>{kuntu, jianli, towsley}@cs.umass.edu, <sup>2</sup>dave\_braines@uk.ibm.com, <sup>3</sup>TurnerL9@cardiff.ac.uk

**Abstract**—Learning network representations has a variety of applications, such as network classification. Most existing work in this area focuses on static undirected networks and does not account for presence of directed edges or temporal changes. Furthermore, most work focuses on node representations that do poorly on tasks like network classification. In this paper, we propose a novel network embedding methodology, *gl2vec*, for network classification in both static and temporal directed networks. *gl2vec* constructs vectors for feature representation using static or temporal network graphlet distributions and a null model for comparing them against random graphs. We demonstrate the efficacy and usability of *gl2vec* over existing state-of-the-art methods on network classification tasks such as network type classification and subgraph identification in several real-world static and temporal directed networks. We argue that *gl2vec* provides additional network features that are not captured by state-of-the-art methods, which can significantly improve their classification accuracy by up to 10% in real-world applications such as detecting departments for subgraphs in an email network or identifying mobile users given their app switching behaviors represented as static or temporal directed networks.

## I. INTRODUCTION

Networks, where elements are denoted as nodes and their interactions are denoted as edges, are fundamental to the study of complex systems [2], [29], including social, communication, and biological networks. Analysis of such networks include network classification, community detection and so on. This often involves applying machine learning techniques to these problems, which requires the network to be represented as a feature vector. However, representing a network is challenging due to high dimensionality and network structure.

Various ways of learning *feature representations* of nodes in networks have been recently proposed to exploit their relations to vector representations [1], [14], [28], [37], [40]. However, most of these are applied to node and edge predictions and fail to fully capture *network structures*. It is still unclear if the result of network classification by these node embedding methods can be improved, since the whole network structure also plays a significant role. Furthermore, typical analysis usually models these systems as static undirected graphs that describe relations between nodes. However, in many realistic applications, these relations are directional and may change over time [17], [20], [35]. Modeling these *directed and temporal* properties is of additional interest as it can provide a richer characterization of relations between nodes in networks.

In this paper, we address the aforementioned issues by proposing a novel network embedding methodology, *gl2vec*,

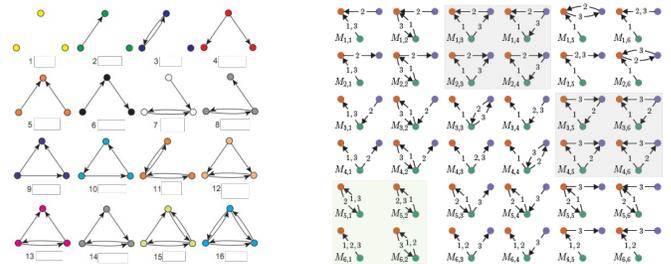


Fig. 1: (Left): All 16 triads in static directed network [11]; (Right): All 2-node and 3-node, 3-edge,  $\delta$ -temporal graphlets as defined in [35]. Edge labels correspond to the ordering of edges. All 36 graphlets are labeled with  $M_{i,j}$  across 6 rows and 6 columns. The first edge in each graphlet is from the green to the orange node. The second edge is the same along each row, and the third edge is the same along each column.

for network classification in both static and temporal directed networks. *gl2vec* constructs vectors for *feature representations* by comparing *static or temporal network graphlet statistics in a network to random graphs generated from different null models* (subgraph ratio profile, i.e., SRP, see Section III). Graphlets are small non-isomorphic induced subgraphs representing connected patterns in a network and their frequency can be used to assess network structures. For example, Figure 1 shows triads (Left), and all possible 2-node and 3-node, 3-edge,  $\delta$ -temporal graphlets (Right). These will be described in detail in Section II. We show that the ratios of occurrences of graphlets in a network to their occurrences in random graphs (SRPs) can be used as a fixed length feature representation to classify and compare networks of varying sizes and periods of time with high accuracy. We apply various well-known machine learning models along with our graph feature representation for network classifications, and make a comparison with state-of-the-art methods, such as different graph kernels [40], *node2vec* [14], *struc2vec* [37], *sub2vec* [1], *graph2vec* [28], for network classification. We argue that *gl2vec* provides additional network features that are not captured by state-of-the-art methods, which can significantly improve classification accuracy.

In particular, we consider two classification problems. First, we study how static and temporal network graphlets can be used to classify the *network type*. A network type is defined as the network domain [35] that a network belongs to, e.g., email networks, Google+ or Twitter in social networks, question

answering networks, or even networks representing switching between mobile apps. Graphs or subgraphs from the same network type often have similar structures [26]. Identifying the type of network further allows us to study interactions between nodes, and predict unobserved network structures. Secondly, we consider the problem of identifying a particular (sub)network within the same network type from its static or temporal topological structure. For example, we predict the community ID for (sub)graphs within a network, such as identifying a department based on the temporal email-exchange pattern or detecting a mobile phone user given their app switching behaviors represented as static or temporal networks.

Given a network topological structure, identifying the network type or a network community ID in a network can be viewed as a (sub)graph classification problem. Many existing methods use different graph embedding techniques to represent graphs in a vector space and apply machine learning methods for classification. Yet, little work has applied network graphlets to real-world application in directed (temporal) network classification. In this paper, we find that a strong relation exists between network type, graphlet distribution and subgraph ratio profile (SRP).

Highlights of our contributions include:

- 1) We propose a novel graphlet feature representation method, *gl2vec*, for network classification in both static and temporal directed networks.
- 2) We empirically evaluate *gl2vec* against state-of-the-art methods on tasks such as network type classification and subgraph identification in several real-world static and temporal datasets. We find that *gl2vec* outperforms state-of-the-art methods in these two tasks.
- 3) More importantly, we show that when *gl2vec* is concatenated with state-of-the-art methods, the concatenation provides a significant improvement on classification accuracy in real-world applications from several domains. This indicates that *gl2vec* provides important network features not captured by state-of-the-art methods.

The rest of the paper is organized as follows. Firstly, we present explicit formulation of the problem in Section II. From this, we present *gl2vec* in Section III and evaluate it in Section IV, followed by a discussion on related work in Section V and a conclusion of our findings in Section VI.

## II. PROBLEM FORMULATION

In this section, we provide definitions used in the rest of the paper and formulate the problem being addressed. For temporal networks and temporal network graphlets, we consider definitions given in [35], although we can equally use definitions in [20]. We present them here for completeness.

**Definition 1.** A temporal directed network [35] is a set of nodes and a collection of directed temporal edges with a timestamp on each edge. Formally, a temporal directed network  $T$  on a set of nodes  $V$  is a collection of tuples  $(u_i, v_i, t_i)$ ,  $i = 1, \dots, N$ , where  $N$  is the number of directed

temporal edges,  $u_i, v_i \in V$  and  $t_i \in \mathbb{R}$  is a timestamp. We refer to  $(u_i, v_i, t_i)$  as a temporal edge.

In order to strictly order the tuples, we assume timestamps  $t_i$  are unique. This assumption can be easily extended to cases where timestamps are not unique at the cost of complex notation.

**Definition 2.** A directed static network  $G(V, E)$  is defined as a set of nodes, denoted as  $V$  and a set of directed edges without timestamps, denoted as  $E \subset V^2 \setminus \{(u, u) : u \in V\}$ .

In the following, we formalize the definitions of (static) graphlet and temporal graphlet.

**Definition 3.** Graphlets are small connected non-isomorphic induced subgraphs of a larger network.

In particular, we focus on triads, shown in Figure 1 (Left). Note that the first three triads are not connected, hence do not satisfy the graphlet definition, but we argue that they are also important in constructing vectors for network feature representation.

**Definition 4.** Temporal network graphlets [35] are defined as induced subgraphs on sequences of temporal edges. Formally, a  $k$ -node,  $l$ -edge,  $\delta$ -temporal graphlet is a sequence of  $l$  edges,  $M = (u_1, v_1, t_1), \dots, (u_l, v_l, t_l)$  that are time-ordered within a duration  $\delta$ , i.e.,  $t_1 < \dots < t_l$  and  $t_l - t_1 \leq \delta$ , such that the induced static graph from edges is connected with  $k$  nodes.

We consider all 2-node and 3-node, 3-edge,  $\delta$ -temporal graphlets, as shown in Figure 1 (Right). Note that [35] used the term network motif.

### A. Problem Formulation

Next, we formulate our problem, which applies to the tasks of network type classification and subgraph identification.

Denote  $\{G_i(V_i, E_i, L_i)\}_{i=1}^N$  as (sub)graphs in different static or temporal networks, where  $V_i$  is a set of nodes and  $E_i$  is a set of edges in  $G_i$ . If  $G_i$  is a temporal network,  $E_i$  is then a temporal edge with a timestamp as defined in Definition 1, otherwise,  $E_i$  is a directed edge. Suppose that graphs can be categorized into  $D$  classes,  $D < N$ . We associate each graph  $G_i$  with a label  $L_i \in \{1, \dots, D\}$ .

Let  $f : \{G_i\} \rightarrow \mathbb{R}^m$  be a mapping function (also called graph embedding function) from  $G_i$  to a  $1 \times m$  feature representation vector defined using SRPs of static or temporal graphlets. We formally define SRP in Section III.

Let  $g : \mathbb{R}^m \rightarrow P \in \mathbb{R}^D$  be a classifier that maps a feature representation to a categorical distribution  $P$  for  $D$  labels. We represent probability distribution of  $G_i$ 's label as  $P_i = [p_{i,1}, \dots, p_{i,D}] = g(f(G_i))$ .

Our goal is to solve this classification problem by designing an embedding function  $f$  and selecting a machine learning model  $g$  that minimizes the sum of cross entropy [9] for all graphs

$$\arg \min_{g, f} \left( - \sum_i \sum_{j=1}^D \mathbf{1}_{L_i=j} \log(p_{i,j}) \right) = \arg \min_{g, f} \left( - \sum_i \log(p_{i,L_i}) \right).$$

We obtain  $g$  by training machine learning models. In the next section, we discuss how to design an embedding function  $f$  for static and temporal networks using graphlets.

### III. NETWORK EMBEDDING USING GRAPHLET

Network embedding has received considerable attention due to its effect on the performance of network classification, see Section V. However, previous work has primarily focused on examining this for undirected static networks. Applying these techniques to directed static networks may lose network structure information, while applying them to temporal networks loses temporal information, and both may result in poor accuracy. Therefore, we introduce a new static (temporal) network embedding technique based on static (temporal) network graphlets.

Graph embeddings need to be independent of network size and, if temporal, the time period the network covers. While previous work has shown that the counting and probability distribution of graphlets are strongly related to network types [35], graphlet counts may differ across networks. Instead, we use subgraph ratio profile (SRP) for network embedding, which is computed using graphlet counts from both the network in question and random graphs produced using a null model.

**Definition 5.** A null model [31] is a generative model used to generate random graphs that matches a specific graph in some of its structural features such as the degrees of nodes or number of nodes and edges.

For static networks, we consider the null model for random graphs with the same number of nodes and edges ( $NE$ ).  $NE$  has been widely used in previous studies since it is easy to generate random graphs [22] and the probability of a node degree in a random graph can be approximated by Poisson distribution in the large limit of graph size [32]. Thus network features and graphlet statistics can be easily modeled<sup>1</sup>.

For temporal networks, since there is no equivalent null model, we consider ensembles of randomized time-shuffled data as a temporal null model [25]. To be more specific, we randomly permute the timestamps on the edges while keeping the node pairs fixed. This model breaks the temporal dependencies between edges but preserves the network structure.

In our study, we use a null model to compare graphlet counts in a network against random graphs. The difference between counts is then used to construct an SRP as a feature representation of the network.

**Definition 6.** Subgraph ratio profile (SRP) [26] for a graphlet  $i$  is defined as

$$SRP_i = \frac{\Delta_i}{\sqrt{\sum \Delta_i^2}}, \quad (1)$$

<sup>1</sup> [4], [32] showed node degrees in a wide range of real-world networks do not necessarily follow a Poisson distribution and suggested a null model with controlled node degree sequence for network study. Thus, we consider other variants, and numerically show that their impacts on network classification accuracy is negligible. Hence are omitted here due to space constraints.

where  $\Delta_i = \frac{N_{ob_i} - \langle N_{rand_i} \rangle}{N_{ob_i} + \langle N_{rand_i} \rangle + \epsilon}$ . Here  $N_{ob_i}$  is the count of graphlet  $i$  observed in an empirical network, and  $\langle N_{rand_i} \rangle$  is the the average count in random networks in a null model. Last,  $\epsilon$  (usually set to four) is an error term to make sure that  $\Delta_i$  is not too large when a graphlet  $i$  rarely appears in both empirical and random graphs.

A large positive value of an SRP indicates that a graphlet occurs much more frequently in a network than would be expected by random chance. Since SRP for a graphlet has been normalized, it can be used to compare different size networks. The network embedding is a vector containing 16 SRPs for static triads. For null models of temporal directed networks, we randomly order of temporal edges. The embedding contains the SRPs for the 36 temporal graphlets illustrated in Figure 1 (Right).

#### A. Algorithm

$gl2vec$  works as follows: given the topological structure of a directed static or temporal network, we first compute its graphlet counts. For static networks, we applied JMotif [43] to compute triad counts for networks and random graphs in different null models. We refer interested readers to [43] for more details. For temporal networks, we use the SNAP package [35] to compute 3-edge,  $\delta$ -temporal graphlet counts.

Then we compute average graphlet counts in null models  $NE$ . For static networks, there are two approaches: simulation based and probability based. The simulation based approach generates a large set of random graphs with the same structure of the given network and a graphlet counts are computed for each random graph. The probability based approach computes the probability of occurrence for each type of graphlet given the in/out degree of the nodes involved. We apply the probability based approach to  $NE$  due to its fast computation speed and high accuracy. For temporal networks, we generate random graphs by shuffling timestamps on edges and then compute their average temporal graphlet counts. Finally, we compute SRPs for corresponding graphlets using (1). The pseudocode is presented in Algorithm 1.

---

#### Algorithm 1: $gl2vec$

---

**Data:** Static or temporal graph edges list  $E$ ,  
Null model  $M$

**Result:** Graph feature vector  $\vec{f}$

- 1  $\vec{N}_{ob} = \text{getGraphletCounts}(E)$  ;
  - 2  $\vec{N}_{rand} = \text{getAvgGraphletCountsInNullModel}(E, M)$
  - 3 **for**  $i = 1 : |\vec{N}_{obs}|$  **do**
  - 4      $\vec{f}_i = \text{getSRP}(\vec{N}_{obs_i}, \vec{N}_{rand_i})$
  - 5 **return**  $\vec{f}$
- 

### IV. EXPERIMENTS

In this section we conduct network classification on several real-world static and temporal directed networks. Experiments include two tasks: network type classification and subgraph

identification. In network type classification, we use *gl2vec* to predict the most likely relation and interaction between nodes, e.g., email communication, question answering or friendship in social networks. In subgraph identification, we predict the community ID for (sub)graphs within the same network. Examples include identifying a department based on email-exchange patterns or detecting a mobile phone user based on their app switching behavior represented as static or temporal networks.

Highlights of our experimental findings include:

- 1) *gl2vec*, constructing vectors for feature representations using static or temporal graphlet SRPs, can significantly outperform state-of-the-art methods in network type classification.
- 2) Adding graphlet features from *gl2vec* to state-of-the-art-methods significantly improves their performance. This suggests that graphlet patterns from SRPs provide substantial information about network type that do not exist in state-of-the-art-methods.
- 3) Both static and temporal graphlets play important roles in temporal network classification.

#### A. Datasets

We use a wide range of real-world network datasets, which only contain topological structure. Attributes of nodes and edges are unknown, except for labels for classification, and timestamps of edges, in temporal networks. These datasets may challenge some current state-of-the-art methods that require attributes of nodes or edges.

1) *Static directed network datasets*: We use different types of static directed networks and perform network classification using their topological structures in our experiments.

**SNAP datasets [24]**: For social Networks, *Twitter dataset* contains 1000 ego-networks with 81,306 nodes and 1,768,149 edges. *Google+ dataset* contains 133 ego-networks with 106,674 nodes and 13,673,453 edges. A directed edge from  $u$  to  $v$  represents that user  $u$  follows  $v$ . The size of ego-networks range from 10 to 4,964 nodes.

*Askubuntu* and *Mathoverflow* datasets are question-answering networks that store interactions between users. The interactions include posting answers to question (a2q), comments to questions (c2q) and comments to answers (c2a). Both datasets contain four directed networks: an a2q network, a c2q network, a c2a network and a network containing all interactions. *Askubuntu (Mathoverflow)* contains 159,316 (24,818) nodes and 596,933 (239,978) edges.

*p2p-Gnutella dataset* contains 9 directed peer-to-peer file sharing networks with 6,301 to 62,586 nodes and 20,777 to 147,892 edges. Nodes represent hosts and edges represent topological connections between hosts.

*Cit-HepPh* and *Cit-HepTh* are two physics paper citation networks. *Cit-HepPh (Cit-HepTh)* contains 34,546 (27,770) nodes and 421,578 (352,807) edges.

*Slashdot* is a friendship network with 77,360 nodes and 905,468 edges, where users tag each other as friends. *WikiVote dataset* contains votes from users in Wikipedia to promote

other users to become administrators. There are a total of 7,115 nodes and 103,689 edges. *Bitcoin OTC trust weighted signed network* contains ratings from Bitcoin users to other users and contains 5,881 nodes and 35,592 edges.

**Other Network Types**: *Epinion social network* [38] is a who-trusts-whom network from a consumer review site Epinions.com containing 75,879 nodes and 508,837 edges. A directed edge represents that a user “trusts” another user. Advice dataset contains advice-seeking between employees in four different companies [7], [21], [23]. Network sizes range from 30 to 60 with number of edges ranging from 200 to 500. Co-sponsorship networks [12] contain US Senate co-sponsorship patterns during the 1995, 2000, 2005, and 2010 congressional terms. Nodes represent senators and a directed edge from  $u$  to  $v$  represents that senator  $u$  cosponsored at least one piece of legislation for which senator  $v$  was the primary sponsor.

2) *Temporal directed network datasets*: We also collect temporal directed networks to test feature representation using temporal graphlets.

**Email Networks**: EmailEU [47] is a directed temporal network constructed from email exchanges in a large European research institution for a 803-day period. It contains 986 email addresses as nodes and 332,334 emails as edges with timestamps. There are 42 ground truth departments in the dataset and we choose 26 departments whose email network sizes are larger than 10. EmailTraffic [34] is a temporal directed network storing email interactions of 819 staff in 23 different departments in BBN for about 7 months. Edges with integer timestamps represent emails sent out at a certain time.

We constructed temporal subgraphs, each lasting 12 weeks for departments in EmailEU networks. This ensures each subgraph becomes a connected network component when converted to an unweighted static graph. We create these graphs at the beginning of every four weeks to avoid too much overlap of edges between graphs. Each department has up to 28 subgraphs as a result. For departments in EmailTraffic, we create subgraphs at the beginning of every week and each subgraph covers four weeks.

**SwitchApp**: (from the Tymer project [44], [45]) contains application switching data for 53 Android users over a 42-day period. We construct a directed temporal network for each user on each day, where a directed edge (denoted as  $e_{uv}$ ) with an integer timestamp  $t$  represents a user switching from an app  $u$  to another  $v$  at time  $t$ .

#### B. Experiment Setup

We compute SRPs for static and temporal graphlets for corresponding static and temporal networks in our datasets. We use three widely used machine learning models that provide good performance using small amounts of training data in multi-class classification: XGBoosting [5], SVM [6], random forest [41]. XGBoosting usually has a superior performance over other classifiers when the dataset is of middle size. SVM is suitable for a small amount of training data. Random forest not only works well for imbalanced data, but also performs

feature selection during training which can help us investigate the usefulness of our feature representation, especially when used in conjunction with other approaches by concatenating the feature vectors.

We use grid search method to search the best hyper-parameters for these models. For XGBoosting algorithm, the learning rate ranges from 0.001 to 1, maximal tree depth range from 4 to 32, minimal child weight is 1 and the subsample ratio of train instances ranges from 0.4 to 1. The regularization weight in SVM ranges from 1 to 8. In random forest, the number of trees ranges from 50 to 400 and the minimal number of samples required to split a tree node from 2 to 10. 10-fold cross-validation is adopted to split the data to select the best parameters. All experiments are conducted using a cluster with 32 Xeon CPU with 256Gb RAM and one Tesla K40 GPU.

We compare the network classification accuracy of *gl2vec* to state-of-the-art methods, including graphlet and Weisfeiler-Lehman kernels [40], and recently developed node and graph embedding methods *node2vec* [14], *struc2vec* [37], *sub2vec* [1], *graph2vec* [28].

For node embedding methods such as *node2vec* and *struc2vec*, we apply sum-based approach [8] to aggregate node embedding vectors to construct a graph embedding. We refer interested readers to [16] for more detail. The length of network embedding (ranging from 50 to 500) is determined using grid search and 10-fold cross-validation. We modify state-of-the-art methods to apply them to directed graphs: we run a random walk on directed graphs in *sub2vec* instead of undirected graph. Some state-of-the-art methods also require node attributes for network embeddings and node degree are suggested for computing undirected graph embedding [16]. For directed networks, we use NetworkX to compute the in/out degree and centralities such as betweenness, closeness and in/out degree centrality for each node. We also consider additional attributes: counts of subgraphs of a specific triad that a node belongs to. These counts are normalized as a distribution indicating the likelihood a node belongs to a specific triad.

### C. Network Types Classification

In network type classification, we are given the topological structure of a subgraph in a network. Our goal is to predict the type of interaction that an edge represents, e.g. email exchange or question answering.

Among all the datasets introduced in Section IV-A, *EmailEU*, *EmailTraffic* and *SwitchApp* datasets have ground truth labels (department ID or user ID) available for each subgraph, which is created from email exchanges in a department or app switch behaviors of a user within a period of time. Hence, we can obtain all subgraphs for these communities in these three networks. For the other datasets, there is no ground truth information on network communities; we detect network communities using modularity [30] to obtain subgraphs. These subgraphs are converted into feature vectors using the previously introduced embedding methods and assigned labels according to network types. Finally, we collect about 10,000

	XGBoost (%)	SVM (%)	RF (%)
GK Graphlet +gl2vec	<b>78.94 ± 3.18</b> 82.18 ± 2.86	72.66 ± 2.79 69.01 ± 2.27	78.72 ± 3.01 81.39 ± 3.36
GK WL +gl2vec	78.26 ± 2.65 82.54 ± 2.85	72.81 ± 2.74 68.59 ± 2.75	<b>78.41 ± 3.02</b> 82.26 ± 3.43
MotifDist +gl2vec	<b>78.08 ± 3.34</b> 81.75 ± 3.48	71.40 ± 2.29 69.70 ± 3.64	78.01 ± 3.56 80.95 ± 3.63
node2vec +gl2vec	<b>74.25 ± 3.07</b> 88.76 ± 1.26	69.03 ± 1.23 73.24 ± 2.92	72.24 ± 1.67 86.14 ± 1.71
graph2vec +gl2vec	72.48 ± 3.99 79.83 ± 4.59	70.81 ± 3.84 66.70 ± 4.04	<b>72.61 ± 3.36</b> 80.03 ± 4.38
sub2vec +gl2vec	<b>81.39 ± 1.70</b> 92.30 ± 2.29	79.69 ± 1.41 83.16 ± 2.62	78.44 ± 2.26 90.01 ± 2.16
struc2vec +nodeTriadDistr +gl2vec	<b>79.15 ± 3.42</b> 81.93 ± 3.53 93.38 ± 1.51	78.22 ± 3.15 79.18 ± 3.55 84.25 ± 0.82	78.94 ± 3.31 82.01 ± 3.42 93.48 ± 1.42
gl2vec	<b>81.58 ± 3.07</b>	71.64 ± 2.13	79.42 ± 3.69

TABLE I: Network type classification accuracy. We use “+” to denote an embedding generated by combining two embedding methods. Bold indicated best performance machine learning model for each embedding.

(sub)graphs from 2,355 real-world networks taken from 15 network types introduced above, which include Google+ and Twitter in social networks, high energy physics theory citation networks, Gnutella P2P networks, SwitchApp and so on.

1) *Static Directed Network*: We use all datasets to evaluate embedding methods on static networks. Note that we convert temporal networks into unweighted static networks by removing the timestamps on the edges. Baseline methods include graphlet graph kernel (GK graphlet), Weisfeiler-Lehman graph kernel (GK WL), feature vector with triad distribution (MotifDist), *node2vec*, *graph2vec*, *sub2vec* and *struc2vec*.

The accuracies of different embedding methods for network type classification are presented in Table I. We make the following observations:

- The graph-based network embedding methods, GK Graphlet, *sub2vec*, *gl2vec*, and *struc2vec* with added subgraph features (triad distribution for a node), have a larger average accuracy compared to other node-based network embedding methods. This further validates the importance of including subgraph information into feature representations for tasks like network classification in which network structure plays a significant role.
- The machine learning methods used also have an impact on the results. For this task, XGBoost provides the best performance on average in network type classification. Although *sub2vec* is robust across all three machine learning models, *gl2vec* achieves the highest accuracy and we can always choose the trained model with the highest accuracy for prediction.
- We also combine *gl2vec* with state-of-the-art methods by directly concatenating their feature representation vectors. We observe a significant improvement on state-of-the-art methods, especially for *sub2vec* and *struc2vec*. This suggests that both our approach and state-of-the-art methods capture important but different features for network type classification. The best approach for the problem is to combine those features. Furthermore, there are also im-

provements on MotifDist and GK Graphlet. This indicates that adding null models to construct feature representation helps improve performance. Since representations from  $gl2vec$ , MotifDist and GK Graphlet construct features from graphlets, the improvement is not as significant as other methods.

2) *Temporal directed network*: We consider the temporal datasets discussed in Section IV-A. We explore if temporal graphlets provide more information than static graphlets in temporal networks. We investigate their effect on predicting whether a temporal (sub)graph is an email exchange network or the app switching behavior of a mobile user. The results are shown in Figure 2. From Figure 2, we observe that temporal information improves network type classification in all models considered here. Therefore, it is important to use temporal graphlets for constructing vectors for feature representations of temporal networks, since temporal graphlets provides more network structure information than static graphlets.

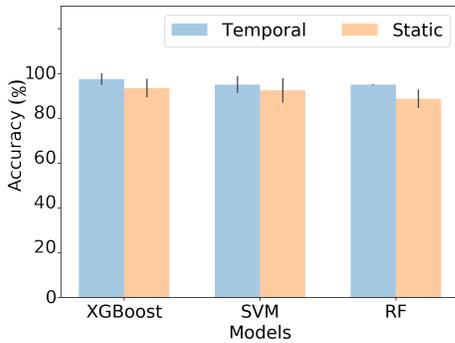


Fig. 2: Classifying email datasets and SwitchApp Temporal Networks.

#### D. Subgraph Identification

In subgraph identification, we are interested in classifying subgraphs within the same network given their topological structure. For example, we can identify which department an email exchange subgraph belongs to or detect a mobile phone user given their app switching behavior.

We use EmailEU, EmailTraffic and SwitchApp datasets since ground truth labels (department ID or user ID) are available for each subgraph, which is created from email exchanges in a department or app switch behavior of a user within a period of time. We first solve this problem using static graph embedding methods. Then we investigate whether the timestamp information of edges can help improve identification accuracy.

1) *Static directed networks*: The results on the accuracy of identifications of departments in emailEu, emailTraffic networks and user ID in app switch network using different methods are illustrated in Tables II, III and IV, respectively. We cannot obtain results from graph2vec due to its insufficient memory in GPU. We also evaluate  $gl2vec$  when combined with state-of-the-art methods. We use “+” to denote these com-

	XGBoost (%)	SVM (%)	RF (%)
<b>MotifDistr</b>	56.68 ± 6.70	45.82 ± 7.38	<b>61.54 ± 10.50</b>
+gl2vec	64.18 ± 6.52	52.20 ± 4.80	63.79 ± 8.94
<b>GK WL</b>	50.96 ± 8.91	47.92 ± 6.15	<b>57.01 ± 6.91</b>
+gl2vec	63.12 ± 5.44	51.95 ± 4.44	65.29 ± 8.81
<b>GK Graphlet</b>	61.22 ± 4.70	52.32 ± 5.16	<b>62.90 ± 4.49</b>
+gl2vec	62.04 ± 5.69	52.22 ± 4.85	64.35 ± 8.86
<b>node2vec</b>	52.08 ± 3.00	57.76 ± 3.11	<b>57.89 ± 2.83</b>
+gl2vec	63.20 ± 3.69	59.20 ± 5.89	63.22 ± 3.40
<b>sub2vec</b>	55.45 ± 3.42	52.02 ± 3.29	<b>59.87 ± 3.77</b>
+gl2vec	73.01 ± 8.93	58.88 ± 9.42	77.69 ± 6.90
<b>struc2vec</b>	60.25 ± 9.40	56.8 ± 11.34	<b>60.59 ± 11.14</b>
+nodeTriadDistr	60.78 ± 9.13	59.86 ± 9.22	61.24 ± 9.88
+gl2vec	69.78 ± 6.20	54.91 ± 9.14	70.30 ± 7.36
<b>gl2vec</b>	61.72 ± 3.09	51.03 ± 3.30	<b>63.09 ± 3.23</b>

TABLE II: Accuracy in correctly identifying 26 EmailEU department in static directed networks.

	XGBoost (%)	SVM(%)	RF (%)
<b>MotifDistr</b>	67.81 ± 7.60	62.60 ± 8.27	<b>70.04 ± 7.48</b>
+gl2vec	78.81 ± 10.87	71.93 ± 6.78	80.19 ± 7.73
<b>GK WL</b>	72.18 ± 5.86	70.73 ± 6.81	<b>75.99 ± 5.82</b>
+gl2vec	77.96 ± 9.03	71.73 ± 7.45	80.58 ± 7.24
<b>GK graphlet</b>	74.39 ± 10.71	70.77 ± 12.35	<b>78.61 ± 8.91</b>
+gl2vec	77.17 ± 11.73	71.52 ± 6.71	80.18 ± 7.23
<b>node2vec</b>	74.02 ± 8.13	70.45 ± 12.25	<b>77.41 ± 6.93</b>
+gl2vec	85.21 ± 6.64	75.36 ± 6.88	87.81 ± 4.93
<b>sub2vec</b>	<b>77.79 ± 3.93</b>	77.80 ± 3.63	77.01 ± 3.83
+gl2vec	83.39 ± 5.43	86.74 ± 5.25	87.00 ± 4.58
<b>struc2vec</b>	<b>73.78 ± 9.40</b>	65.33 ± 9.34	72.16 ± 9.14
+nodeTriadDistr	74.35 ± 10.72	66.84 ± 10.08	77.17 ± 10.44
+gl2vec	79.85 ± 14.01	56.23 ± 14.88	81.43 ± 12.38
<b>gl2vec</b>	76.80 ± 6.24	71.13 ± 6.49	<b>80.78 ± 5.65</b>

TABLE III: Accuracy in correctly identifying EmailTraffic department in static directed networks.

binations. For example, MotifDistr+gl2vec combines feature vectors from MotifDistr and  $gl2vec$  for feature representation.

We notice that the addition of graphlet SRP features to state-of-the-art methods can significantly improve performance of the corresponding state-of-the-art methods. This indicates that our  $gl2vec$  provides new information not present in state-of-the-art methods.

2) *Algorithm performance with graphlet features* : One observes from Tables II, III and IV that random forest (RF) is usually more accurate for graph embeddings that include

	XGBoost (%)	SVM (%)	RF (%)
<b>MotifDistr</b>	11.82 ± 2.03	11.62 ± 2.02	<b>12.33 ± 2.28</b>
+gl2vec	16.16 ± 1.85	12.95 ± 2.31	15.34 ± 1.45
<b>GK WL</b>	11.50 ± 1.65	<b>14.59 ± 0.97</b>	13.43 ± 2.26
+gl2vec	16.01 ± 2.43	13.15 ± 1.44	17.31 ± 1.81
<b>GK graphlet</b>	13.89 ± 1.26	15.24 ± 1.67	<b>15.29 ± 2.28</b>
+gl2vec	16.52 ± 1.77	13.98 ± 2.51	15.95 ± 1.61
<b>node2vec</b>	<b>10.15 ± 1.50</b>	7.91 ± 1.32	9.98 ± 1.66
+gl2vec	16.33 ± 2.04	12.94 ± 2.71	16.21 ± 1.97
<b>sub2vec</b>	16.27 ± 2.20	16.19 ± 4.37	<b>16.54 ± 1.64</b>
+gl2vec	31.74 ± 3.58	23.43 ± 2.33	33.94 ± 4.58
<b>struc2vec</b>	<b>14.18 ± 2.21</b>	9.75 ± 2.49	12.17 ± 2.64
+nodeTriadDistr	15.23 ± 2.57	7.81 ± 2.02	13.16 ± 2.09
+gl2vec	19.53 ± 3.13	9.30 ± 1.60	20.70 ± 3.07
<b>gl2vec</b>	16.17 ± 1.80	13.56 ± 1.60	<b>16.82 ± 1.31</b>

TABLE IV: Accuracy in correctly identifying 53 SwitchApp user in static directed networks.

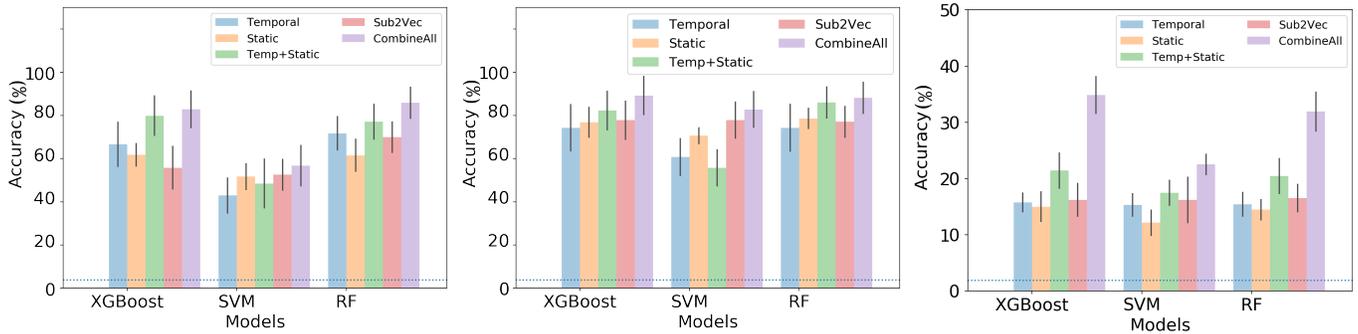


Fig. 3: (Left): Department identification in EmailEU dataset; (Middle): BBN department identification in EmailTraffic; (Right): User identification in SwitchApp. Dashed line represents the accuracy of a random selection model.

our SRP feature vectors in baselines. This is because RF automatically performs feature selection during training and adapts to the change in number of features. As a result, it is easier for RF to achieve better results given a similar amount of effort fine-tuning the hyper-parameters. Finally, the improvements on all machine learning models confirm that it is worth combining our graph embedding with other methods to achieve better performance.

3) *Temporal directed networks*: In temporal networks from EmailEU and EmailTraffic, we attempt to identify which department emails belong to. For the SwitchApp dataset, we attempt to identify a particular user based on their daily app switching behavior represented as a temporal network.

For the EmailEU and EmailTraffic dataset, multiple temporal and static networks are constructed for each department from email exchanges as described in Section IV-A2. For the SwitchApp dataset, 42 temporal and static networks are generated for each person from their app switching behaviors every day. XGBoosting, SVM and random forest are implemented using different network feature representations: subgraph ratio profile (SRP) with temporal (“Temporal”) and with static (“Static”) graphlets, combined SRPs with both temporal and static graphlet (“Temp+Static”). We illustrate the result from sub2vec representation (“Sub2Vec”) because it performs best among the baseline methods. Finally, we create a combination of all three representations (“CombineAll”).

The results for EmailEU, EmailTraffic and SwitchApp are shown in Figure 3. The dashed line is the accuracy of a random selection model. The accuracy achieved by temporal graphlet embedding is slightly better than that of static graphlet embedding in both emailEU and SwitchApp datasets. However, static graphlet embedding performs better than temporal graphlets in EmailTraffic dataset. This shows that static graphlets are still useful for temporal network classification and can capture useful features even better than temporal graphlets in some datasets. Hence, we combine both static and temporal graphlet features (“Temp+Static”) and observe that this achieves a significant improvement in accuracy, which suggests that both temporal and static graphlets are useful for network identification (of departments or personal app switching behavior). Furthermore, our graphlet-based network

embeddings are competitive with the state-of-the-art method, *sub2vec*. Finally, combining all three graph embedding vectors for classification yield the best accuracy. This suggests that both our static and temporal embedding approaches capture useful features to boost the performances of state-of-the-art methods. We also find that the accuracy for SwitchApp is much lower than other two datasets. This is because all users have structurally similar app switch networks [45].

## V. RELATED WORK

The primary focus of related works in classifying networks involves examining the topological structure of the graph. The work most related to our method is graph kernel, which has been used to calculate similarities between static undirected graphs [13], [18], [46]. However, the corresponding computational complexity grows significantly with increase in network size. Moreover, studies in graphlet kernel do not consider features generated by comparing graphlet count between an empirical network and random graphs from different null models, which turn out to lead to a significant improvement in network classification in our experiments.

Different node embedding techniques have been proposed in recent years, such as node2Vec [15], DeepWalk [36], Line [42] and Local Linear Embedding [39] that use feature vectors to embed nodes into high-dimensional space and empirically perform well. However, these methods can only be applied to node classification but not graph classification. Graph neural network (GCN) [10], [19] recently obtain competitive results against kernel-based methods and graph-based regularization techniques, but they are computationally expensive and used for small scale tasks.

Additionally, several approaches have been proposed to aggregate node feature vectors to a feature vector for networks. For example, graph-coarsening approach [10] computes a hierarchical structure containing multiple layers, nodes in lower layers are clustered and combined as node in upper layers using element-wise max-pooling. However, this has high computational complexity. Some approaches [33] define an order of nodes and concatenate their feature vectors for a convolutional neural network for classification, however, this can only be applied to undirected static networks. Recently,

some subgraph embedding based approaches were proposed. *struc2vec* [37] applied sum-based approach such as mean-field [8] and loopy belief propagation [27] to aggregate node embedding to graph representation. *sub2vec* [1] embedded subgraphs with arbitrary structure, while *graph2vec* [28] was proposed based on a *doc2vec* framework to learn data-driven distributed representations of arbitrary sized graphs. But these embedding do not fully capture network structures to the best performance. Similar to ours is [3], which uses motif frequencies, while we use graphlet distributions and SRP. Furthermore, [3] requests node labels, but we do not.

## VI. CONCLUSION

We proposed *gl2vec* to classify static and temporal directed networks based on their topological structure. Experiments with real-world datasets showed that both temporal and static graphlets are important for network type classification and subgraph identification. Furthermore, we have illustrated that concatenating these two embedding with many state-of-the-art methods yield the best accuracy for real-world applications such as identifying network types, predicting community ID for subgraphs and detecting mobile phone users based on their app-switching behaviors. In future work, we will investigate if graphlet census information can serve as features for nodes in a network. Specifically, we will investigate whether embedding nodes with the numbers of graphlets that it belongs to in a network can improve node and network classification.

## REFERENCES

- [1] B. Adhikari, Y. Zhang, N. Ramakrishnan, and B. A. Prakash. Sub2vec: Feature Learning for Subgraphs. In *PAKDD*, 2018.
- [2] R. Albert and A.-L. Barabási. Statistical Mechanics of Complex Networks. *Reviews of modern physics*, 74(1):47, 2002.
- [3] E. G. Allan Jr, W. H. Turckett, and E. W. Fulp. Using network motifs to identify application protocols. In *IEEE GLOBECOM*, 2009.
- [4] L. A. N. Amaral, A. Scala, M. Barthlmey, and H. E. Stanley. Classes of small-world networks. *PNAS*, 97(21):11149–11152, 2000.
- [5] T. Chen and C. Guestrin. Xgboost: A Scalable Tree Boosting System. In *ACM SIGKDD*, 2016.
- [6] C. Cortes and V. Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [7] R. L. Cross, R. L. Cross, and A. Parker. *The hidden power of social networks: Understanding how work really gets done in organizations*. Harvard Business Press, 2004.
- [8] H. Dai, B. Dai, and L. Song. Discriminative embeddings of latent variable models for structured data. In *ICML*, 2016.
- [9] P.-T. De Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein. A Tutorial on the Cross-Entropy Method. *Annals of operations research*, 134(1):19–67, 2005.
- [10] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering. In *NIPS*, 2016.
- [11] M. Doroud, P. Bhattacharyya, S. F. Wu, and D. Felmlee. The evolution of ego-centric triads: A microscopic approach toward predicting macroscopic network properties. In *IEEE SocialCom*, 2011.
- [12] J. H. Fowler. Connecting the congress: A study of cosponsorship networks. *Political Analysis*, 14(4):456–487, 2006.
- [13] Gaüzère, Benoit and Grenier, Pierre-Anthony and Brun, Luc and Villemin, Didier. Treelet Kernel Incorporating Cyclic, Stereo and Inter Pattern Information in Chemoinformatics. *Pattern Recognition*, 48(2):356–367, 2015.
- [14] A. Grover and J. Leskovec. node2vec: Scalable Feature Learning for Networks. In *ACM SIGKDD*, 2016.
- [15] A. Grover and J. Leskovec. node2vec: Scalable Feature Learning for Networks. In *ACM SIGKDD*, 2016.
- [16] W. L. Hamilton, R. Ying, and J. Leskovec. Representation learning on graphs: Methods and applications. *arXiv:1709.05584*, 2017.
- [17] P. Holme and J. Saramäki. Temporal Networks. *Physics reports*, 519(3):97–125, 2012.
- [18] R. Kaspar and B. Horst. *Graph Classification and Clustering based on Vector Space Embedding*, volume 77. World Scientific, 2010.
- [19] T. N. Kipf and M. Welling. Semi-supervised classification with graph convolutional networks. *arXiv:1609.02907*, 2016.
- [20] L. Kovanen, M. Karsai, K. Kaski, J. Kertész, and J. Saramäki. Temporal Motifs in Time-Dependent Networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2011(11):P11005, 2011.
- [21] D. Krackhardt. Cognitive social structures. *Social networks*, 9(2):109–134, 1987.
- [22] M. Kretzschmar and M. Morris. Measures of concurrency in networks and the spread of infectious disease. *Mathematical biosciences*, 133(2):165–195, 1996.
- [23] E. Lazega et al. *The collegial phenomenon: The social mechanisms of cooperation among peers in a corporate law partnership*. Oxford University Press on Demand, 2001.
- [24] J. Leskovec and A. Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- [25] A. Mellor. Classifying Conversation in Digital Communication. *arXiv:1801.10527*, 2018.
- [26] R. Milo, S. Itzkovitz, N. Kashtan, R. Levitt, S. Shen-Orr, I. Ayzenshtat, M. Sheffer, and U. Alon. Superfamilies of Evolved and Designed Networks. *Science*, 303(5663):1538–1542, 2004.
- [27] K. P. Murphy, Y. Weiss, and M. I. Jordan. Loopy Belief Propagation for Approximate Inference: An Empirical Study. In *UAI*, 1999.
- [28] A. Narayanan, M. Chandramohan, R. Venkatesan, L. Chen, Y. Liu, and S. Jaiswal. graph2vec: Learning Distributed Representations of Graphs. *Arxiv preprint arXiv:1707.05005*, 2018.
- [29] M. Newman. *Networks: an Introduction*. Oxford university press, 2010.
- [30] M. E. Newman. Modularity and community structure in networks. *PNAS*, 103(23):8577–8582, 2006.
- [31] M. E. Newman and M. Girvan. Finding and Evaluating Community Structure in Networks. *Physical review E*, 69(2):026113, 2004.
- [32] M. E. Newman, S. H. Strogatz, and D. J. Watts. Random graphs with arbitrary degree distributions and their applications. *Physical review E*, 64(2):026118, 2001.
- [33] M. Niepert, M. Ahmed, and K. Kutzkov. Learning Convolutional Neural Networks for Graphs. In *ICML*, 2016.
- [34] C. Olsson, P. Petrov, J. Sherman, and A. Perez-Lopez. Finding and explaining similarities in linked data. In *STIDS*, 2011.
- [35] A. Paranjape, A. R. Benson, and J. Leskovec. Motifs in temporal networks. In *ACM WSDM*, 2017.
- [36] B. Perozzi, R. Al-Rfou, and S. Skiena. Deepwalk: Online Learning of Social Representations. In *ACM SIGKDD*, 2014.
- [37] L. F. Ribeiro, P. H. Saverese, and D. R. Figueiredo. struc2vec: Learning Node Representations from Structural Identity. In *ACM SIGKDD*, 2017.
- [38] M. Richardson, R. Agrawal, and P. Domingos. Trust management for the semantic web. In *semantic Web conference*. Springer, 2003.
- [39] S. T. Roweis and L. K. Saul. Nonlinear Dimensionality Reduction by Locally Linear Embedding. *Science*, 290(5500):2323–2326, 2000.
- [40] N. Shervashidze, P. Schweitzer, E. J. v. Leeuwen, K. Mehlhorn, and K. M. Borgwardt. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(Sep):2539–2561, 2011.
- [41] V. Svetnik, A. Liaw, C. Tong, J. C. Culberson, R. P. Sheridan, and B. P. Feuston. Random Forest: a Classification and Regression Tool for Compound Classification and QSAR Modeling. *Journal of chemical information and computer sciences*, 43(6):1947–1958, 2003.
- [42] J. Tang, M. Qu, M. Wang, M. Zhang, J. Yan, and Q. Mei. Line: Large-Scale Information Network Embedding. In *WWW*, 2015.
- [43] K. Tu. Jmotif, 2018.
- [44] K. Tu, J. Li, D. Towsley, D. Braines, and L. Turner. Network classification in temporal networks using motifs. In *ECML/PKDD-AALTD*, 2018.
- [45] L. D. Turner, R. M. Whitaker, S. M. Allen, D. E. Linden, K. Tu, J. Li, and D. Towsley. Evidence to support common application switching behaviour on smartphones. *Royal Society Open Science*, 6(3), 2019.
- [46] P. Yanardag and S. Vishwanathan. Deep Graph kernels. In *ACM SIGKDD*, 2015.
- [47] H. Yin, A. R. Benson, J. Leskovec, and D. F. Gleich. Local Higher-Order Graph Clustering. In *ACM SIGKDD*, 2017.