

Learning to count: A deep learning framework for graphlet count estimation

Xutong Liu, Yu-Zhen Janice Chen, John Chi Shing Lui, Konstantin Avrachenkov

► To cite this version:

Xutong Liu, Yu-Zhen Janice Chen, John Chi Shing Lui, Konstantin Avrachenkov. Learning to count: A deep learning framework for graphlet count estimation. Network Science, Cambridge University Press, 2020, pp.30. 10.1017/nws.2020.35 . hal-02942321

HAL Id: hal-02942321

<https://hal.inria.fr/hal-02942321>

Submitted on 17 Sep 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

LEARNING TO COUNT: A DEEP LEARNING FRAMEWORK FOR GRAPHLET COUNT ESTIMATION

A PREPRINT

Xutong Liu

The Chinese University of Hong Kong
Shatin, NT, Hong Kong
liuxt@cse.cuhk.edu.hk

Yu-Zhen Janice Chen

University of Massachusetts Amherst
MA 01002, USA
yuzhenchen@cs.umass.edu

John C. S. Lui

The Chinese University of Hong Kong
Shatin, NT, Hong Kong
cslui@cse.cuhk.edu.hk

Konstantin Avrachenkov

INRIA Sophia Antipolis
Valbonne 06902, France
k.avrachenkov@inria.fr

September 17, 2020

ABSTRACT

Graphlet counting is a widely-explored problem in network analysis and has been successfully applied to a variety of applications in many domains, most notably bioinformatics, social science and infrastructure network studies. Efficiently computing graphlet counts remains challenging due to the combinatorial explosion, where a naive enumeration algorithm needs $O(N^k)$ time for k -node graphlets in a network of size N . Recently, many works introduced carefully designed combinatorial and sampling methods with encouraging results. However, the existing methods ignore the fact that graphlet counts and the graph structural information are correlated. They always consider a graph as a new input and repeat the tedious counting procedure on a regular basis even if it is similar or exactly isomorphic to previously studied graphs. This provides an opportunity to speed up the graphlet count estimation procedure by exploiting this correlation via learning methods. In this paper, we raise a novel Graphlet Count Learning (GCL) problem: given a set of historical graphs with known graphlet counts, how to learn to estimate/predict graphlet count for unseen graphs coming from the same (or similar) underlying distribution. We develop a deep learning framework which contains two *convolutional neural network* (CNN) models and a series of data *preprocessing techniques* to solve the GCL problem. Extensive experiments are conducted on three types of synthetic random graphs and three types of real world graphs for all 3,4,5-node graphlets to demonstrate the accuracy, efficiency and generalizability of our framework. Compared with state-of-the-art exact/sampling methods, our framework shows great potential, which can offer up to two orders of magnitude speedup on synthetic graphs and achieves on par speed on real world graphs with competitive accuracy.

This is a preprint of the article accepted in Network Science journal, Cambridge University Press.

Keywords Graphlet count estimation · Convolutional neural networks · Deep learning on graph · Network analysis

1 Introduction

Graphlets are defined as k -node connected induced subgraph patterns. For an undirected graph, 3-node graphlets include open triangle (\triangle) and close triangle (\blacktriangle). There are 6 different types of graphlets for $k = 4$ and 21 types of graphlets for $k = 5$, where Figure 3 and Table 1 depicts all possible 3, 4, 5-node graphlets. The number of each graphlet, called *graphlet count*, is a signature which characterizes the *local network structure* of a given graph. Global features, such as diameter, degree distribution, in contrast, are less representative, because graphs with similar global

features can have totally different local structures. Therefore, graphlet count is widely used and plays a prominent role in network analysis of many fields such as bioinformatics [1], social science [2] and infrastructure network studies [3]. In bioinformatics domain, graphlet counts were applied for protein classification and function prediction [4], disease gene identification [5], and biological network comparison [1]. In social science discipline, researchers utilized graphlet counts for analyzing social network structure [6] and interpersonal relations [2]. In infrastructure network studies, graphlet counts were adopted for anomaly detection [3] and spam detection [7].

While graphlet count is valuable for network analysis, computing exact graphlet count is inherently difficult and computationally expensive. The main challenge is the combinatorial explosion, where the count of k -node graphlets in a network grows k -th order polynomially large as the network size grows. For example, a moderate size arXiv collaboration network with 12K nodes can have more than 3.3 million triangles (\blacktriangle) and 0.2 billion 4-paths (\blacktriangleright). Experiments show that even the state-of-the-art exact counting method [8] cannot finish within a week calculating 4-node graphlet counts for a Twitter graph with 21.3M nodes [9].

To deal with this difficulty, many frameworks and algorithms were proposed to achieve a more practical goal: estimate or approximate graphlet count with acceptable error in exchange for computational efficiency. These methods [10, 11, 12, 13, 14, 15, 16] leverage on sampling techniques to estimate graphlet count with bounded error and significantly outperforms exact counting methods in space and time.

In the real world, many graph instances have some underlying structures or characteristics, and can be thought as coming from a specific distribution/domain. For example, social networks, are often scale-free networks whose degrees follow power law distributions [17]; Brain networks, have been identified as small-world networks [18]; Sub-graphs or snapshots of a given network and streaming graphs can also be viewed as coming from some distribution. Recent work [19] confirmed this fact, and build classifiers to predict (with 95.7% accuracy) the category of a given network, with only simple features including density, average degree, assortativity and etc.

Since graphlets act as fingerprints of a given graph, graph’s structural information can be naturally reflected in their graphlet patterns as well as their graphlet counts. In other words, graph structural information and graphlet counts are correlated. Moreover, this correlation could be further enhanced for graphs from the same (or similar) distribution. For instance, graphs from some biochemistry datasets, such as MUTAG, NCI1 and NCI109, only have tree-like graphlets. Nevertheless, conventional graphlet count frameworks, including both exact counting and sampling methods, do not leverage on this fact and always repeat the onerous counting procedures even if a new graph is similar to or exactly isomorphic to previously studied graphs. They neglect the connection between the structural information and the graphlet count, and merely treat them as inputs and outputs separately in the counting procedure. Ideally, this connection can be modeled as a complicated function mapping from graph structures to their graphlet counts. For instance, given a graph’s adjacency matrix A and all its eigenvalues which contains important structural information [20], the triangle count can be written as one sixth of the sum of cubes of the eigenvalues [21]. But note that, for any other graphlets, there is no literature showed the explicit form of the function. Intuitively, if we can learn to approximate this function, graphlet counts can be efficiently estimated by applying the learned function to any given graph. Essentially, we need to solve the *Graphlet Count Learning (GCL) problem*: assume we have a set of networks with known graphlet counts, how can we learn to estimate/predict graphlet counts in new-coming networks from the same (or similar) underlying distribution. Note that we assume the historical graphs and new-coming graphs are coming from the same (or similar) distribution because the function we want to learn could be related to the underlying graph distributions and has totally different properties for different distributions.

Learning paradigm has been applied in other graph problems, e.g., graph classification, community detection, etc. Yet, frameworks for learning graphlet count cannot be found in literature. The GCL problem lies in the intersection of graphlet count problem and learning paradigm. To the best of our knowledge, the problem described above has never been studied in both graphlet counting and computational learning literature. We are the *first* to study this problem and our research goal is to develop a learning framework to solve the GCL problem.

Proposed Work. To achieve this goal, we adopt a powerful model in deep learning literature as the core component of our learning framework, namely the Convolutional Neural Network (CNN). CNN is a class of deep feed forward artificial neural networks which has become an important tool in addressing a wide variety of machine learning problems, including image recognition [22], video recognition [23] and natural language processing [24]. CNN turns out to be very good at discovering intricate structures and offers attractive features including significant representation power, great learning capacity, constant testing speed and efficiency in space. However, since CNN is natively developed for regular, low-dimensional grid-structured data such as texts, images, and videos; applying CNN to graph data for graphlet count estimation is far from trivial. In this work, we first propose a straight-forward 2-D CNN model which applies vanilla CNN on the adjacency matrices of input graphs. Then we point out and analyze some drawbacks of the 2-D CNN and propose a 1-D CNN model. The latter largely improves the learning ability and reduces the number of

parameters to be trained by tailoring the receptive fields of filters. Experiments show 1-D CNN reduce the relative error up to 50% with less training and testing time compared to 2-D CNN.

Moreover, we propose four preprocessing techniques to further enhance the performance and generalization ability of the CNN model. *Adjacency Matrix Zero Padding* standardize the input adjacency matrices and make them validate for training. *Swapping Augmentation* and *Node Ordering* leverage on graph isomorphism to solve overfitting and underfitting problem. *Random Node Sketching* solves the memory restriction problem for large graphs with provable error bound. Extensive experiments on synthetic random graphs and real world graphs show the effectiveness of these preprocessing techniques.

Contributions

- **Novel Framework.** To the best of our knowledge, we are the first to formulate the Graphlet Count Learning problem (GCL), which bridges the gap between classic graphlet count problem and the learning paradigm. We are also the first to propose a deep learning based 2-D CNN model and its improved version, 1-D CNN model, to solve the GCL problem for all 3, 4, 5-node graphlets. Together with four novel preprocessing techniques, we show our framework’s potential to compete with the state-of-the-art methods, and is accurate, efficient and general given different graph topological structures for different k -node graphlets.
- **Extensive Experiments.** We conduct extensive experiments on synthetic random graphs, real world biochemistry networks, arXiv Collaboration networks and Facebook social networks. We also compare our framework with a series of state-of-the-art graphlet counting frameworks to validate the effectiveness of our framework.
- **Accuracy & Efficiency.** Our framework can offer up to two orders of magnitude speedup on synthetic graphs and achieves on par speed for real world graphs with highly competitive accuracy compared to state-of-the-art exact/sampling methods.
- **Generalizability.** Our framework gives accurate estimation given different graph structures or underlying distributions (e.g. sparse/dense graphs, small/large graphs, graphs from different graph models etc.). Moreover, our framework can be easily extended to estimate other k -node graphlet counts (e.g., 6-node graphlets) without changing the learning model.

The rest of the paper is organized as follows. The preliminary for graphlet counting, Convolutional Neural Networks and problem formulation are given in Section 2. The proposed CNN-based graphlet count learning framework is described in Section 3. The framework evaluation experiment results are presented in Section 4. Related works are discussed in Section 5. Finally, Section 6 concludes the paper and points out the future direction.

2 Preliminary

2.1 Notations and Definitions

Our input network can be modeled as graph $G = (V, E)$, where V is the set of nodes and E is the set of edges of this graph. In this work, we study undirected and unweighted graph that has no self-loop or multiple edges between any two nodes, but we do not require the graph to be connected.

- **Induced Subgraph.** A subgraph can be either an *induced subgraph* or a *non-induced subgraph*. A k -node induced subgraph $G_k = (V_k, E_k)$ has k nodes in V_k and contains all the edges in the original graph whose both endpoints are in V_k . Formally, an induced subgraph in a graph $G = (V, E)$ can be written as $G_k = (V_k, E_k)$, where $V_k \subset V$, $|V_k| = k$, and $E_k = \{(u, v) : \forall u, v \in V_k \wedge (u, v) \in E\}$. In contrast, a non-induced subgraph does not require that all edges connecting nodes in V_k are also present in E_k . For example, in Figure 1, when $V_4 = \{2, 3, 4, 5\}$, edge set $\{(2, 3), (2, 5), (4, 5)\}$ form a non-induced subgraph because it does not contain the edge $(2, 4)$, while the same vertex set V_4 with edge set $E_4 = \{(2, 3), (2, 4), (2, 5), (4, 5)\}$ represents an induced subgraph.
- **Isomorphism.** For two graphs, $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$, graph isomorphism is defined as a bijection mapping $\psi : V_1 \rightarrow V_2$ and $(u, v) \in E_1 \Leftrightarrow (\psi(u), \psi(v)) \in E_2, \forall u, v \in V_1$. If there exists such a mapping between G_1 and G_2 , these two graphs are isomorphic. For instance, the two graphs in Figure 2 are isomorphic to each other, where the mapping is $\psi(a) = d, \psi(b) = a, \psi(c) = c, \psi(d) = b$. The graph isomorphism (GI) problem is the computational problem of determining whether two finite graphs are isomorphic.
- **Graphlet and Graphlet Count.** A graphlet is a connected induced subgraph in a graph G . We denote a family of k -node graphlets as $\mathcal{G}^k = \{g_1^k, \dots, g_m^k\}, k = 3, 4, 5$. For example, there are two different 3-node graphlets

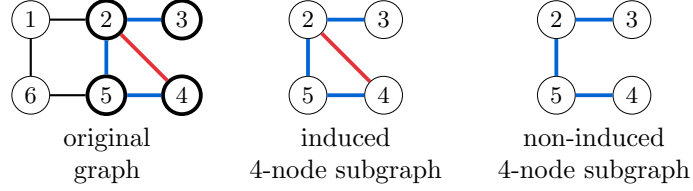


Figure 1: An Example of Induced/Non-induced Subgraph

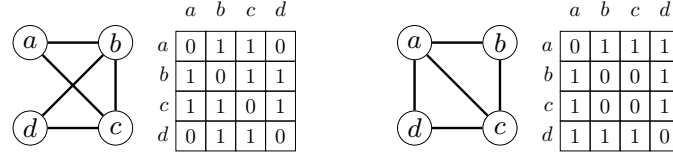


Figure 2: An Example of Isomorphic Graphs

and six different 4-node graphlets. Figure 3 depicts all the 3-node and 4-node graphlets, and Table 1 depicts all 5-node graphlets.

Graphlet count $c_i^k \in \mathbb{Z}_{\geq 0}$ is defined as the number of induced subgraphs that are isomorphic to graphlet g_i^k in a graph G . Figure 3 also gives instances of graphlet counts. In this example graph, $c_1^4 = 2$ because there are two induced subgraphs isomorphic to 4-path g_1^4 ($\bullet\text{---}\bullet\text{---}\bullet\text{---}\bullet$). One is induced by nodes $\{1, 2, 3, 5\}$ and the other is induced by nodes $\{1, 2, 4, 5\}$.

Graphlets	g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Shape								
Count	4	2	2	0	0	1	1	0
GFD	0.4	0.2	0.2	0.0	0.0	0.1	0.1	0.0

Figure 3: An Example of Graphlet Counts for All 3-node and 4-node Graphlets

- **Graphlet Frequency Distribution (GFD).** Given a set of graphlets $S = \{g_1, \dots, g_m\}$, graphlet frequency distribution (GFD) can be denoted as a vector $\mathbf{f} = [f_1, \dots, f_m]$, where f_i represents the frequency of the graphlet g_i , and $f_i = \frac{f_i}{\sum_{i=1}^m f_i}$. For the example in Figure 3, we calculate the frequency of each graphlet among all 3-node and 4-node graphlets.
- **Relationship between GFD and Graphlet Count.** If we have the graphlet counts $\{c_1, \dots, c_m\}$ for graphlets $\{g_1, \dots, g_m\}$, we can obtain graphlet frequency distribution (GFD) as $f_i = \frac{c_i}{\sum_{i=1}^m c_i}$. On the other hand, if we have GFD, we can also derive the graphlet counts with the help of node degree information d_v [9, 13], where degree d_v is the number of neighbor nodes of node v . For example, $\sum_{v \in V} \binom{d_v}{2} = c_1^3 + 3 \cdot c_2^3$, $\sum_{v \in V} \binom{d_v}{3} = c_2^4 + c_4^4 + 2 \cdot c_5^4 + 4 \cdot c_6^4$, $\sum_{v \in V} \binom{d_v}{3} = c_2^4 + c_4^4 + 2 \cdot c_5^4 + 4 \cdot c_6^4$ and $\sum_{v \in V} \binom{d_v}{4} = c_5^5 + c_6^5 + c_9^5 + c_{11}^5 + 2c_{14}^5 + c_{15}^5 + c_{16}^5 + c_{18}^5 + 2c_{19}^5 + 3c_{20}^5 + 5c_{21}^5$. For every family of k -node graphlets \mathcal{G}^k , we can find a linear relationship as such. As an example, in Figure 3, we have $\sum_{v \in V} \binom{d_v}{2} = 10$ and $f_{g_1^3} = 0.4$, $f_{g_2^3} = 0.2$. We can derive the open triangle ($\bullet\text{---}\bullet\text{---}\bullet$) count and the close triangle ($\bullet\text{---}\bullet\text{---}\bullet$) count by letting $c_1^3 = f_{g_1^3} \cdot x$ and $c_2^3 = f_{g_2^3} \cdot x$ and calculating x via the equation $\sum_{v \in V} \binom{d_v}{2} = f_{g_1^3} \cdot x + 3 \cdot f_{g_2^3} \cdot x$.
- **Computational Complexity.** Exact computation of graphlet counts is inherently expensive as the graphlet size and/or the graph size grow. First, the number of non-isomorphic k -node graphlets grows exponentially large as k increases. To illustrate, the number of different types of 4, 5-node graphlets are 6 and 21; for 6-node graphlets, it becomes $|\mathcal{G}^6| = 112$; $|\mathcal{G}^7| = 853$. Second, it is critical to decide an induced subgraph is isomorphic to which k -node graphlet, which can be reduced to graph isomorphism problem. Unfortunately, this problem is not known to have polynomial solutions [25]. These two facts make graphlet counting highly challenging when graphlet size k is large. However, even k is small (i.e., $k = 3, 4, 5$), a specific k -node graphlet count can

Graphlets	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
Shape											
Graphlets	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
Shape											

Table 1: All 5-node Graphlets

be up to $\binom{N}{k}$, which is k -th order polynomial or superlinear with respect to the graph size N . This means, counting graphlets using naive enumeration method is computationally intensive. For example, a moderate size arXiv collaboration network with 12K nodes can have more than 3.3 million triangles and 0.2 billion 4-paths. Straight-forward exhaustive enumeration approach has to count each graphlet at least once, which has prohibitively high cost. Experiments show that even the state-of-the-art exact counting method [8] can not finish within a week to calculate 4-node graphlet counts for a Twitter graph with 21.3M nodes [9].

2.2 Convolutional Neural Network (CNN)

Convolutional neural network (CNN) is a type of feedforward artificial neural networks, which were initially developed for modeling biological neural systems [26]. CNN has been widely used on learning regular grid data and has led to breakthroughs in image recognition [22], video recognition [23] and natural language processing [24] problems.

In a CNN architecture, multiple neural network layers are stacked in sequence. Typically, a CNN consists of two main types of layers: convolutional layer and fully connected layer [27]. In the context of image task, an image of height 32 pixels, width 32 pixels, and three color channels, R, B, G, will be presented as a 3-D tensor of shape $(32, 32, 3)$. Taking such a tensor as the input, a convolutional layer computes output values by performing dot products between its learnable weight filters and local regions, or receptive fields, of the input tensor. For example, from left to right, from top to down throughout every location in the input tensor, one weight filter of size $(5, 5, 3)$ can have at most $28 \times 28 \times 1$ different dot products by sliding the weight filter by $(1, 1)$ each time. Each dot product usually follows by a non-linear activation (e.g., $\text{ReLU}(x) = \max(0, x)$). Hence, in this case, if the first convolutional layer applies 64 weight filters, the output tensor would be of shape $(28, 28, 64)$. And the next neural network layer will take it as an input. In a fully connected layer, each output value is a weighted sum of all elements of the layer's input. Unlike convolutional layer which only focuses on one small local region at a time, fully connected layer connects to the whole input tensor. For instance, suppose the next layer is a fully connected layer with 16 output values. To produce one output value, fully connected layer takes the $28 \times 28 \times 64$ elements as the input, sum them up with $28 \times 28 \times 64$ learnable weights and a bias, and applies a non-linear activation on the sum. This activation layer produces a size $(1, 1, 16)$ output.

In general, a CNN architecture serves as a transformation function which maps input tensors to values or labels, and this function is parameterized by non-linear activations and learnable weights and biases. Parameters such as weights and biases will be trained in an iterative manner, i.e., via backpropagation refinement.

2.3 Problem Formulation

In this work, we are solving the following Graphlet Count Learning (GCL) problem. Suppose we are given a particular type of k -node graphlet g_i^k , $k = 3, 4, 5$, and a dataset \mathbb{G} consisting of m undirected unweighted graphs from the same underlying distribution \mathcal{D} , i.e., $\mathbb{G} = \{G_1, \dots, G_m\}$ and $\mathbb{G} \sim \mathcal{D}$. We assume the dataset can be divided into two parts $\mathbb{G}_1 = \{G_1, \dots, G_{m_1}\}$ and $\mathbb{G}_2 = \{G_{m_1+1}, \dots, G_m\}$, where graphlet counts c_i^k are known in \mathbb{G}_1 and unknown in \mathbb{G}_2 . Our goal is to develop a framework that can learn from \mathbb{G}_1 and estimate unknown graphlet count c_i^k for every graph in \mathbb{G}_2 . Mathematically, for any graphlet g_i^k , we want to learn a function $y_j = f_i^k(G_j|\mathcal{D})$ mapping from the graph structure G_j to the graphlet count y_j for $j \in \{m_1 + 1, \dots, m\}$, where the function value $f_i^k(G_j|\mathcal{D})$, $j = 1, \dots, m_1$ are given.

Note that we assume our data are coming from the same (or similar) distribution because the function we want to learn could be related to the underlying graph distributions and has totally different properties for different distributions. Conventional graphlet count frameworks can directly estimate or compute exact graphlet counts in \mathbb{G}_2 . But they fail to utilize the ground truth \mathbb{G}_1 , which will suffer loss in accuracy or efficiency or even both.

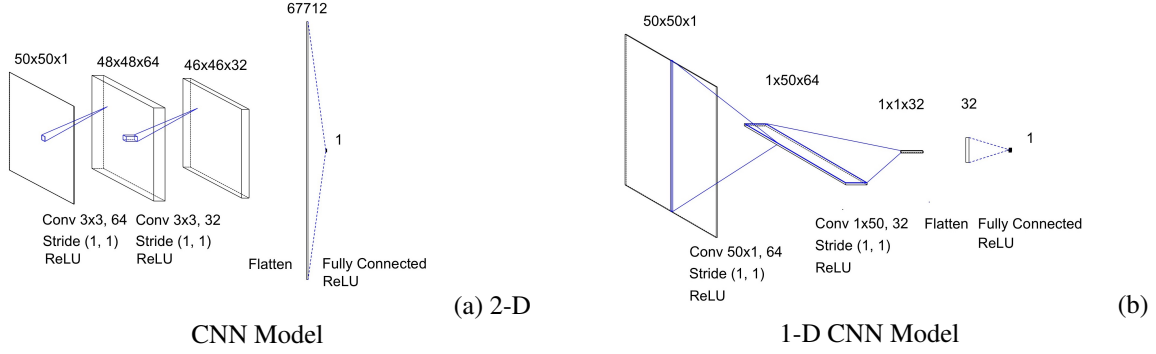


Figure 4: CNN Model Architectures

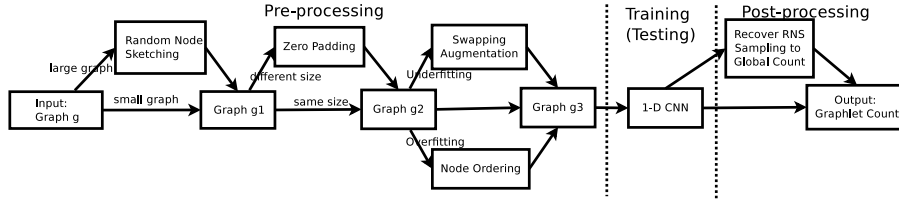


Figure 5: Framework Architecture Diagram

3 Methodology

In this section, we introduce the convolutional neural network (CNN) based framework to address the GCL problem. This framework consists of two parts: the core CNN architecture and the preprocessing techniques.

In Section 3.1, we propose the basic 2-D CNN model and its improved version, 1-D CNN model, and discuss the main idea and detailed architecture of both CNN models. In Section 3.2, four preprocessing techniques are designed to (a) standardize the dataset to create valid inputs for CNN models, (b) enhance the learning power and improve the model accuracy, (c) solve the memory restriction when dealing with large graphs. In Section 3.3, we assemble above models and techniques as a whole and show the pipeline of our framework.

3.1 Convolutional Neural Network Models

The basic idea of CNN models is to learn a transformation function which maps the input graph to its graphlet counts. This function is parameterized by learnable weights and biases which are associated with different layers in CNN, from shallow to deep. These parameters carry local structural information in the shallow layers, stack and refine themselves through deeper layers by some non-linear activation functions. In the final layer, the CNN model combines results from previous layers and produces an estimated graphlet count.

In our work, both 2-D CNN model and 1-D CNN model are developed for taking a graph adjacency matrix ($N \times N$ matrix, where N is the number of nodes) and/or its variants as input and giving a graphlet count (non-negative integer number) as output. And the five data preprocessing methods are also for handling graph adjacency matrices. Note that we choose the adjacency matrix as our input data structure because the CNN model heavily relies on matrix(tensor) operations, e.g., matrix(tensor) addition and multiplication and etc. Other graph representations, such as edge sequences and adjacency lists, are not meaningful or well defined under these operations. However, one major drawback of the adjacency matrix is that it requires $O(N^2)$ space for a graph with N nodes and exhausts memory resources when N is large. We will talk about how to solve this problem in Section 3.2.

3.1.1 2-D CNN Model

The main idea of 2-D CNN model is to think of the adjacency matrix as an "image" and apply vanilla CNN to it. Intuitively, each entry in the adjacency matrix forms a pixel-like base element in this "image". The graphlets, such as triangles, paths and etc., are analogous to higher order patterns, e.g. textures or pattern related objects in an image. Once

our first few shallow layers of 2-D CNN can recognize the graphlet patterns, the deep layers can sweep the previous output of shallow layers and combine them to estimate the graphlet count of the adjacency matrix.

Based on this idea, we design a 3-layer 2-D CNN model whose architecture is illustrated in Figure 4(a). For an input of shape $(N, N, 1)$, this 2-D CNN model deploys two convolution-activation layers, both with filters of shape $(3, 3)$, followed by a fully connected layer. In the first layer, 64 filters are employed in order to capture more variety of feature patterns. At the second layer, the number of filter decreases to 32. Shape $(3, 3)$ filter has been shown more effective in many deep learning researches, and same for our case. We choose the ReLU function (i.e., $\text{ReLU}(x) = \max(0, x)$) for all the activations in this model because the graphlet count estimation is a regression task for which we expect the output to be a real number in $[0, \infty)$.

Let us define some notations for our CNN models. Let $\mathbf{O}^{(l)} \in \mathbb{R}^{N^{(l)} \times N^{(l)} \times C^{(l)}}$ be the output tensor at layer l , where $l = 0, 1, 2, 3$, $N^{(l)}$ denotes the width (and height) along each channel and $C^{(l)}$ denotes the channel size. Let $\mathbf{O}_{i,j,t}^{(l)}$ be the $(i, j)^{\text{th}}$ element along the t^{th} channel. We assign $\mathbf{O}^{(0)}$ as the graph adjacency matrix. Mathematically, our 2-D CNN structure can be described as follows:

$$\mathbf{O}_{i,j,t}^{(l)} = \text{ReLU}(\mathbf{W}_t^{(l)} \cdot \mathbf{O}^{(l-1)}[i : i + H_1^{(l)} - 1, j : j + H_2^{(l)} - 1, :] + b_t^{(l)}), \quad (1)$$

$$\mathbf{O}^{(3)} = \text{ReLU}(\text{Flatten}(\mathbf{O}^{(2)})^T \cdot \mathbf{W}^{(3)} + b^{(3)}). \quad (2)$$

Equation (1) corresponds to two convolution layers, $l = 1, 2$. Each layer applies $C^{(l+1)}$ filters over the input feature map $\mathbf{O}^{(l-1)}$, and the t^{th} filter is parameterized by a trainable 3-D weight tensor $\mathbf{W}_t^{(l)} \in \mathbb{R}^{H_1^{(l)} \times H_2^{(l)} \times C^{(l)}}$, where $H_1^{(l)}$, $H_2^{(l)}$ denotes the width and height of the filter. $[a : b, c : d, :]$ is a slicing function which extracts subset of elements indexing from a to b (include both ends) in width, c to d (include both ends) in height and all in channel to form a new tensor. The operator \cdot is the sum of element wise product of two tensors. After adding bias term $b_t^{(l)}$, we apply ReLU ($\max(0, x)$) as the activation function to obtain the output feature map $\mathbf{O}^{(l)}$. Equation (2) is associated with the fully connected layer. It flattens the output $\mathbf{O}^{(2)}$ into a column vector, applies $\mathbf{W}^{(3)}$, $b^{(3)}$ and ReLU to obtain the estimated graphlet count. Our 2-D CNN is trained with back propagation and mean squared error as the loss function. Finally, we choose the shape of filters as $(H_1^{(1)}, H_2^{(1)}) = (3, 3)$ and $(H_1^{(2)}, H_2^{(2)}) = (3, 3)$. The channel size is set to be $C^{(0)} = 1, C^{(1)} = 64, C^{(2)} = 32, C^{(3)} = 1$.

3.1.2 1-D CNN Model

The straightforward 2-D CNN offers an efficient architecture to extract meaningful statistical patterns. In this section, we adapt the structure of the 2-D CNN and propose a 1-D CNN model, which is inspired by Johnson and Zhang's work [24] of applying the CNN for 1-D text data for region embedding. The 1-D CNN improves the efficiency by reducing the number of trainable parameters and enhance the learning performance by feeding the *complete* 1-hop neighborhood connection information into the first layer filters.

Concretely, we consider the input graph with N nodes and model its adjacency matrix as a tensor of shape $(N, N, 1)$. Our 1-D CNN model also deploys two convolution-activation layers and one fully connected layer, and adopts the ReLU activation function. The difference is the shapes of the filters. We apply 64 $(N, 1)$ -shaped filters instead of $(3, 3)$ -shaped filters at the first convolutional layer. Also, 32 $(1, N)$ -shaped filters is used at the second layer rather than another 32 $(3, 3)$ -shaped filters. The Figure 4(b) depicts this structure when $N = 50$.

Under such a paradigm, the first convolution layer cover all 1-hop neighbors for each node, and the function of this neural network layer is analogous to computing an embedding for each node. At the second layer, the $(1, N)$ -shaped filters learn the higher level latent features based on the first level embedding of the nodes. Mathematically, our 1-D CNN structure can be described using Equation (3), (4) and (5). It differs from the 2-D CNN model by setting the shape of filter $(H_1^{(1)}, H_2^{(1)}) = (N, 1)$ and $(H_1^{(2)}, H_2^{(2)}) = (1, N)$ in Equation (1) and (2).

$$\mathbf{O}_{i,j,t}^{(1)} = \text{ReLU}(\mathbf{W}_t^{(1)} \cdot \mathbf{O}^{(0)}[0 : N - 1, j : j, :] + b_t^{(1)}), \quad i = 0, j = 0, \dots, N - 1, t = 0, \dots, 63, \quad (3)$$

$$\mathbf{O}_{i,j,t}^{(2)} = \text{ReLU}(\mathbf{W}_t^{(2)} \cdot \mathbf{O}^{(1)}[0 : 0, 0 : N - 1, :] + b_t^{(2)}), \quad i = 0, j = 0, t = 0, \dots, 31 \quad (4)$$

$$\mathbf{O}^{(3)} = \text{ReLU}(\text{Flatten}(\mathbf{O}^{(2)})^T \cdot \mathbf{W}^{(3)} + b^{(3)}). \quad (5)$$

Compared with the 2-D CNN model which only applies small sized (i.e. (3, 3)-shaped) filters, the 1-D CNN allows the filters to take in more local connection information. Moreover, the total number of trainable parameters is only $O(N)$, whereas 2-D CNN has $O(N^2)$ in the last fully connected layer.

3.2 Preprocessing Techniques

Given the core CNN architecture, we further develop five preprocessing techniques to improve the performance and generalization ability of our framework. In Section 3.2.1, we propose Adjacency Matrix Zero Padding technique to standardize the size of input adjacency matrices and make them valid for training. In Section 3.2.2 and Section 3.2.3, we utilize graph isomorphism property to enhance the learning ability of CNN model by Node Ordering and Swapping Augmentation. In Section 3.2.4, we develop the Random Node Sketching to solve the memory restriction when dealing with large graphs and provide a theoretical guarantee for the accuracy of this approach.

3.2.1 Adjacency Matrix Zero Padding

Due to the constraint of the fixed fully connected units in CNN model, the size of input data has to be the same. However, in many real world graph datasets, graph samples not necessarily have the same size (same number of nodes). In the context of learning images, common solutions are to crop images or pad the images with a solid color so that they will have the same sizes. In our situation, cropping is not appropriate as it causes the loss of edge connectivity information that is essential for graphlet counting.

To preserve edge connectivity information of all training graphs, we consider the largest graph in the training set, and use its dimension (say N) as the dimension of the input adjacency matrix ($N \times N$). For other graphs in the training set, we take each adjacency matrix and pad it with zero till we have an input matrix of dimension $N \times N$. This solves the varying input size problem. Moreover, the physical meaning of padding adjacency matrix with zero is analogous to adding isolated nodes to a graph, which do not form any extra graphlet.

3.2.2 Node Ordering

In many cases, the nodes in a graph are not precisely labeled and strictly ordered; hence, an adjacency matrix is actually not a unique representation of a given graph. The very same graph can be represented by adjacency matrices that, on the face of it, look totally different. This is unfavorable for learning as it introduces more uncertainty/variance to the distribution of the dataset. If a CNN model is not powerful enough, this property makes it harder for neural network to learn the graphlet count estimation.

The node ordering technique works as follows. First, we label each node $v \in V$ in the input graph $G(V, E)$ by a specific labeling strategy and denoted the label as $l(v)$. In this work, we propose four different labeling strategies: (1) degree centrality: $l(v) = \deg(v)$, where $\deg(v)$ is the degree of v (2) closeness centrality: $l(v) = 1 / \sum_{u \in V} d(u, v)$, where $d(u, v)$ is the distance between node u and v . (3) betweenness centrality: $l(v) = \sum_{s \neq v \neq t \in V} \sigma_{st}(v) / \sigma_{st}$, where σ_{st} is total number of shortest paths from node s to t and $\sigma_{st}(v)$ is the number of those paths that pass through v . (4) Weisfeiler-Lehman algorithm: $l(v) = \text{WL}(v)$, where $\text{WL}(v)$ is the label given by Weisfeiler-Lehman algorithm [28]. Then we reorder the nodes according to the label value $l(v)$ in decreasing (increasing order) and settle down its adjacency matrix representation accordingly. This transformation are applied to all input adjacency matrices in the dataset. By node ordering, we obtain more consistent representations (canonical representations) of graphs to reduce the variance of the dataset.

3.2.3 Swapping Augmentation

As we mentioned in previous subsection, an unattributed graph which has no node labels can have many different representatives in adjacency matrices (as isomorphic graphs). This may be harmful if a CNN model is not powerful enough to fit the data. However, when the dataset contains too few samples or the model tends to overfit the data, we can utilize this isomorphism property to acquire sufficient data for training.

The approach is straightforward: we randomly pick two indices i and j , and then swap the i^{th} row with j^{th} row j and swap the i^{th} column with j^{th} column. We can repeat the swapping operation for each graph m times to create m more training data. Analogous to flipping or rotation of images, via augmenting the dataset, we can improve the generalization ability of CNN models and, thereby, improve the accuracy.

3.2.4 Random Node Sketching

Real world graphs can be extremely large. For instance, social networks like Sina Weibo can produce up to 21M-node, with 261M edges dataset [29]. In practice, our GPU and CPU resources are limited. Faced with the memory constraint, it is inefficient and impractical to build a CNN model which takes in the whole adjacency matrix as the input. Consequently, we save all our adjacency matrices using the sparse matrix format (i.e. COOrdinate format ¹), and unpack them when necessary. Moreover, we propose the *Random Node Sketching* (RNS) approach to further improve the our scalability. Intuitively, we are training our model on randomly sampled subgraphs (i.e., sub-matrices) instead of the whole (sparse) adjacency matrix. Note that there are other sketching/sampling methods, such as *Random Edge Sketching* (RES) which samples a subgraph by preserving each edge uniformly at random. Though RES may give more accurate estimation (especially for some of sparse 3-node and 4-node graphlets, e.g., g_1^3, g_1^4, g_2^4), we choose RNS over RES for the following reasons. First, RNS is more space efficient given the same node/edge sampling probability p , i.e., the expected edges preserved is p^2 for RNS, and p for RES. Second, RNS requires less trainable parameters, because our 1-D CNN model has $O(N)$ trainable parameters for a network of size N , and RNS reduces the number of nodes to pN while RES preserves almost all the nodes. Take the social network socfb-FSU53 [30] as an example, when RES samples 2% edges, 75% nodes are preserved, and RNS only preserve 2% nodes. Third, as to be explained below, for any graphlet, the count in the sampled RNS graphs has one-to-one correspondence from graphlet counts in RNS samples to the original large graph's. However, RES needs all other graphlet counts (e.g. graphlet counts for all 6 types of 4-node graphlets) in RES samples to infer the graphlet counts for the original graph [31], so the accuracy may decrease due to the inaccurate estimation of other graphlet counts.

The basic idea of RNS is to randomly generate a set of small subgraphs which can easily satisfy the memory constraint for each large graph. After that, we apply our CNN model to these subgraphs, combine the results together and estimate the graphlet count of the original large graph. We also provide a theoretical guarantee that this process will not amplify the relative error.

To be concrete, each RNS graph is generated as follows. For each large graph $G(V, E)$ in our dataset, we choose a node preserving probability p and determine m as the total number of RNS graphs to be generated for G . Then, we select a set of $p \cdot |V|$ nodes from G uniformly at random, and denote them as V' . One RNS graph $G'(V', E')$ is exactly the subgraph induced from G using this node set V' . Finally, We repeat this process m times to generate m RNS graphs $\mathbb{G}' = \{G'_1, G'_2, \dots, G'_m\}$.

After generating the RNS set \mathbb{G}' , we establish an equation between the expected graphlet count on G' and the truth count on G , with one-to-one correspondence. Let c and c' denote the graphlet count of a specific k -node graphlet g for the large graph G and any RNS graph G' respectively. Let the function $\mathbb{I}(x)$ be the indicator function and $g(j)$ be the j^{th} k -node graphlet instance of large graph G , where $j \in \{1, \dots, c\}$. We have:

$$\begin{aligned}
 E[c'] &= E\left[\sum_{j=1}^c \mathbb{I}(\text{Every node of } g(j) \text{ is preserved})\right] \\
 &= \sum_{j=1}^c E[\mathbb{I}(\text{Every node of } g(j) \text{ is preserved})] \\
 &= \sum_{j=1}^c Pr(\text{Every node of } g(j) \text{ is preserved}) \\
 &= p^k c
 \end{aligned} \tag{6}$$

We introduce a RNS based estimator of c :

$$\hat{c} = p^{-k} \sum_{i=1}^m \tilde{c}(i) / m \tag{7}$$

, where $\tilde{c}(i)$ is the estimated graphlet count by applying the CNN model on the i^{th} RNS graph G'_i . We define the relative error for our CNN model on each RNS graph G'_i to be $e(i) = |\tilde{c}(i) - c(i)|/c(i)$ and the relative error for the RNS estimator of large graph G as $e = |\hat{c} - c|/c$. Note that the RNS based estimator \hat{c} is not unbiased because we can not guarantee our CNN model is unbiased. We have the following theorem to bound the relative error with $(1 - \delta)$ probability:

¹<https://docs.scipy.org/doc/scipy/reference/sparse.html>

Theorem 3.1. *Given any graph G , its any k -node graphlet count c , its RNS set $\mathbb{G}' = \{G'_1, \dots, G'_m\}$ with node preserving probability p and any parameters $\alpha, \delta > 0$, if the relative error of CNN is at most ϵ , i.e., $e(i) \leq \epsilon$ for each $i = \{1, \dots, m\}$, then the relative error on G is at most $\epsilon + \alpha(1 + \epsilon)$ with $(1 - \delta)$ probability, i.e. $e \leq \epsilon$, when $m \geq \frac{1}{2}\alpha^{-2}p^{-2k} \log \frac{2}{\delta}$.*

Proof of Theorem 3.1. We first use the Hoeffding's inequality [32] to give a high probability bound of the relative error $|c(i) - \mathbb{E}[c(i)]|/\mathbb{E}[c(i)]$, given $m \geq \frac{1}{2}\alpha^{-2}p^{-2k} \log \frac{2}{\delta}$. Because $c(i)$ are independent RNS samples and bounded by $[0, c]$, we define $\bar{c} = \sum_{i=1}^m c(i)/m$ and use the Hoeffding's inequality,

$$\begin{aligned} Pr(|\bar{c} - \mathbb{E}[\bar{c}]| \leq \alpha \mathbb{E}[\bar{c}]) &\geq 1 - 2 \exp\left(-\frac{2m^2\alpha^2\mathbb{E}[\bar{c}]^2}{mc^2}\right) \\ &= 1 - 2 \exp\left(-\frac{2m^2\alpha^2p^{2k}c^2}{mc^2}\right) \\ &\geq 1 - \delta \end{aligned} \tag{8}$$

The relative error on the large graph G is then bounded with probability at least $(1 - \delta)$:

$$\begin{aligned} e &= |\hat{c} - c|/c \\ &= |p^{-k} \sum_{i=1}^m \tilde{c}(i)/m - c|/c \\ &\leq |p^{-k} \sum_{i=1}^m (1 \pm \epsilon)c(i)/m - c|/c \\ &\leq |p^{-k} p^k c(1 \pm \epsilon)(1 \pm \alpha) - c|/c \\ &\leq \epsilon + (1 + \epsilon)\alpha \end{aligned} \tag{9}$$

The first and second equality holds by definition. The first inequality holds because the relative error of CNN for each instance i is at most ϵ . The second inequality holds because of the inequality (8) and the fact $\mathbb{E}[\bar{c}] = p^k c$. \square

Theorem 3.1 can be viewed as the sensitivity analysis of RNS method. In the high probability error bound (Inequality (9)), the first part of the error term ϵ is brought by the inaccurate estimation of the CNN model and the second part $(1 + \epsilon)\alpha$ is brought by inadequate RNS samples. When we have sufficient number of RNS graphs, i.e. $m \rightarrow \infty$, the relative error $e = \epsilon$, which is no larger than the relative error of CNN. This shows that Random Node Sketching method is effective in leveraging CNN to estimate graphlet counts on large graphs without amplifying the relative error. Note that this theoretical high probability bound is conservative whereas fewer RNS samples (or smaller preserving probability p) are needed to achieve accurate estimation in practice.

3.3 Overall CNN based Framework

Figure 5 presents the overall architecture of the proposed framework. Given a set of graphs with known graphlet count, some preprocessing techniques are applied before inputting these graphs into CNN for training. First, we determine if the size of graphs can fit into the memory, if not, we use *random node sketching* to obtain some sample graphs. Second, we apply *zero padding* if the sizes of graphs in the dataset are not consistent. Third, we use *swapping augmentation* to enlarge the training dataset if the graph samples are insufficient. On the contrary, if the data is adequate, we use *node ordering* to enhance the accuracy of count estimation. Finally, we use the preprocessed data to train a CNN model which learns a mapping from graph structures to graphlet counts. For large graphs, we recover the global graphlet count from the graphlet count estimations for random node sketching samples in the end.

4 Experiments

In this section, we present extensive experiments to evaluate the performance of our framework. We aim to answer the following questions:

- **Q1 Accuracy:** Is the graphlet count estimation of our framework accurate? Is 1-D CNN always more accurate than 2-D CNN?

- **Q2 Efficiency:** Is our framework faster than other existing methods?
- **Q3 Generality:** Does our framework give accurate estimation on any graph structures (e.g. sparse v.s. dense graphs, graphs of different random graph models)?
- **Q4 Scalability:** Can our framework handle large graph?
- **Q5 Practicality:** Is our framework effective on real world networks?

4.1 Experimental Setup

4.1.1 Performance Evaluation Metrics

We consider the following statistics and metrics, which provide a comprehensive picture of the graph dataset property, graphlet count estimation errors, and running speed, for evaluating the performance of the proposed framework. All results are averaged over multiple runs (i.e. 10 runs) to reduce the random effects in the experiments.

- **Mean and Standard Deviation (STD) of the Ground Truth.** Given a graph dataset with S graph samples. Let c_i be the ground truth graphlet count of sample graph i , $i = 1, \dots, S$. We calculate the mean of ground truth counts as $\mu = \sum_{i=1}^S c_i / S$ and the standard deviation (STD) of ground truth counts as $\delta = \sqrt{\sum_{i=1}^S (c_i - \mu)^2 / S}$. These two statistics provide information about the underlying distribution of the dataset. The former shows the average count of each graphlet. The later measures the heterogeneity/diversity of a dataset by showing the derivation from the average graphlet count.
- **Mean of Absolute Error (MAE).** Given a graph dataset with S graph samples. Let c_i be the ground truth graphlet count and c'_i be the estimated graphlet count of sample graph i , $i = 1, \dots, S$. We compute the mean of absolute error as $MAE = \sum_{i=1}^S |c'_i - c_i| / S$. This metric measures how accurate our estimation is.
- **Relative Error (RE).** Although mean of absolute error (MAE) measures the accuracy of estimation, it does not show us if such error is acceptable for the given dataset. For example, $MAE = 10$ may be a good result if the mean of ground truth over the dataset is 1000; on the contrary, if the mean of ground truth is 20, this error is less acceptable. In order to measure the effectiveness of the estimation, we compute the ratio between mean of absolute error (MAE) and mean of ground truth. We take relative error as $e = MAE / \mu$.
- **Running Time.** In order to evaluate the efficiency of our methods, we choose the running time as the evaluation metric. For our CNN framework, we provide both training time and testing time. The training time records the total running time to load the dataset into the memory, preprocess the dataset (e.g. node ordering, RNS sampling, etc.) and train the model parameters. The testing time records the total running time to predict the testing graphs in the dataset. For baseline methods, we provide the total running time to count/estimate testing graphs after they have been preprocessed and loaded into the memory. While training time is an important criteria to evaluate the efficiency, we care more about the testing time and use testing time to compete with other methods. The reason is that our CNN method has the “once for all” property. Concretely, once we have trained our model, we do not need to retrain it for all the new coming data from the same distribution \mathcal{D} . Since the training time can be amortized over each testing graph, it can be neglected when we have considerable number of testing graphs.

Moreover, we separately report the testing/testing time for each graphlet, because the user may want to separately estimate the count for a specific graphlet. Notice that our CNN model are fully decoupled for different graphlets and can be highly paralleled. Therefore, the time for each graphlet can also represent the total time for counting all graphlets. As for competing methods, they depend on counts of other graphlets to count a specific graphlet, and we use total time for each graphlet unless they provide separate modules to count a specific graphlet.

4.1.2 Dataset

We synthesize datasets with three random graph models: Barabasi-Albert (BA) graph, Erdos-Renyi (ER) graph, and random geometric graph (RGG). The size of synthetic graphs in experiments, if not specified otherwise, is 50-node. In each dataset, we have 3000 training graphs, 300 validation graphs, and 300 testing graphs. The ground truth of graphlet counts is calculated by PGD library [33].

- **Barabasi-Albert (BA) Graph.** Barabasi-Albert model generates random graphs with the preferential attachment mechanism. Specifically, starting from the base connected graph, every newly added node connects to m existing nodes, and the probability of connecting to a particular existing node is proportional to its degree. Hence, the degrees in a Barabasi-Albert random graph is exponentially distributed.

Table 2: Biochemistry Dataset Details

Dataset	MUTAG	NCI1	NCI109
Maximum # of Nodes	28	111	111
Average # of Nodes	17.93	29.87	29.86
Number of Graphs	188	4110	4127
# of Training Set Graphs	160	3800	3800
# of Validation Set Graphs	8	100	100
# of Testing Set Graphs	20	210	227

Table 3: Collaboration & Social Network Dataset Details

Graph	ca-hep-ph	ca-hep-th	socfb-MSU24	socfb-FSU53
# of Node	28k	23k	32k	28k
# of Edge	3.1M	2.4M	1.1M	1.0M
Triangle Count	195.7M	191.3M	6.5M	7.9M

- **Erdos-Renyi (ER) Graph.** In Erdos-Renyi random graph, the edge between every two nodes exists with probability q . The degree distribution in Erdos-Renyi random graph is mostly evenly distributed.
- **Random Geometric Graph (RGG).** A random geometric graph is constructed by placing nodes uniformly at random in a unit cube and connecting two nodes by an edge if and only if their distance is within a given radius r . Different from Barabasi-Albert graph or Erdos-Renyi graph, random geometric graph has more clear community properties in its structure.

We also evaluate our proposed framework on three categories of real world networks: biochemistry graph datasets, collaboration networks and social networks.

- **Biochemistry Datasets.** We conduct experiments on three biochemistry datasets: the MUTAG dataset [34] which contains the graphs of mutagenic aromatic and heteroaromatic compounds, and the NCI1 and NCI109 dataset [35] which contains the graphs of chemical compounds tested on lung cancer cells and ovarian cancer cells respectively. We divide each dataset into three parts, one for training, one for validation, and one for testing. Details are listed in Table 2.
- **Collaboration Networks.** We evaluate our framework with two collaboration graphs: ca-hep-ph and ca-hep-th [36]. These two graphs are the author networks of arXiv where a node denotes an author and an edge denotes a common publication. ca-hep-ph and ca-hep-th are the networks of the High Energy Physics Phenomenology section and the High Energy Physics Theory section respectively. Details of these two graphs are listed in Table 3.
- **Social Networks.** We also evaluate our framework with two sparse social networks: socfb-MSU24 and socfb-FSU53 [30]. These two graphs are social friendship networks extracted from Facebook consisting of people as nodes and friendship ties as edges. The detailed information can be found in Table 3

4.1.3 CNN Training Setup

All adjacency matrices are saved using the sparse matrix format to save space, and unpacked when necessary. When space are limited, we use sparse matrix multiplication to train (or test) our model, instead of unpacking them during training (or testing). We train our CNN models on training datasets using the mean squared error loss function and the Adam optimizer [37]. Hyperparameters are tuned with validation datasets: the batch size was selected from $\{16, 32, 64\}$ examples and the learning rate is selected from $\{0.0001, 0.0005, 0.001, 0.005, 0.01\}$. Our models are trained for at most 600 epochs (with early stopping), which usually takes up to 20 minutes (depends on the total number and the size of training graphs) on one 12GB memory GeForce GTX TITAN X GPU. Finally, we test our model on testing datasets using best hyperparameters. Moreover, in all our experiments, we find different hyperparameters does not affect our results very much and it is always safe to choose the learning rate as 0.001 and batch size as 32. For some rare cases, the 0.001 learning rate is too large and we then use 0.0001 instead to get stable results. For each set of the experiment, we repeat 10 times, present the average of relative errors e_1, \dots, e_{10} as final result, and use the standard deviation of these ten relative errors as the error bar in our figures.

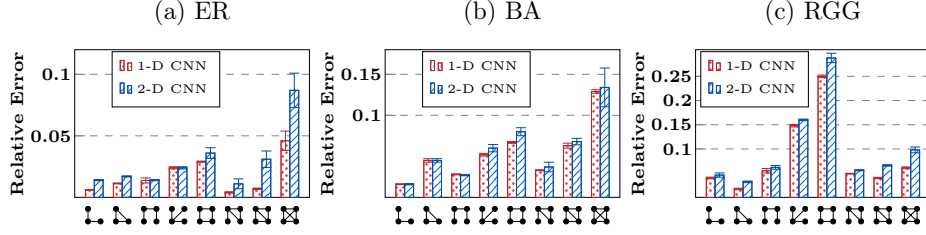


Figure 6: 1-D CNN v.s. 2-D CNN Performance

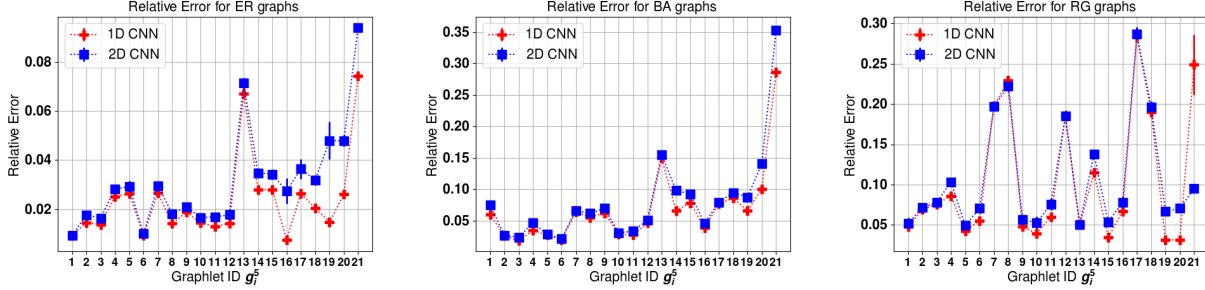


Figure 7: 1-D CNN v.s. 2-D CNN for 5-node graphlets

4.2 1-D CNN v.s. 2-D CNN

In Section 3, we propose two CNN models and discuss the physical meaning of 1-D and 2-D convolutions. In this subsection, we validate our hypothesis by comparing the accuracies of 1-D CNN and 2-D CNN on estimating 3,4-node graphlets in three kinds of random graphs. In the experiments presented in Figure 6, 1-D CNN is always more accurate than 2-D CNN. Although the accuracy improvements on estimating tree-like graphlets, such as 4-path (g_1^4) and 3-star (g_2^4), are slight, the accuracy improvements on approximating dense graphlets, such as 4-chordal cycle (g_5^4) and 4-clique (g_6^4), go up to 50%. For 5-node graphlets, Figure 7 gives the similar result, which shows that 1-D CNN outperforms 2-D CNN for all random graphs. Another advantage of 1-D CNN is that it reduces the number of weights to be trained compared with the 2-D CNN as we mentioned in Section 3.1.

4.3 Comparison with Competing Methods

We empirically compare the speed of our framework with that of state-of-the-art graphlet count or graphlet frequency distribution (GFD) estimation algorithms. To compare with GFD estimation works, we derive graphlet counts from the estimated GFD using the relationship as introduced in Section 2. In this experiment, we evaluate our model and previous works by estimating all eight types of 3, 4, 5-node graphlets on three synthetic datasets: 50-node Erdos-Renyi random graphs (ER graphs) with edge existing probability $q = 0.5$, 50-node Barabasi-Albert graphs (BA graphs) with $m = 6$ and 50-node Random Geometric graphs (RGG) with $r = 0.45$. Each dataset contains 3000 training graphs, 300 validation graphs and 300 testing graphs. We benchmark with six sampling frameworks: GRAFT [11], CC2 [15], GUISE [38], LGE [29], ApproxG [39] and MOSS-5 [40]. We also provide the running time for exact counting method PGD EXACT [8] and ESCAPE [41] for reference. Following is a brief synopsis of these works:

- **GRAFT** [11]. GRAFT first samples a set of edges uniformly randomly. Then for each sampled edge, it counts the number of occurrence of each graphlet that uses this edge. In the end, the obtained sample graphlet counts are scaled up to approximate the global 3,4,5-node graphlet counts. Note that GRAFT are hard wired to compute statistics for all 3,4,5-node graphlets, and do not provide separate modules to compute 3,4-node graphlets.²
- **CC2** [15]. Color coding is a general approach that prunes the enumeration search tree by sampling. Yet, general color coding was proposed only for estimating treelet counts. [42] This work adopts the general color

²We use the implementation of GRAFT from <http://dmgroup.cs.iupui.edu/mmrahman-GRAFT.php>.

Table 4: Running Time of 1-D CNN and Competing Methods for 3,4-node graphlets

		(a) ER Graph							
		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Training (in Minutes)	CNN (Ours)	1.55	2.73	3.92	3.22	3.47	3.36	4.05	2.96
Testing (in Seconds)	CNN (Ours)	0.012	0.02	0.01	0.015	0.019	0.02	0.025	0.02
	GRAFT	3186	3036	2818	2471	1763	2632	1751	2106
	CC2	-	-	-	-	-	-	-	-
	GUISE	2978	2978	3349	2709	3349	2632	2709	4579
	LGE	8.58	8.58	8.58	8.58	8.58	8.58	8.58	8.58
	ApproxG	8.30	8.30	9.94	9.81	8.99	8.97	9.94	9.02
	MOSS-5	4.98	5.02	5.66	6.30	5.65	5.34	3.25	5.59
	PGD EXACT	10.20	10.20	10.20	10.20	10.20	10.20	10.20	10.20
	ESCAPE	2.82	2.82	2.82	2.82	2.82	2.82	2.82	2.82
		(b) BA Graph							
		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Time (in Minutes)	CNN (Ours)	0.49	0.36	0.60	1.82	0.35	1.40	0.99	0.38
Testing (in Seconds)	CNN (Ours)	0.012	0.02	0.01	0.015	0.019	0.02	0.025	0.02
	GRAFT	471	413	513	291	196	401	196	179
	CC2	-	-	-	-	-	-	-	-
	GUISE	1620	1620	2077	58	854	112	221	1620
	LGE	22.65	24.23	19.76	16.51	26.89	21.46	19.76	22.65
	ApproxG	9.10	9.09	9.20	9.18	9.18	9.20	8.99	8.99
	MOSS-5	6.19	6.19	6.39	6.39	6.21	6.27	6.27	6.39
	PGD EXACT	8.62	8.62	8.62	8.62	8.62	8.62	8.62	8.62
	ESCAPE	2.61	2.61	2.61	2.61	2.61	2.61	2.61	2.61
		(c) RGG Graph							
		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Time (in Minutes)	CNN (Ours)	1.75	2.51	2.56	1.00	1.72	2.56	2.53	2.51
Testing (in Seconds)	CNN (Ours)	0.02	0.015	0.015	0.03	0.015	0.02	0.02	0.015
	GRAFT	605	1052	605	202	138	605	605	400
	CC2	-	-	-	-	-	-	-	-
	GUISE	1037	1214	896	339	691	691	529	206
	LGE	7.05	7.01	6.91	22.85	7.46	22.85	12.98	3.22
	ApproxG	22.56	22.56	23.50	22.51	22.56	22.56	23.15	22.75
	MOSS-5	17.79	17.79	19.59	17.44	17.41	19.98	17.22	19.97
	PGD EXACT	9.28	9.28	9.28	9.28	9.28	9.28	9.28	9.28
	ESCAPE	2.91	2.91	2.91	2.91	2.91	2.91	2.91	2.91

Note: a hyphen (-) indicates that the method did not terminate within 6 hours. Bold faces are the best results.

coding technique and further extends it to estimate the frequency distribution of all kinds of graphlets up to 7-node.³

- **GUISE** [38]. GUISE employs a Markov Chain Monte Carlo method to sample 3,4,5-node graphlets uniformly and construct the approximate graphlet frequency distribution of a graph. Similar to GRAFT, its implementation is hard wired to compute statistics for all 3,4,5-node graphlets.⁴
- **LGE** [29]. The LGE framework derives 3,4-node global graphlet count estimation by computing graphlets in a set of sampled localized neighborhoods (egonet). Here, we stick to the settings in [29], sampling edge neighborhoods and applying the PGD [8] framework for the edge-centric graphlet counting.⁵
- **ApproxG** [39]. ApproxG applies edge sampling and parallel computing to reduce the 3,4-node graphlet counting runtime.⁶
- **MOSS-5** [40]. MOSS-5 is a 5-node graphlets counting method based on subgraph sampling. In their implementation, they separately provide MOSS-4 to count 4-node graphlets, which are used for estimating 4-node graphlet counts in our experiments.⁷

³We use the implementation of CC2 from <https://github.com/Steven-/graphlets>.

⁴We use the implementation of GUISE from <http://dmgroup.cs.iupui.edu/mmrahman-GUISE.php>.

⁵We use the implementation of PGD from <https://github.com/nkahmed/PGD>.

⁶We use the implementation of ApproxG from <https://bitbucket.org/approxg/approxg/src/master/>.

⁷We use the implementation of MOSS-5 from <http://nskeylab.xjtu.edu.cn/dataset/phwang/code/mosscode.zip>.

Table 5: Running Time of 1-D CNN and Competing Methods for 5-node graphlets

(a) ER Graph										
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5
CNN Train (Ours)	4.582	5.725	5.722	5.687	5.716	5.699	5.774	5.546	5.744	5.752
CNN Test (Ours)	0.022	0.027	0.024	0.024	0.024	0.024	0.024	0.025	0.026	0.026
MOSS-5	965.4	414.2	414.2	216.5	73.7	414.2	73.7	414.2	216.5	73.7
GRAFT	888.5	2668.4	2668.4	2074.6	1780.7	2819.2	2074.6	2668.4	2369.7	2369.7
ESCAPE	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34
	g_{11}^5	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5
CNN Train (Ours)	5.749	5.609	4.987	5.601	5.675	5.680	5.374	5.692	5.600	4.504
CNN Test (Ours)	0.027	0.026	0.025	0.023	0.025	0.023	0.026	0.022	0.026	0.024
MOSS-5	73.7	16.9	73.7	73.7	73.7	16.9	73.7	73.7	16.9	16.9
GRAFT	2369.7	888.5	1780.7	1780.7	2819.2	1780.7	2369.7	2668.4	2074.6	596.6
ESCAPE	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34	89.34
(b) BA Graph										
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5
CNN Train (Ours)	4.876	6.219	6.253	5.892	3.248	5.444	5.613	3.907	3.607	4.550
CNN Test (Ours)	0.023	0.024	0.022	0.024	0.024	0.022	0.023	0.026	0.025	0.025
MOSS-5	14.1	25.5	57.6	25.5	25.5	25.5	0.6	25.5	14.1	14.1
GRAFT	207.1	308.5	308.5	274.1	274.1	308.5	142.6	175.5	207.1	274.1
ESCAPE	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3
	g_{11}^5	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5
CNN Train (Ours)	4.516	3.525	4.004	3.814	3.986	3.402	3.153	3.937	3.349	2.874
CNN Test (Ours)	0.026	0.025	0.024	0.024	0.024	0.023	0.023	0.024	0.023	0.024
MOSS-5	5.6	3.0	5.6	3.0	5.6	5.6	5.6	0.6	3.0	5.6
GRAFT	207.1	75.1	175.5	207.1	274.1	175.5	175.5	240.1	240.1	75.1
ESCAPE	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3	136.3
(c) RGG Graph										
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5
CNN Train (Ours)	3.127	6.049	6.238	4.665	6.000	6.255	4.034	3.381	6.121	6.211
CNN Test (Ours)	0.026	0.027	0.024	0.026	0.026	0.025	0.028	0.025	0.025	0.025
MOSS-5	4858.6	336.4	226.8	336.4	15.1	15.1	336.4	226.8	15.1	7.1
GRAFT	1485.7	993.6	500.4	500.4	833.7	668.3	500.4	169.1	833.7	833.7
ESCAPE	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31
	g_{11}^5	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5
CNN Train (Ours)	3.997	1.607	4.267	6.082	6.125	3.160	3.557	6.115	5.742	4.364
CNN Test (Ours)	0.026	0.028	0.025	0.026	0.026	0.026	0.027	0.024	0.026	0.030
MOSS-5	7.1	4858.6	15.1	7.1	7.1	336.4	15.1	7.1	7.1	7.1
GRAFT	169.1	169.1	336.4	993.6	500.4	500.4	336.4	993.6	1321.6	169.1
ESCAPE	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31	78.31

Note: Bold faces are the best results. Time for training are in minutes, others are in seconds.

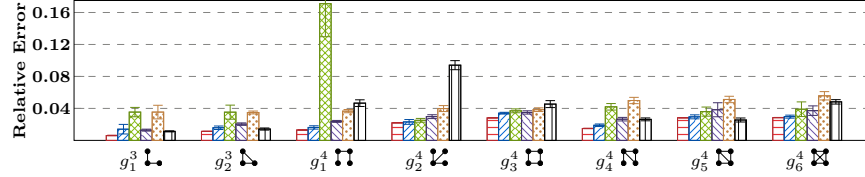
- **PGD EXACT** [33]. PGD EXACT represents the Parallel Parameterized Graphlet Decomposition (PGD) framework which counts 3,4-node graphlets by leveraging on the combinatorial relationships among these graphlets.
- **ESCAPE**. [41]. ESCAPE is the state-of-the-art algorithmic framework that can be used to exactly count any 5-node graphlets. Similar to MOSS-5, ESCAPE provides a separate module for 4-node graphlets, which are used when counting the 4-node graphlets.⁸

Our comparison results are summarized in Table 4 and Figure 8 for 3,4-node graphlets, and Table 5 and Figure 9 for 5-node graphlets. For a fair comparison, we **tune the number of iterations** for all benchmarking sampling methods, so

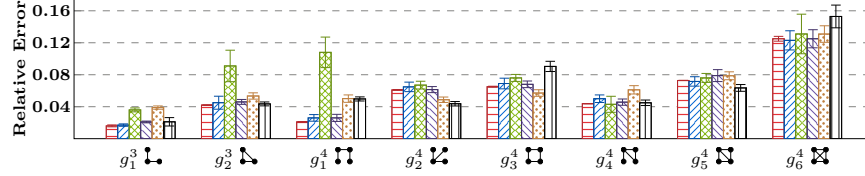
⁸Note that PGD EXACT and ESCAPE gives the exact graphlet count, and thus their relative errors are always 0.

1-D CNN GRAFT GUISE LGE APPROXG MOSS-5

(a) ER Graph



(b) BA Graph



(c) RGG Graph

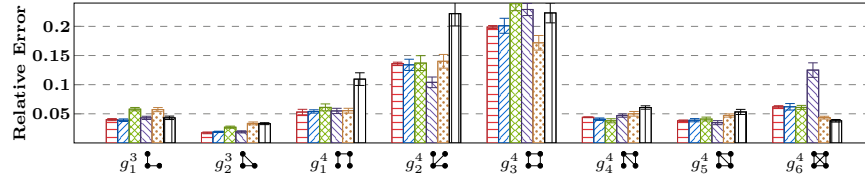


Figure 8: Relative Errors of 1-D CNN and Competing Methods for 4-node graphlets

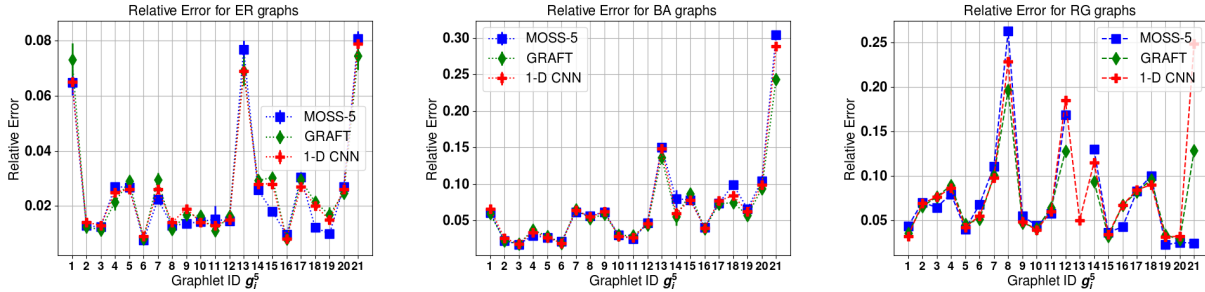


Figure 9: Relative Errors of 1-D CNN and Competing Methods for 5-node graphlets

that they obtain as close relative errors to that of 1-D CNN as possible. Take counting 4-clique (g_6^4) on the ER graphs as an example. Our 1-D CNN makes 2.8% relative error. In order to obtain similar accuracy, GRAFT samples 60% of edges, GUISE runs random walk for 400K iterations, and LGE samples 13% of edges (and their neighborhoods) for 450 runs. For MOSS-5 and ApproxG, the accuracy is controllable and we can directly set the error bound. The relative errors are shown in the Figure 8, where 1-D CNN and other competing methods achieve similar relative errors. For 3-4 node graphlets, Table 4 (a) shows that the running time (for total 300 testing graphs) used by estimation methods, i.e., GRAFT, CC2, GUISE, LGE, MOSS-5 and APPROXG, on random graphs are significantly more than that used by 1-D CNN. The fastest sampling method (MOSS-5) still uses up to two orders of magnitudes longer than ours. Comparing to exact counting methods, 1-D CNN also shows better running time, and provides at least two orders of magnitudes speed-up. The BA graphs and RGG show the similar results, which are presented in Table 4 (b) and (c).

For 5-node graphlets, we do not show the results for CC2 and GUISE, because they can not produce accurate results within reasonable time. The running time of 1-D CNN is almost the same as the 4-node graphlets', while competing methods requires longer running time. Thus, our 1-D CNN provides more speed-up for 5-node graphlets. These results demonstrate that our CNN based graphlet count estimation approach offers remarkable speedup on predicting graphlet counts on synthetic graphs while still maintaining high accuracy.

4.4 Estimation on Different Graph Structures

We evaluate the effectiveness of our model on different graph structures by testing on three different random graphs, i.e., Barabasi-Albert (BA) graph, Erdos-Renyi (ER) graph, and random geometric graph (RGG), for 3,4-node graphlets. For studying the performance of our model on sparse graphs and dense graphs, we experiment on different graph densities. To show that our model can tolerate heterogeneity in dataset, we also test on datasets with different ranges of graph densities.

4.4.1 Experiments on Different Graph Densities

We train models for different graph densities: for BA graph, we alter the *edge attachment rate* m from 3 to 7; for ER graph with test on graphs with the *edge existing probability* q from 0.1 to 0.7; for RG graph, we experiment with the *radius* r ranging from 0.25 to 0.5.

We record results for all 3, 4-node graphlets in Table 6. The results are summarized as follows. In Table 6, our model’s mean absolute errors (MAE) of estimating 3, 4-node graphlet counts on three kinds of random graphs all maintain in a small magnitude while the graph densities increase and the ground truth graphlet counts grow rapidly. Take estimating open triangle count (g_1^3) on Erdos-Renyi (ER) graph as an example. As indicated in Table (a), the mean of truth count grows from 524.15 to 7328.57 as the edge existing probability grows from 0.1 to 0.5. Nevertheless, our MAE rises from 10.31 to only 46.77. Looking at the results for 5-node graphlets in Figure 10, we have similar trends for all random graphs. In our experiments, we find existing graphlet counting methods generally take more time counting denser graphs than sparse graphs of the same size. This experiment not only demonstrates that our framework learns on different graph structures, but also shows that our model has advantage on dense graphs (with smaller relative error).

Table 6: Estimation on Graphs of Different Random Graph Models with Different Densities

(a) ER Graph									
Configu- ration		g_1^3 	g_2^3 	g_1^4 	g_2^4 	g_3^4 	g_4^4 	g_5^4 	g_6^4
0.1	RE(%)	1.9±0.1	17.6±1.0	3.6±0.2	5.2±0.3	11.0±1.2	13.4±0.7	36.3±1.9	100.0±0.0
	MEAN	524.15	19.23	1990.21	662.95	54.21	218.60	12.09	0.22
0.3	RE(%)	0.8±0.0	2.0±0.1	1.0±0.0	2.6±0.1	3.5±0.4	1.1±0.1	2.4±0.1	10.7±1.1
	MEAN	3731.11	535.22	25818.33	8591.37	2787.63	11104.68	2392.26	170.21
0.5	RE(%)	0.6±0.03	0.7±0.03	0.4±0.01	2.7±0.16	2.7±0.30	0.3±0.02	0.5±0.03	2.8±0.33
	MEAN	7328.57	2433.45	43140.88	14377.56	10734.33	42973.66	21378.36	3554.31
(b) BA Graph									
Configu- ration		g_1^3 	g_2^3 	g_1^4 	g_2^4 	g_3^4 	g_4^4 	g_5^4 	g_6^4
3	RE(%)	1.3±0.1	6.7±0.4	3.0±0.2	4.6±0.2	12.0±0.7	5.3±0.2	9.9±0.6	35.0±3.6
	MEAN	875.33	55.07	3345.27	2464.36	123.19	1022.29	114.44	4.26
5	RE(%)	1.0±0.1	3.1±0.2	1.9±0.1	3.1±0.1	7.2±0.4	2.2±0.1	4.0±0.2	12.0±1.2
	MEAN	1821.71	200.81	8884.61	5300.15	555.42	3968.82	732.81	58.06
7	RE(%)	0.8±0.0	1.7±0.1	1.4±0.1	2.3±0.1	5.3±0.3	1.2±0.1	1.9±0.1	6.0±0.7
	MEAN	2841.96	442.12	15129.29	8534.06	1303.88	8865.07	2219.75	244.66
(c) RGG Graph									
Configu- ration		g_1^3 	g_2^3 	g_1^4 	g_2^4 	g_3^4 	g_4^4 	g_5^4 	g_6^4
0.25	RE(%)	8.2±0.4	3.7±0.2	15.6±0.9	22.6±2.3	42.1±4.1	8.8±0.4	5.6±0.2	11.6±1.2
	MEAN	550.49	319.85	1701.14	186.89	15.09	1452.20	551.01	326.45
0.35	RE(%)	3.9±0.2	1.9±0.1	6.2±0.3	10.6±1.2	25.8±2.4	3.7±0.2	3.8±0.2	5.3±0.5
	MEAN	1626.34	1057.27	7982.69	925.16	84.86	7857.52	3210.44	2051.50
0.45	RE(%)	2.5±0.1	1.0±0.1	4.0±0.2	7.2±0.9	19.3±2.2	2.1±0.1	3.1±0.2	2.8±0.3
	MEAN	3291.14	2549.46	18538.42	2446.09	277.45	24321.43	11258.74	8125.60

4.4.2 Experiments on mixed graph configurations

To evaluate the generality of our models, we trained 1-D CNNs with mixed graph configuration (density) in the dataset: for BA graph, we adjust the mixing *edge attachment rate* m range from 7 ± 1 to 7 ± 5 ; for ER graph, we alter the mixing *edge existing probability* q range from 0.4 ± 0 to 0.4 ± 0.2 ; for RGG, we experiment with the mixing *radius* r range from 0.4 ± 0.05 to 0.4 ± 0.25 .

Table 7 records the results for all 3, 4-node graphlets. Our model’s mean absolute errors (MAE) of estimating 3, 4-node graphlet counts on three kinds of random graphs all maintain in a small magnitude while the graph diversity increase.

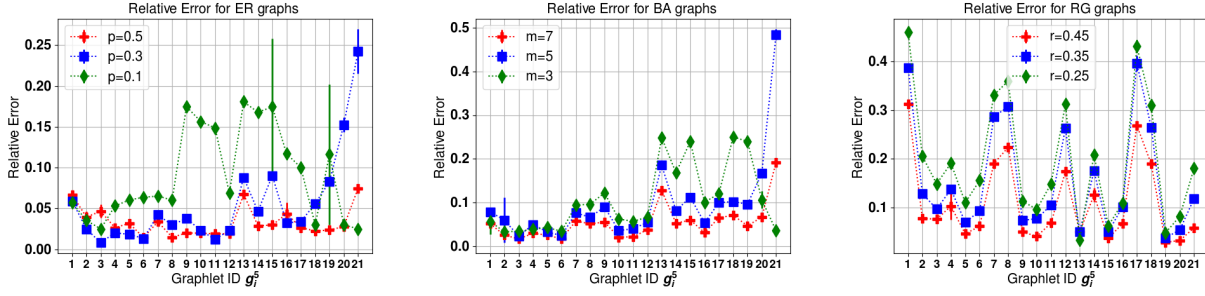


Figure 10: Relative Errors for 5-node Graphlets with Different Densities

For example, for estimating triangle count in BA graph, our model’s MAE raises from 9.26 to 9.66 only, while the standard deviation of ground truth counts in the testing dataset raises from 119.89 to 452.62. Similar results can be found for 5-node graphlets in Figure 11. This means our framework is accurate and has excellent generalization ability on mixed graph configurations.

Table 7: Estimation on Datasets with Different Density Ranges

(a) ER Graph									
Configuration		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
0.4 ± 0.0	RE(%)	0.7 ± 0.0	1.1 ± 0.0	0.6 ± 0.0	2.2 ± 0.1	2.8 ± 0.2	0.4 ± 0.0	1.0 ± 0.1	4.4 ± 0.3
	MEAN	5636.8	1252.9	38143.9	12701.1	6348.3	25433.5	8472.6	944.1
	STD	265.70	134.46	1467.43	573.45	631.39	2379.13	1322.36	214.73
0.4 ± 0.1	RE(%)	0.9 ± 0.0	1.1 ± 0.1	3.2 ± 0.1	3.7 ± 0.2	3.0 ± 0.2	0.5 ± 0.0	0.9 ± 0.0	4.4 ± 0.2
	MEAN	5472.2	1259.0	36314.1	12110.9	6185.0	24798.5	8891.1	1111.3
	STD	1050.80	534.27	5476.89	1888.30	2323.41	9210.34	5348.80	926.91
0.4 ± 0.2	RE(%)	0.9 ± 0.1	0.8 ± 0.0	10.5 ± 0.6	3.3 ± 0.2	3.0 ± 0.1	1.1 ± 0.1	0.6 ± 0.0	2.6 ± 0.1
	MEAN	5641.7	1671.9	33977.3	11313.1	7153.2	28745.5	14100.9	2570.7
	STD	2026.38	1225.85	9717.21	3227.68	4411.23	17768.77	12955.35	3052.00
(b) BA Graph									
Configuration		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
$7+1$	RE(%)	0.1 ± 0.0	2.0 ± 0.1	1.5 ± 0.1	2.9 ± 0.1	5.2 ± 0.2	1.6 ± 0.1	2.6 ± 0.1	7.1 ± 0.4
	MEAN	2879.2	462.5	15333.0	8598.3	1362.4	9222.7	2391.3	277.6
	STD	414.86	119.89	2615.75	1532.10	387.68	2371.10	846.62	129.45
$7+3$	RE(%)	1.0 ± 0.1	2.2 ± 0.1	2.2 ± 0.1	2.9 ± 0.2	6.1 ± 0.3	1.5 ± 0.1	3.0 ± 0.1	6.4 ± 0.3
	MEAN	2708.4	449.7	14107.8	8129.6	1303.2	8864.6	2468.4	316.9
	STD	951.73	270.15	5632.95	3061.57	819.18	5243.69	1940.83	311.72
$7+5$	RE(%)	0.8 ± 0.0	1.8 ± 0.1	2.8 ± 0.2	3.0 ± 0.1	4.9 ± 0.3	1.2 ± 0.1	2.2 ± 0.1	6.8 ± 0.4
	MEAN	2776.9	543.2	14375.8	8274.2	1594.3	10404.1	3455.0	522.8
	STD	1569.41	452.62	8864.47	4857.48	1345.40	8452.05	3441.41	603.70
(b) RGG Graph									
Configuration		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
0.4 ± 0.05	RE(%)	3.8 ± 0.2	1.4 ± 0.1	5.7 ± 0.3	8.4 ± 0.4	20.8 ± 2.3	3.3 ± 0.1	3.4 ± 0.2	4.0 ± 0.4
	MEAN	2433.7	1753.8	13282.8	1643.5	173.1	15323.7	6758.0	4668.1
	STD	665.18	592.16	4119.02	690.95	92.10	6520.78	3403.68	2717.33
0.4 ± 0.15	RE(%)	3.6 ± 0.2	1.3 ± 0.1	5.1 ± 0.3	9.4 ± 0.4	22.3 ± 2.0	2.8 ± 0.2	3.4 ± 0.2	3.1 ± 0.3
	MEAN	2403.3	1853.2	12356.5	1690.6	196.8	16797.0	8076.9	5982.6
	STD	1387.38	1393.33	7556.02	1285.88	192.68	14378.81	8281.93	6999.81
0.4 ± 0.25	RE(%)	0.3 ± 0.0	0.9 ± 0.0	3.9 ± 0.2	9.0 ± 0.4	20.6 ± 2.1	2.4 ± 0.1	2.7 ± 0.1	2.2 ± 0.2
	MEAN	2568.6	2310.7	12084.4	1781.3	284.8	20793.9	11645.3	9695.2
	STD	1987.53	2257.41	9347.08	1593.88	332.00	20840.74	13613.36	12868.71

4.4.3 Experiments on mixed graph distributions

In section 4.4.2, we train three separate 1-D CNNs for three synthetic datasets with mixed graph configurations. In this section, we study whether training a **single** 1-D CNN with these three datasets mixed together can get accurate

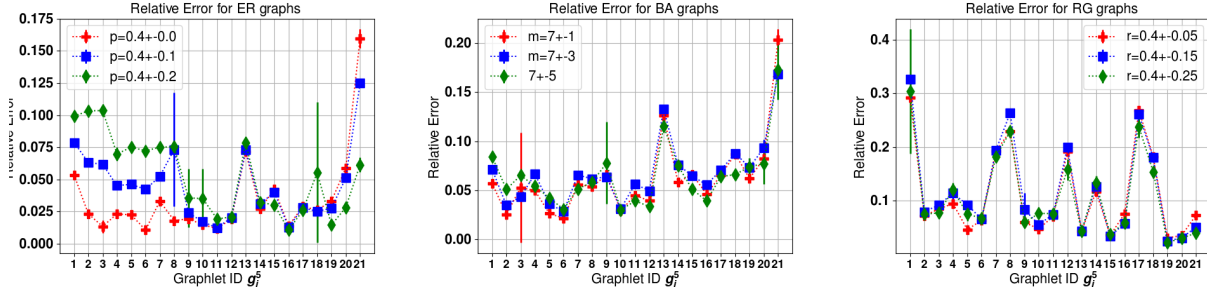


Figure 11: Relative Errors for 5-node Graphlets with Different Density Ranges

Table 8: Estimation on Mixed Graph Distributions

	g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Relative Error(%)	6.98	2.82	6.47	9.58	12.90	2.87	3.34	7.27
STD of the Relative Error (%)	0.67	0.25	0.55	0.11	1.72	0.17	0.10	0.08

predictions for graphs generated from this mixed graph distribution. We randomly generate 3000 graphs as training data, including 1000 ER graphs with p ranges from 0.05 to 0.7, 1000 BA graphs with m ranges from 1 to 13 and 1000 RGG with r ranges from 0.15 to 0.65. We randomly generate another 450 graphs from the above mixed distribution for testing. As shown in the Table 8, our model can adapt itself to this mixed distribution setting and predict unseen graphs generated from this mixed distribution with fairly acceptable accuracy (all below 10% except for 4-cycle). As for the 5-node graphlets, Figure 12 shows more than half of the relative errors are below 10% and 18/21 are below 15%. This shows our model has the potential to be applicable to different underlying distributions with a single trained model.

4.5 Estimation on Large Random Graphs

To demonstrate the scalability of our framework and how our CNN with RNS behave on large graphs, we do experiments on 10K-node random graphs. We randomly generate two 10k-node random ER graphs with edge existing probability $q = 0.1$, two 10k-node random BA graphs with edge attachment rate $m = 200$ and two 10k-node random RGG graphs with radius $r = 0.15$. For ER graphs, we apply RNS with node preserving probability $p = 0.015$ for one graph to generate 3000 training graphs and 3000 for the other graph with the same p for testing. For BA graphs, we apply RNS with node preserving probability $p = 0.035$ for one graph to generate 5000 training graphs and 2000 testing graphs with the same p for the other graph. For RGG, we apply RNS with $p = 0.015$ to generate 3000 graphs for training and 10000 graphs for testing. Figure 13 shows the results of our experiment for 3,4-node graphlets, where blue lines denote the ground truth graphlet counts (note that there is no blue bars as the relative error of ground truth is 0); green lines present the graphlet count estimation for large random graphs by RNS with exact counting on RNS samples and green bars denote the relative error incurred by RNS; red lines and bars represent the result by applying RNS with our CNN framework on RNS samples. In Figure 13, our estimations are accurate as they (red lines) almost coincide with the ground truths (blue lines) and the relative errors of RNS with CNN are close to RNS with exact counting on all random graphs. We also provide the training time and testing time for RNS + CNN, the running time for RNS + Exact Count and the running time for PGD exact to calculate the ground truth of all 10k graphs in Table 9. For ER and BA graphs, we can finish testing in less than 0.62 seconds. For RGG, we spend more time (less than 1.21 seconds) because we sample more RNS samples than ER and BA graphs. We show results for 5-node graphlets on the relative error in Figure 14. Again, RNS+CNN shows consistent accuracy (around 6-8% relative error) for most 5-node graphlets as 4-node graphlets', except for g_{18}^5 - g_{21}^5 for BA graphs. For these 5-node graphlets, the large error suggests we may need more RNS samples to produce more accurate results. For the running time shown in Table 10, our RNS+CNN shows clear advantage over the ESCAPE methods and can finish within 1.2 seconds.

4.6 Preprocessing Techniques Evaluation

We conduct experiments to analyze the effects of our proposed preprocessing techniques.

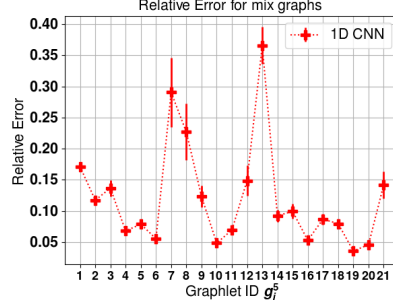
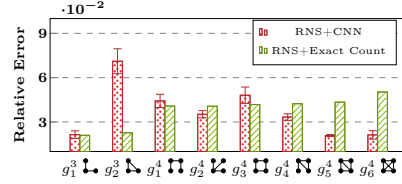
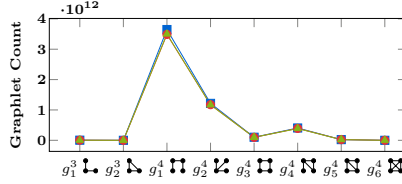


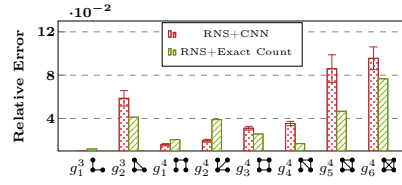
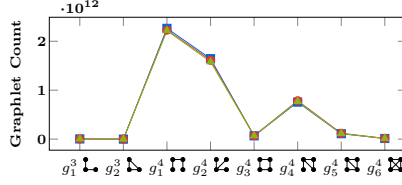
Figure 12: Relative Errors for 5-node Graphlets for Mixed Random Graphs

- ■ - GLOBAL GROUND TRUTH - ● - RNS+CNN - ▲ - RNS+EXACT COUNT

(a) 10000-node ER Graph



(b) 10000-node BA Graph



(c) 10000-node RG Graph

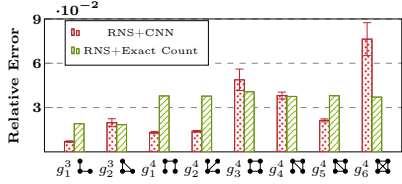
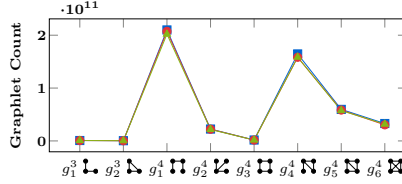


Figure 13: 3-4 node Graphlet Counts on Large Random Graphs

4.6.1 Node Ordering

We evaluate four kinds of node ordering methods: betweenness centrality, closeness centrality, Weisfeiler-Lehman algorithm [43], and degree centrality on all 3,4-node graphlets. We show the results in Table 11 and mark the most accurate results with bold faces. In our experiments, we observe that betweenness centrality and closeness centrality are not very effective. Among all the evaluated ordering methods, degree centrality provides the biggest improvements on estimating almost all 3 and 4-node graphlet counts on Erdos-Renyi (ER) graphs, Barabasi-Albert (BA) graphs. For estimating 4-clique count on RGGs, the degree centrality ordering reduces 5.73% on the relative error. The Weisfeiler-Lehman ordering achieves the best performance on random geometric graphs (RGG) and is also shown to be effective on ER graphs and BA graphs. However, given that the complexity of Weisfeiler-Lehman ordering is much higher than degree centrality, it is more suitable to use degree centrality as the first choice. We also validate the above observation for 5-node graphlets. Figure 15 shows degree centrality almost always outperforms other node orderings in all graphs.

Table 9: Running Time of 3,4-node Graphlets for Large Random Graphs

(a) ER Graph		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Training (in Minutes)	CNN	1.25	1.63	6.35	4.20	1.72	2.48	1.28	0.55
Testing (in Seconds)	RNS + CNN	0.285	0.355	0.355	0.349	0.349	0.269	0.294	0.424
	RNS + EXACT	180.2	180.2	180.2	180.2	180.2	180.2	180.2	180.2
	PGD EXACT	456.9	456.9	456.9	456.9	456.9	456.9	456.9	456.9
	ESCAPE	881.01	881.01	881.01	881.01	881.01	881.01	881.01	881.01
(b) BA Graph		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Training (in Minutes)	CNN	6.93	6.28	12.13	12.03	5.86	32.23	17.45	9.55
Testing (in Seconds)	RNS + CNN	0.770	0.759	0.620	0.630	0.770	0.640	0.640	0.629
	RNS + EXACT	180.2	180.2	180.2	180.2	180.2	180.2	180.2	180.2
	PGD EXACT	131.9	131.9	131.9	131.9	131.9	131.9	131.9	131.9
	ESCAPE	310.89	310.89	310.89	310.89	310.89	310.89	310.89	310.89
(b) RGG		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
Training (in Minutes)	CNN	1.50	1.26	1.80	1.03	0.78	1.73	1.85	1.80
Testing (in Seconds)	RNS + CNN	0.934	0.890	1.170	1.165	0.894	1.164	1.195	1.210
	RNS + EXACT	259.4	259.4	259.4	259.4	259.4	259.4	259.4	259.4
	PGD EXACT	228.7	228.7	228.7	228.7	228.7	228.7	228.7	228.7
	ESCAPE	2372.20	2372.20	2372.20	2372.20	2372.20	2372.20	2372.20	2372.20

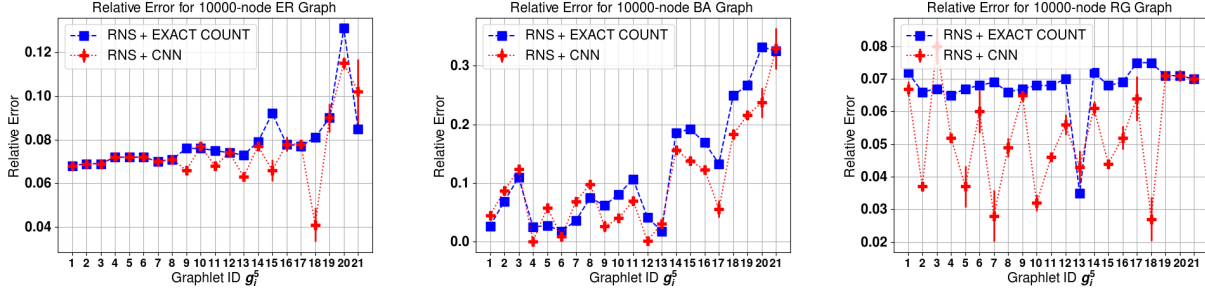


Figure 14: Relative Errors for 5-node Graphlets for Large Random Graphs

4.7 Estimation on Real World Graph

In this subsection, we demonstrate that our framework can be applied to real world scenarios by experimenting on three categories of real world graphs: biochemistry graphs, collaboration graphs and social network graphs.

4.7.1 Biochemistry Dataset

We train our 1-D CNN models for estimating **non-zero** graphlet counts, including, g_1^4 ($\bullet\bullet$), g_2^4 ($\bullet\bullet$), g_1^5 ($\bullet\bullet$), g_2^5 ($\bullet\bullet$), on three biochemistry datasets: MUTAG, NCI1, and NCI109. Table 2 gives the detailed information of these datasets.

Because the graph samples in these real world datasets have different number of nodes, we preprocess the input graph adjacency matrices with *adjacency matrix zero padding* technique as stated in Section 3.2.1. For NCI1 and NCI109, since the size of their graph samples vary over a large range, e.g., the average number of nodes in a graph is around 30 but the biggest graph has 111 nodes, padding the graph adjacency matrices to size (111, 111) adds many useless information to the inputs. Hence, we further apply *node ordering* to obtain more consistent representation of the graph samples. We summarize our results in Table 12 as follows. Our model performs well on MUTAG dataset. The relative error (MAE over Mean) of all estimations are below 6.3%, and for estimating open-triangle (g_1^3 $\bullet\bullet$), our model achieves 0.9% relative error. On NCI1 and NCI109, the relative error 1 of our estimations are all less than 7.7%. Our model performs better on estimating open-triangle (g_1^3 $\bullet\bullet$) on both NCI1 and NCI109, making only 0.09% and 1% relative error respectively. This experiment shows that our framework is effective on real world biochemistry datasets.

Table 10: Running Time for Large Random Graphs for 5-node graphlets

(a) ER Graph											
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
CNN Train (Ours)	3.441	3.889	3.981	3.888	3.795	3.700	3.574	4.951	3.668	4.042	3.581
RNS + CNN (Ours)	0.170	0.185	0.170	0.190	0.175	0.190	0.170	0.195	0.210	0.195	0.165
RNS + EXACT	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04
ESCAPE	242820	242820	242820	242820	242820	242820	242820	242820	242820	242820	242820
	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
CNN Train (Ours)	3.744	4.140	3.513	3.782	3.734	3.605	3.500	4.295	3.797	2.188	
RNS + CNN (Ours)	0.165	0.175	0.185	0.180	0.170	0.200	0.195	0.180	0.200	0.180	
RNS + EXACT	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	298.04	
ESCAPE	242820	242820	242820	242820	242820	242820	242820	242820	242820	242820	
(b) BA Graph											
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
CNN Train (Ours)	6.390	7.229	6.970	6.766	5.393	6.033	6.419	6.116	5.654	6.194	6.263
RNS + CNN (Ours)	0.615	0.640	0.670	0.655	0.575	0.620	0.595	0.645	0.630	0.710	0.670
RNS + EXACT	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69
ESCAPE	36255	36255	36255	36255	36255	36255	36255	36255	36255	36255	36255
	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
CNN Train (Ours)	5.715	5.559	5.966	5.949	5.644	5.998	5.496	6.073	6.022	5.446	
RNS + CNN (Ours)	0.595	0.575	0.665	0.645	0.610	0.680	0.530	0.700	0.705	0.575	
RNS + EXACT	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	258.69	
ESCAPE	36255	36255	36255	36255	36255	36255	36255	36255	36255	36255	
(c) RGG Graph											
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
CNN Train (Ours)	2.833	1.593	2.612	2.481	1.581	2.550	2.581	2.438	2.414	2.605	2.402
RNS + CNN (Ours)	0.875	1.120	1.120	1.125	1.140	0.960	1.040	1.050	0.875	0.915	0.920
RNS + EXACT	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4
ESCAPE	273810	273810	273810	273810	273810	273810	273810	273810	273810	273810	273810
	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
CNN Train (Ours)	2.607	1.554	3.068	2.626	1.555	2.557	2.523	2.551	2.756	2.520	
RNS + CNN (Ours)	1.110	1.124	1.110	0.925	1.124	1.035	1.110	0.925	0.875	0.875	
RNS + EXACT	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	1289.4	
ESCAPE	273810	273810	273810	273810	273810	273810	273810	273810	273810	273810	

Note: Bold faces are the best results. Time for training are in minutes, others are in seconds.

4.7.2 Collaboration Dataset

To demonstrate our framework's capability of handling large real world graph, we evaluate our framework using two collaboration networks: ca-hep-ph and ca-hep-th. Table 3 gives details about these two graphs. Following the problem formulation in Section 2.3 to generate a dataset \mathbb{G} , we use the ca-hep-ph graph to generate 4000 training/validation graphs as \mathbb{G}_1 to train our model. Then we generate 900 (unseen) testing graphs from the ca-hep-th as \mathbb{G}_2 . Specifically, we apply the *random node sketching* approach introduced in Section 3.2.4 to sample 4000 graphs from ca-hep-ph with the *node preserving probability* $p = 0.01$ to form a training set with size 4000 (among them 400 graphs are reserved for validation) ; and sample 900 graphs for testing from ca-hep-th with the node preserving probability $p = 0.009$. Since the sizes of the sampled small graphs vary, we further apply *adjacency matrix zero padding* to pad adjacency matrices to the same size, i.e., (280, 280) for both training graphs and testing graphs. Also, we apply degree centrality based Node Ordering technique to preprocess the training and testing graphs. Finally, as introduced in Section 3.2.4, we first take the CNN model's estimation on 900 testing graphs as inputs and reconstruct the global graphlet count estimation for ca-hep-th by using the RNS estimator described by Equation (7). We denote the results as "RNS + CNN" (red bars in Figure 16(a)).

Table 11: Relative Error for Different Node Orderings

(a) ER Graph								
Relative Error(%)	g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
No ordering	1.42	3.51	2.98	4.88	5.84	1.59	4.00	13.21
Betweenness	1.41	3.68	2.82	5.21	5.97	1.74	3.99	13.34
Closeness	1.44	3.46	2.82	4.94	5.77	1.69	4.22	12.27
WL	1.43	3.65	2.45	4.61	5.20	1.65	4.29	11.52
Degree	1.35	3.27	2.33	4.24	5.19	1.65	3.57	11.01
(b) BA Graph								
Relative Error(%)	g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
No ordering	1.53	3.15	5.05	8.19	10.63	1.59	4.74	14.02
Betweenness	1.52	2.86	4.86	8.63	10.70	1.64	3.74	11.29
Closeness	1.66	2.82	4.50	8.00	10.00	1.97	3.66	11.45
WL	1.41	2.54	3.20	4.73	9.02	1.65	4.07	10.14
Degree	1.31	2.31	2.36	3.26	7.13	1.83	2.71	8.29
(c) RGG								
Relative Error(%)	g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
No ordering	9.35	5.80	18.90	29.00	54.90	11.56	7.68	19.25
Betweenness	9.76	5.74	18.92	30.89	50.62	11.09	7.62	18.65
Closeness	9.50	5.69	19.17	29.73	50.14	10.60	7.66	19.35
WL	8.17	4.33	16.47	27.19	50.24	9.08	7.81	13.80
Degree	8.91	4.53	17.71	28.03	51.16	9.61	6.88	14.85

Note: Bold faces are the best results.

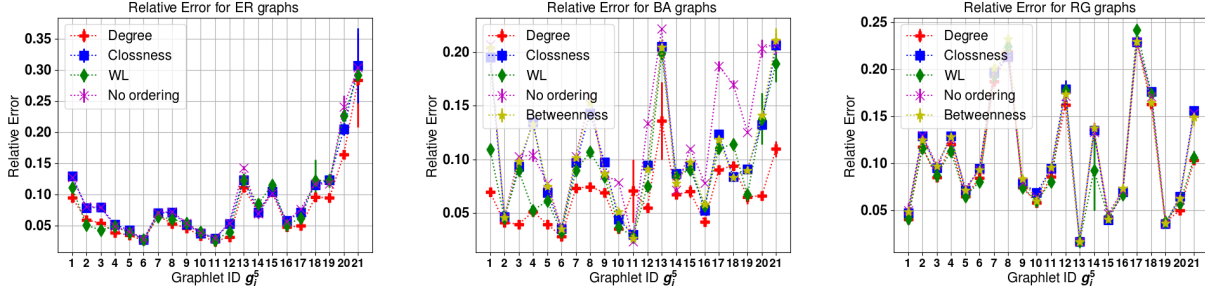


Figure 15: Relative Errors for 5-node Graphlets with Different Node Orderings

For 3,4-node graphlets, we first calculate the results by reconstructing the global count estimation using the ground truth count on 900 testing graphs, which is denoted as "RNS + Exact Count" (green bars). Also, we use LGE with 660 edges (including its 2-hop egonets) for 7 runs and denote the result as "LGE" (purple bars). We also run MOSS-5 and ApproxG with controlled accuracy close to the RNS+CNN, which are denoted as "MOSS-5" (black bars) and "ApproxG" (orange bars). Note that we tried GRAFT, GUISE and CC2 on ca-hep-th graph as well, but we do not present them because they did not terminate within 72 hours even with very low sampling ratio.

In Figure 16(a), we have larger error on 3-star (g_2^4 \star) count, we analyze the reason as follows. First, as the results recovered from sample ground truth (green bars in Figure 16(a)) also have rather big error on 3-star, the 900 small testing graphs we sample seem to be less statistically representative for 3-star count in ca-hep-th. Second, we observe, in the testing set, the variance of 3-star graphlet count is large, i.e., a small number of graph samples with significantly larger 3-star counts than most of the graph samples in the set. The error mainly comes from the underestimation on these outlying samples with large count. This is also the reason why LGE does not perform well in estimation 3-star graphlet count. Except for 3-star, Figure 16(a) shows that, for all other 3-node and 4-node graphlets, the relative error of our estimations are all below 8.7%. These results are close to the best results we can achieve theoretically (the green line and bars) given that we only sample 900 graphs. The results can be further improved by taking more samples from ca-hep-th. Compared with competing methods, RNS+CNN is superior in the accuracy for graphlet g_2^3 (\star), g_1^4 (\star), g_5^4 (\star) and g_6^4 (\star) and on par or slightly worse for g_1^3 (\star), g_3^4 (\star) and g_4^4 (\star). As for the running time shown in Table 13 (a), RNS+CNN spends less than 0.31 seconds for all 3 and 4-node graphlets and is the fastest among competing methods.

Table 12: Graphlet Count Estimation Results on Real World Biochemistry Datasets

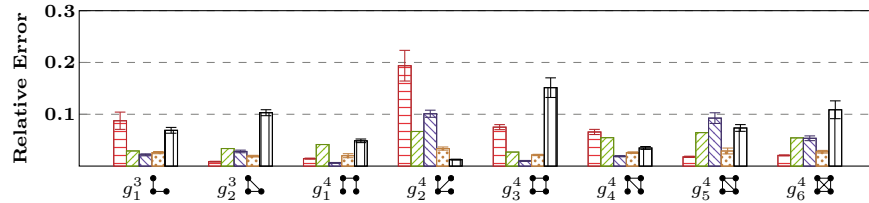
(a) MUTAG					
	g_1^3	g_1^4	g_2^4	g_1^5	g_2^5
MAE	0.29	1.37	0.41	2.62	2.05
Mean	30.75	43.25	7.8	58.7	32.2
STD	9.62	16.10	2.73	27.49	13.92
Relative Error	0.009 ± 0.0012	0.031 ± 0.0030	0.052 ± 0.0028	0.045 ± 0.0017	0.063 ± 0.0014
(b) NCI1					
	g_1^3	g_1^4	g_2^4	g_1^5	g_2^5
MAE	0.51	2.27	0.37	6.97	4.20
Mean	51.26	70.27	13.84	90.94	54.32
STD	24.46	35.53	8.53	50.76	36.10
Relative Error	0.009 ± 0.0006	0.032 ± 0.0010	0.027 ± 0.0029	0.076 ± 0.0038	0.077 ± 0.0018
(c) NCI109					
	g_1^3	g_1^4	g_2^4	g_1^5	g_2^5
MAE	0.48	2.17	0.42	6.51	3.84
Mean	48.04	65.56	13.01	85.01	50.81
STD	25.08	36.48	8.39	51.66	34.48
Relative Error	0.010 ± 0.00097	0.032 ± 0.00048	0.032 ± 0.0016	0.076 ± 0.0032	0.075 ± 0.0031

MOSS-5 is the closest to ours and the PGD exact method even consumes 199.92 seconds, and ours is about 318 times faster than theirs.

For 5-node graphlets, we compare our 1-D CNN with MOSS-5 and ESCAPE. Note that results for GRAFT, CC2 and GUISE are not shown because they can not finish within reasonable time. Look at the relative error shown on the left in Figure 17, most 5-node graphlets have relative error less than 7.5%. As for the running time shown in Table 14 (a), our RNS+CNN achieves at least 3.4 times speed-up compared with MOSS-5, and five orders of magnitude faster than ESCAPE. More importantly, the running time remains almost the same as the 3,4-node graphlets', while exact counting methods (i.e., ESCAPE) grows exponentially with the graphlet size. This experiment demonstrates that our framework is effective in real world scenario for large collaboration graphs.

■ RNS+CNN ■ RNS+EXACT COUNT ■ LGE ■ APPROXG ■ MOSS-5

(a) Errors on Collaboration Networks



(b) Errors on Social Networks

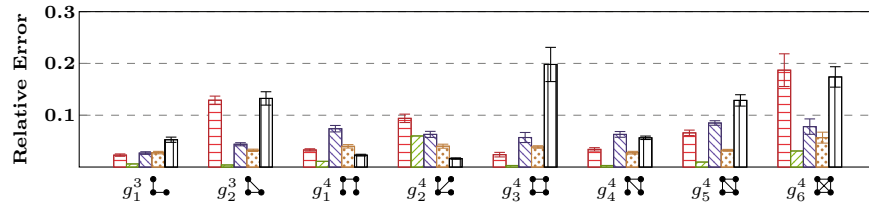
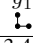
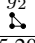
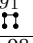
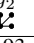
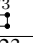
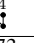
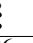
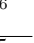
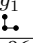
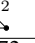




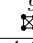



Figure 16: 3,4-node Graphlet Count Estimation Results on Real World Graphs

Table 13: Running Time for 3,4-node Graphlet Count Estimation on Real World Graphs

(a) Collaboration Networks		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
									
Training (in Minutes)	RNS+CNN (Ours)	3.46	5.20	2.98	5.93	2.23	6.73	3.86	2.75
Testing (in Seconds)	RNS+CNN (Ours)	0.23	0.31	0.31	0.25	0.31	0.31	0.23	0.28
	RNS + EXACT	167.7	167.7	167.7	167.7	167.7	167.7	167.7	167.7
	GRAFT	-	-	-	-	-	-	-	-
	CC2	-	-	-	-	-	-	-	-
	GUISE	-	-	-	-	-	-	-	-
	LGE	12.53	11.75	12.72	13.55	12.72	12.72	12.72	12.72
	ApproxG	0.98	0.98	1.02	1.02	1.25	1.02	1.25	1.02
	MOSS-5	0.55	0.55	0.63	0.62	0.78	0.62	0.62	0.88
	PGD EXACT	199.92	199.92	199.92	199.92	199.92	199.92	199.92	199.92
	ESCAPE	872.92	872.92	872.92	872.92	872.92	872.92	872.92	872.92
(b) Social Networks		g_1^3	g_2^3	g_1^4	g_2^4	g_3^4	g_4^4	g_5^4	g_6^4
									
Time (in Minutes)	RNS+CNN (Ours)	3.96	2.72	4.13	2.88	5.20	6.07	3.51	4.45
Testing (in Seconds)	RNS+CNN (Ours)	0.487	0.476	0.440	0.456	0.564	0.310	0.445	0.483
	RNS + EXACT	18.25	18.25	18.25	18.25	18.25	18.25	18.25	18.25
	GRAFT	-	-	-	-	-	-	-	-
	CC2	-	-	-	-	-	-	-	-
	GUISE	-	-	-	-	-	-	-	-
	LGE	0.865	0.827	0.865	0.865	0.865	0.865	0.827	0.865
	ApproxG	0.49	0.49	0.56	0.56	0.49	0.56	0.49	0.49
	MOSS-5	2.19	2.23	2.60	2.22	5.21	2.75	5.35	4.99
	PGD EXACT	3.42	3.42	3.42	3.42	3.42	3.42	3.42	3.42
	ESCAPE	10.82	10.82	10.82	10.82	10.82	10.82	10.82	10.82

Note: a hyphen (-) indicates that the method did not terminate within 4 hours. Bold faces are the best results.

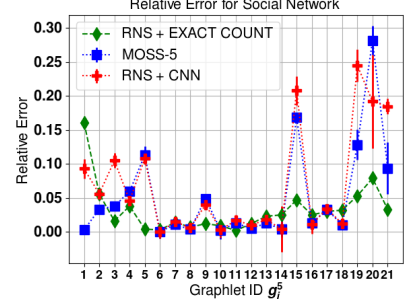
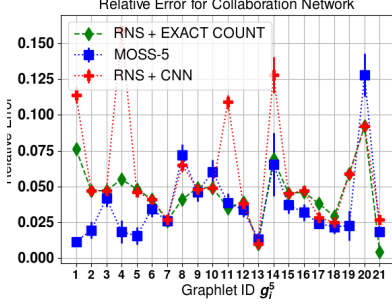


Figure 17: 5-node Graphlet Count Estimation Results on Real World Graphs

4.7.3 Social Network Dataset

We also evaluate our framework using two sparse social networks: socfb-MSU24 and socfb-FSU53. Similar to Section 4.7.2, we first use RNS to generate 3300 graphs (3000 for training, 300 for validation) with node preserving probability $p = 0.02$ from socfb-MSU24. Then 500 testing graphs are generated using socfb-FSU53 with the same node preserving probability. Next, we pad training, validation and testing graphs to (560, 560) and apply the degree centrality based node ordering technique. Finally, we use Equation (7) to calculate the estimated graphlet counts for socfb-FSU53.

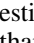
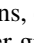
For 3,4-node graphlets, we show the relative error results for our method and competing methods in Figure 16(b). The RNS+Exact Count method achieves the best accuracy among competing methods and our RNS+CNN method gets fairly accurate estimations, except for 3-triangle () and 4-clique (). The reason is that 3-triangle and 4-clique counts are much less than other graphlets in sparse graphs. With inadequate samples, our CNN model can not learn well on structural patterns of the 3-triangle and 4-clique. For the running time (shown in Table 13(b)), our method is still the fastest for most 3, 4-node graphlets. But this time, we provide less than 10 times speed up compared to the PGD exact method, and the ApproxG, MOSS-5 are also very close to RNS+CNN. In our experiments, we find most existing

Table 14: Running time for 5-node Graphlet Count Estimation on Real World Graphs

(a) Collaboration Networks											
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
CNN Train (Ours)	3.675	3.694	3.450	3.714	3.435	3.739	3.679	3.262	3.729	3.712	3.449
RNS + CNN (Ours)	0.232	0.248	0.190	0.243	0.210	0.237	0.248	0.258	0.238	0.240	0.207
RNS + EXACT	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3
MOSS-5	0.800	0.790	0.660	0.800	0.800	0.940	1.030	1.030	1.030	0.670	0.800
ESCAPE	130733	130733	130733	130733	130733	130733	130733	130733	130733	130733	130733
	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
CNN Train (Ours)	3.427	3.285	3.456	3.730	3.748	3.390	3.158	3.495	3.705	3.465	
RNS + CNN (Ours)	0.243	0.245	0.203	0.235	0.242	0.240	0.195	0.205	0.240	0.202	
RNS + EXACT	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	153.3	
MOSS-5	0.940	1.750	0.800	0.670	0.660	1.750	0.790	0.670	0.800	0.940	
ESCAPE	130733	130733	130733	130733	130733	130733	130733	130733	130733	130733	
(b) Social Networks											
	g_1^5	g_2^5	g_3^5	g_4^5	g_5^5	g_6^5	g_7^5	g_8^5	g_9^5	g_{10}^5	g_{11}^5
CNN Train (Ours)	6.269	5.513	5.250	6.279	8.470	5.311	6.288	10.416	5.286	10.556	5.342
RNS + CNN (Ours)	0.525	0.435	0.475	0.550	0.805	1.160	0.548	1.138	0.480	0.780	1.202
RNS + EXACT	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2
MOSS-5	1.370	0.330	0.310	0.610	0.310	13.740	1.730	13.740	1.370	13.740	1.370
ESCAPE	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5
	g_{12}^5	g_{13}^5	g_{14}^5	g_{15}^5	g_{16}^5	g_{17}^5	g_{18}^5	g_{19}^5	g_{20}^5	g_{21}^5	
CNN Train (Ours)	5.371	9.912	6.296	10.614	6.375	5.309	6.485	6.579	6.683	6.803	
RNS + CNN (Ours)	0.518	0.817	0.625	0.457	0.565	0.463	0.527	0.605	0.558	0.418	
RNS + EXACT	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	125.2	
MOSS-5	1.370	4.620	55.790	0.620	1.370	0.610	55.790	0.610	0.620	0.520	
ESCAPE	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	584.5	

Note: Bold faces are the best results. Time for training are in minutes, others are in seconds.

methods are more efficient for sparser graphs than denser ones given their graph sizes do not vary too much. So for sparse graphs, the competing methods are more competitive.

For 5-node graphlets, we show the relative error in the right hand side of Figure 17. We have less than 5% error for half of the graphlets, and the largest error are clique-like graphlets, i.e., g_{19}^5 , g_{20}^5 and g_{21}^5 . The reason is also due to the inadequate samples for these graphlets. As for the running time, we show the results in Table 14. For most graphlets, our RNS+CNN method are superior than the competing methods. Again, the running time This experiment shows our method is effective and competitive against state-of-the-art methods for sparse social networks.

Interestingly, we find the relative error curve of the social network is very similar to large random BA graphs in Figure 14, and the collaboration network are very similar to large RG graphs. This further demonstrates the performance of our model are related to the underlying distribution. Moreover, we find the results for real-world networks are slightly worse than random graphs shown in Figure 14. We believe the reason is related to the generalizability of our CNN model. For random graphs, we generate training and testing graphs from exactly the same distribution, i.e., using the same or similar configurations. For real-world networks, the training and testing graph do not come from the same distribution, so our CNN model has larger error. To solve this problem, we can train a mixture model using more networks, as what we do in the section 4.4.3, to enhance the generalizability of our model.

5 Related Work

5.1 Existing Graphlet Counting Framework

Existing techniques can be roughly categorized into two groups. One aims at obtaining the exact graphlet counts; the other utilizes sampling techniques and focuses more on graphlet frequency distribution (GFD).

For exact graphlet counts, the most straightforward approach is to traverse through the graph, pick any k nodes, check the edge connectivity between them and count every occurrence of a specific graphlet. Many existing works propose to accelerate this exhaustive enumeration process by leveraging on parallelism, distributed systems and combinatorial relationships among graphlets. For example, Schank et al.’s work [44] gives a MapReduce based algorithm to count

triangle on distributed systems. Ahmed et al. [33] propose the Parallel Parameterized Graphlet Decomposition (PGD) framework which counts 3-node graphlets, 4-clique, and 4-cycle first, and then derive the counts of the rest of 4-node graphlets directly from the combinatorial relationships between 3-node and 4-node graphlets. Pinar et al. [41] build upon Ahmed et al.’s idea and propose the ESCAPE framework for counting all 5-node graphlets.

The second group of works aim at a less ambitious but more practical goal: estimate graphlet count or GFD with high accuracy. They utilize various sampling techniques and provide fast and accurate algorithms with statistically provable error bounds. For example, Rahman et al. propose an edge sampling based framework called GRAFT [11]. Seshadhri et al. [12] propose a wedge sampling technique and Jha et al. [13] extend this idea by introducing path sampling with a special pruning schema to estimate all 4-node graphlets. Rossi et al. [29] extend the PGD framework and develop a Localized Graphlet estimation (LGE) Framework by computing graphlets on sampled localized neighborhoods. Mawhritter et al. [39] propose the distributed ApproxG framework to count 3, 4-node graphlets by the cost-aware task sampling with controllable accuracy. For estimating k -node graphlets (for any k), Bhuiyan et al. [38] first apply the Markov Chain Monte Carlo (MCMC) sampling method, and Chen et al. [14] adapt the MCMC sampling and estimate the ratio between the occurrence of different k -node graphlets. Wang et al. [40] divide 5-node graphlets into two groups, propose T-5 and Path-5 sampling methods for each group and develop the MOSS-5 framework to estimate 3,4,5-node graphlets. Bressan et al. [15] point out the long mixing time of MCMC methods and propose a Color Coding based method to estimate up to 7-node graphlet counts. For special graphlets, such as k -cliques, Jain and Seshadhri’s work [16] incorporates Chiba and Nishizaki’s subgraph listing algorithm [10] and Turan’s Theorem for approximating the number of k -clique (for any k).

However, existing exact counting and sampling methods only take graphs as input and mechanically output their graphlet counts, while ignoring the information from the results they just computed for previous graphs. As we discussed in Section 2.3, failing to utilize the information of historical graphs, these methods suffer loss in accuracy or efficiency or even both when estimating graphlet counts in new-coming graphs.

5.2 Learning on Graphs

Graphs are natural representations for a variety of real-world applications, including social analysis [45], traffic prediction [46], recommendation systems [47] and computer vision [48]. By representing data as graphs, on one hand, we can encode relationships among entities and give more insights in the underlying data. On the other hand, the complex structures of graphs pose a challenge to learning the true insights and extracting useful information [49]. Many traditional techniques such as graph kernel [34, 4] and graph embedding [50, 51] methods are proposed to deal with the learning problem on graphs. Graph kernel method, in general, can be used to measure similarity among graphs by computing shortest paths, subtrees or cycles [4]. However, it is often hard to priorly know which graph kernel captures key structural property better [52]. Graph embedding method, which learns a graph representation in a low dimensional space, is more task-specific. This technique has been adopted to address graph learning problems such as link prediction and node classification [50]. But its learning ability is always limited for discovering intricate network patterns because of its shallow learning mechanism.

Lately, convolutional neural networks (CNNs) have arisen in many fields and yielded promising empirical results. CNNs have demonstrated their significant representation ability, impressive learning capacity, and efficiency in space compared with traditional kernel or embedding approaches. However, most prevailing CNN models are originally developed for regular 1-D, 2-D, 3-D grid data such as text, image or video; while graph is generally of high dimensional and irregular. CNN’s key components such as convolutions and filtering cannot be generalized to the graph domain in a trivial manner. Thus, a variety of recent works were proposed to solve this problem. These works can be further divided into two categories: the spectral models and the spatial models. Spectral methods [53, 54, 55, 56] are based on the spectral graph theory: they define parameter filters and convolution operators by the graph Fourier transform. As the spectral methods relies on specific eigenfunctions of Laplacian matrix, the learned parameters on one graph are difficult to be transferred to another graph with different eigenvalues. The spatial models [57, 58, 48], which focus on generalizing CNN concepts on spatial domain of graphs, are more flexible. For example, Hamilton et al. [57] propose the GraphSAGE framework to learn node embeddings by aggregating feature information from nodes’ local neighborhoods. Velickovic et al. [58] follow Hamilton et al.’s idea and propose a novel attention layer which gives different learnable importance weights on aggregating local neighbors of a node. Besides, Monti et al. [48] also introduce a generic graph convolution network, Monet, by designing a general patch operator which integrates the neighborhood information.

The objective of our work is different from the existing graph learning frameworks. Existing frameworks focus on learning representations of data points, i.e., node embeddings, network embeddings. For these representation learning frameworks, node or edge attributes are their main learning subjects; the underlying graph structure is viewed as an auxiliary information. For our work, however, the underlying graph structure is the key to solve the GCL problem.

We aim to apply learning model to learn the mapping from underlying graph structure to the structural attributes, i.e. graphlet counts.

6 Conclusion

In this paper, we point out that, for existing graphlet counting and estimation framework, the cost for repeated computation in testing can be further reduced, if we learn from previously studied graphs and leverage on the correlation between the graph structural information and graphlet counts. we abstract the above finding, and formulate it as a novel Graphlet Count Learning (GCL) problem, which bridges the graphlet counting and the learning paradigm. We then propose a convolutional neural network (CNN) based framework with a series of data preprocessing methods to solve the GCL problem. To show the effectiveness of our CNN framework, we conduct extensive experiments on three kinds of random graphs and three categories of real world graphs, namely, biochemistry graphs, collaboration graphs and social network graphs for 3,4,5-node graphlet counting. The results demonstrates that our framework has great potential to compete with state-of-the-art methods, and is accurate, efficient, general and can achieve consistent results to handle any sized graphlets (i.e., 6-node graphlets) without changing the model.

For the future work, extending the framework to local graphlet count on directed graph is a natural next step. Also, it will be interesting to further improve the accuracy and scalability of our method to handle larger graphs (e.g. graphs with millions of nodes and edges). Finally, generalizability is an important issue for learning methods. Our Graphlet Count Learning problem assumes graph data comes from the same (or similar) distribution. Later extensions can eventually relax this assumption, and learn a more general model which allows us to handle graphs from totally different distributions.

7 Acknowledgement

The work of John C.S. Lui is supported in part by the grant GRF 14201819.

References

- [1] Nataša Pržulj. Biological network comparison using graphlet degree distribution. *Bioinformatics*, 23(2):e177–e183, 2007.
- [2] Paul W Holland and Samuel Leinhardt. Local structure in social networks. *Sociological methodology*, 7:1–45, 1976.
- [3] Leman Akoglu, Hanghang Tong, and Danai Koutra. Graph based anomaly detection and description: a survey. *Data Mining and Knowledge Discovery*, 29(3):626–688, 2015.
- [4] Nino Shervashidze, SVN Vishwanathan, Tobias Petri, Kurt Mehlhorn, and Karsten Borgwardt. Efficient graphlet kernels for large graph comparison. In *Artificial Intelligence and Statistics*, pages 488–495, 2009.
- [5] Tijana Milenković, Vesna Memišević, Anand K Ganesan, and Nataša Pržulj. Systems-level cancer gene identification from protein interaction network topology applied to melanogenesis-related functional genomics data. *Journal of the Royal Society Interface*, 7(44):423–437, 2010.
- [6] Johan Ugander, Lars Backstrom, and Jon Kleinberg. Subgraph frequencies: Mapping the empirical and extremal geography of large graph collections. In *Proceedings of the 22nd international conference on World Wide Web*, pages 1307–1318. ACM, 2013.
- [7] Luca Becchetti, Paolo Boldi, Carlos Castillo, and Aristides Gionis. Efficient semi-streaming algorithms for local triangle counting in massive graphs. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 16–24. ACM, 2008.
- [8] Nesreen K Ahmed, Jennifer Neville, Ryan A Rossi, and Nick Duffield. Efficient graphlet counting for large networks. In *Data Mining (ICDM), 2015 IEEE International Conference on*, pages 1–10. IEEE, 2015.
- [9] Xiaowei Chen and John Lui. Mining graphlet counts in online social networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 12(4):41, 2018.
- [10] Norishige Chiba and Takao Nishizeki. Arboricity and subgraph listing algorithms. *SIAM Journal on Computing*, 14(1):210–223, 1985.
- [11] Mahmudur Rahman, Mansurul Alam Bhuiyan, and Mohammad Al Hasan. Graft: An efficient graphlet counting method for large graph analysis. *IEEE Transactions on Knowledge and Data Engineering*, 26(10):2466–2478, 2014.

- [12] Comandur Seshadhri, Ali Pinar, and Tamara G Kolda. Triadic measures on graphs: The power of wedge sampling. In *Proceedings of the 2013 SIAM International Conference on Data Mining*, pages 10–18. SIAM, 2013.
- [13] Madhav Jha, C Seshadhri, and Ali Pinar. Path sampling: A fast and provable method for estimating 4-vertex subgraph counts. In *Proceedings of the 24th International Conference on World Wide Web*, pages 495–505. International World Wide Web Conferences Steering Committee, 2015.
- [14] Xiaowei Chen, Yongkun Li, Pinghui Wang, and John Lui. A general framework for estimating graphlet statistics via random walk. *Proceedings of the VLDB Endowment*, 10(3):253–264, 2016.
- [15] Marco Bressan, Flavio Chierichetti, Ravi Kumar, Stefano Leucci, and Alessandro Panconesi. Counting graphlets: Space vs time. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 557–566. ACM, 2017.
- [16] Shweta Jain and C Seshadhri. A fast and provable method for estimating clique counts using turán’s theorem. In *Proceedings of the 26th International Conference on World Wide Web*, pages 441–449. International World Wide Web Conferences Steering Committee, 2017.
- [17] Mark EJ Newman. Power laws, pareto distributions and zipf’s law. *Contemporary physics*, 46(5):323–351, 2005.
- [18] Olaf Sporns, Dante R Chialvo, Marcus Kaiser, and Claus C Hilgetag. Organization, development and function of complex brain networks. *Trends in cognitive sciences*, 8(9):418–425, 2004.
- [19] Ryan A. Rossi and Nesreen K. Ahmed. Complex networks are structurally distinguishable by domain. *Social Network Analysis and Mining*, 9(1), 2019.
- [20] Steven Kay Butler. *Eigenvalues and structures of graphs*. PhD thesis, UC San Diego, 2008.
- [21] Charalampos E Tsourakakis. Fast counting of triangles in large real networks without counting: Algorithms and laws. In *Data Mining, 2008. ICDM’08. Eighth IEEE International Conference on*, pages 608–617. IEEE, 2008.
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [23] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [24] Rie Johnson and Tong Zhang. Semi-supervised convolutional neural networks for text categorization via region embedding. In *Advances in neural information processing systems*, pages 919–927, 2015.
- [25] Uwe Schöning. Graph isomorphism is in the low hierarchy. *Journal of Computer and System Sciences*, 37(3):312–323, 1988.
- [26] David H Hubel and Torsten N Wiesel. Receptive fields and functional architecture of monkey striate cortex. *The Journal of physiology*, 195(1):215–243, 1968.
- [27] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436, 2015.
- [28] Muhan Zhang and Yixin Chen. Weisfeiler-lehman neural machine for link prediction. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 575–583. ACM, 2017.
- [29] Ryan A Rossi, Rong Zhou, and Nesreen K Ahmed. Estimation of graphlet counts in massive networks. *IEEE transactions on neural networks and learning systems*, (99):1–14, 2018.
- [30] Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [31] Pinghui Wang, John Lui, and Don Towsley. Minfer: Inferring motif statistics from sampled edges. *arXiv preprint arXiv:1502.06671*, 2015.
- [32] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):13–30, 1963.
- [33] Nesreen K Ahmed, Jennifer Neville, Ryan A Rossi, Nick G Duffield, and Theodore L Willke. Graphlet decomposition: Framework, algorithms, and applications. *Knowledge and Information Systems*, 50(3):689–722, 2017.
- [34] S Vichy N Vishwanathan, Nicol N Schraudolph, Risi Kondor, and Karsten M Borgwardt. Graph kernels. *Journal of Machine Learning Research*, 11(Apr):1201–1242, 2010.
- [35] Nikil Wale, Ian A Watson, and George Karypis. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and Information Systems*, 14(3):347–375, 2008.

- [36] Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2, 2007.
- [37] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [38] Mansurul A Bhuiyan, Mahmudur Rahman, and M Al Hasan. Guise: Uniform sampling of graphlets for large graph analysis. In *Data Mining (ICDM), 2012 IEEE 12th International Conference on*, pages 91–100. IEEE, 2012.
- [39] Daniel Mawhirter, Bo Wu, Dinesh Mehta, and Chao Ai. Approxg: fast approximate parallel graphlet counting through accuracy control. In *2018 18th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGRID)*, pages 533–542. IEEE, 2018.
- [40] Pinghui Wang, Junzhou Zhao, Xiangliang Zhang, Zhenguo Li, Jiefeng Cheng, John CS Lui, Don Towsley, Jing Tao, and Xiaohong Guan. Moss-5: A fast method of approximating counts of 5-node graphlets in large graphs. *IEEE Transactions on Knowledge and Data Engineering*, 30(1):73–86, 2017.
- [41] Ali Pinar, C Seshadhri, and Vaidyanathan Vishal. Escape: Efficiently counting all 5-vertex subgraphs. In *Proceedings of the 26th International Conference on World Wide Web*, pages 1431–1440. International World Wide Web Conferences Steering Committee, 2017.
- [42] Noga Alon, Raphael Yuster, and Uri Zwick. Color-coding. *Journal of the ACM (JACM)*, 42(4):844–856, 1995.
- [43] Boris Weisfeiler and AA Lehman. A reduction of a graph to a canonical form and an algebra arising during this reduction. *Nauchno-Technicheskaya Informatsia*, 2(9):12–16, 1968.
- [44] Thomas Schank and Dorothea Wagner. Finding, counting and listing all triangles in large graphs, an experimental study. In *International workshop on experimental and efficient algorithms*, pages 606–609. Springer, 2005.
- [45] Lars Backstrom and Jure Leskovec. Supervised random walks: predicting and recommending links in social networks. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 635–644. ACM, 2011.
- [46] Bowu Zhang, Kai Xing, Xiuzhen Cheng, Liusheng Huang, and Rongfang Bie. Traffic clustering and online traffic prediction in vehicle networks: A social influence perspective. In *Infocom, 2012 Proceedings IEEE*, pages 495–503. IEEE, 2012.
- [47] Zi-Ke Zhang, Tao Zhou, and Yi-Cheng Zhang. Personalized recommendation via integrated diffusion on user-item-tag tripartite graphs. *Physica A: Statistical Mechanics and its Applications*, 389(1):179–186, 2010.
- [48] Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proc. CVPR*, volume 1, page 3, 2017.
- [49] Stefano Boccaletti, Vito Latora, Yamir Moreno, Martin Chavez, and D-U Hwang. Complex networks: Structure and dynamics. *Physics reports*, 424(4-5):175–308, 2006.
- [50] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864. ACM, 2016.
- [51] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 701–710. ACM, 2014.
- [52] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutikov. Learning convolutional neural networks for graphs. In *International conference on machine learning*, pages 2014–2023, 2016.
- [53] Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- [54] Mikael Henaff, Joan Bruna, and Yann LeCun. Deep convolutional networks on graph-structured data. *arXiv preprint arXiv:1506.05163*, 2015.
- [55] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, pages 3844–3852, 2016.
- [56] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [57] Will Hamilton, Zhitaoying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pages 1024–1034, 2017.
- [58] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.