

Saliency-Aware Regularized Graph Neural Network

Wenjie Pei^{a,*}, WeiNa Xu^{b,*}, Zongze Wu^c, Weichao Li^d, Jinfan Wang^b,
Guangming Lu^a, Xiangrong Wang^{c,**}

^a*Department of Computer Science, Harbin Institute of Technology at
Shenzhen, Shenzhen, 518172, China*

^b*Institute of Future Networks, Southern University of Science and
Technology, Shenzhen, 518055, China*

^c*College of Mechatronics and Control Engineering, Shenzhen
University, Shenzhen, 518060, China*

^d*Peng Cheng Laboratory, Shenzhen, 518066, China*

Abstract

The crux of graph classification lies in the effective representation learning for the entire graph. Typical graph neural networks focus on modeling the local dependencies when aggregating features of neighboring nodes, and obtain the representation for the entire graph by aggregating node features. Such methods have two potential limitations: 1) the global node saliency w.r.t. graph classification is not explicitly modeled, which is crucial since different nodes may have different semantic relevance to graph classification; 2) the graph representation directly aggregated from node features may have limited effectiveness to reflect graph-level information. In this work, we propose the Saliency-Aware Regularized Graph Neural Network (*SAR-GNN*) for graph classification, which consists of two core modules: 1) a traditional graph neural network serving as the backbone for learning node features and 2) the Graph Neural Memory designed to distill a compact graph representation from node features of the backbone. We first estimate the global node saliency by measuring the semantic similarity between the compact graph representation and node features. Then the learned saliency distribution is leveraged to regularize the neighborhood aggregation of the backbone, which facilitates the message passing of features for salient nodes and suppresses the less relevant nodes. Thus, our model can learn more effective graph represen-

*Equal contribution.

**Corresponding author.

tation. We demonstrate the merits of *SAR-GNN* by extensive experiments on seven datasets across various types of graph data. Code will be released.

Keywords: Graph Neural Network, Graph Classification.

1. Introduction

Successful graph classification requires effective representation learning not only for each node but also for the entire graph. Existing methods based on graph neural networks [1, 2] have achieved significant progress in node representation learning, benefiting from the excellent feature refinement through neighborhood aggregation for each node. Most of these methods focus on modeling the local dependencies between neighboring nodes when performing neighborhood aggregation. Prominent examples include GCN [3] which aggregates neighborhood features by considering the node degree distribution, GraphSAGE [4] as well as GIN [5] which seek to learn effective aggregating functions, and GAT [6] which employs an attention mechanism to learn aggregating weights. While these methods are able to learn effective node representations in a graph, they ignore the global saliency of each node, formulated as the semantic similarities between features of each node and the graph representation in a latent space. Nevertheless, it is crucial to perceive global node saliency for graph classification since different nodes potentially have different relevance. Noisy or irrelevant nodes may have adverse effects on graph classification.

A straightforward way of modeling such global saliency is to employ attention mechanism to calculate the attention distribution, as SAGPool [7] and ChebyGIN [8] behave. However, such implicit modeling way has limited effectiveness since it estimates the attention scores solely based on local node features and lacks the explicit modeling of semantic similarities between each node and the global graph representation. Other typical methods of modeling global node saliency include Graphormer [9] interpreting node degrees as centrality, DGM [10] utilizing node connectivities to compute a lens function, and DGCNN [11] equating the graph convolution output with continuous Weisfeiler-Lehman (WL) [12] for sorting nodes (termed ‘SortPooling’). Nevertheless, all these methods still suffer from the limitation that they do not explicitly measure the semantic similarities between each node and the global graph representation, which essentially characterizes the saliency.

Another essential concern of graph classification is how to learn the effective representation for the entire graph. A commonly adopted approach by most existing methods based on graph neural networks [13, 14, 11, 15, 16], is to first learn node features and then aggregate all node features to produce the graph representation via a readout function, which is typically implemented by various pooling operations. Despite the simplicity of such approach, the obtained graph representation may have limited effectiveness. This is because most model capacity is allocated for learning node features rather than the graph representation, and the quality of obtained graph representation depends heavily on that of the node features.

To address the above two potential limitations, in this work we propose the Saliency-Aware Regularized Graph Neural Network (*SAR-GNN*) for graph classification which is built upon traditional graph neural networks. It first measures the saliency of each node in a global view, namely the semantic similarities between each node and the graph representation, and then leverages the learned saliency distribution to regularize the graph neural network. Specifically, the *SAR-GNN* consists of two core modules: 1) a traditional graph neural network serving as the backbone for learning node features and 2) Graph Neural Memory which is designed to distill a compact graph representation for the entire graph from node features of the backbone. Both the backbone network and the Graph Neural Memory are iteratively stacked to refine node features and the graph representation progressively. The learned compact graph representation is used to estimate the global saliency of each node by measuring the semantic similarities between the graph representation and node features. Then the obtained saliency distribution is used to regularize the aggregating weights of the backbone network when performing neighborhood aggregation for each node, which facilitates the message passing of features for the salient nodes while suppressing the feature propagation for the less relevant nodes to graph classification. As a result, the Graph Neural Memory can distill more relevant features to graph classification from the salient nodes when learning the compact graph representation. Two modules work interdependently with interactions between each other in each stacked layer. The obtained compact graph representation is finally utilized for graph classification.

To conclude, we highlight following contributions.

- We propose the Saliency-Aware Regularized Graph Neural Network (*SAR-GNN*), a novel framework for graph classification, which consists

of two core modules: a backbone network for learning node features and the Graph Neural Memory for distilling the compact graph representation. Two modules are optimized interdependently with interactions between each other in each stacked layer and thus both the node features and the graph representation can be refined iteratively. Such framework enables the proposed *SAR-GNN* to model the global node saliencies in an explicit manner by measuring the semantic similarities between each node and the graph representation.

- We design an effective saliency-aware regularization mechanism, which utilizes the learned global node saliencies to regularize the feature aggregation for each node in such a way that it facilitates the message passing of features for the salient nodes while suppressing the feature propagation for the less relevant nodes to graph classification. As a result, the Graph Neural Memory can distill more relevant features to graph classification from the salient nodes and thereby produce more effective graph representation for graph classification.
- The proposed *SAR-GNN* can be readily applied to most of the existing graph neural networks which perform neighborhood aggregation for refining node features. In particular, we instantiate the backbone of our model with four classical types of graph neural networks, namely GCN, GraphSAGE, GIN and GAT respectively. Then we conduct extensive experiments both quantitatively and qualitatively on seven challenging datasets across various types of graph data, which demonstrate 1) the effectiveness of our method via thorough ablation studies and complexity analysis, and 2) the favorable performance of our model compared to the state-of-the-art methods for graph classification.

2. Related Work

Traditional approaches to graph classification, preceding the popularity of graph neural networks, typically employ graph kernel functions to measure the similarity between pairs of graphs [17, 18, 19]. Our *SAR-GNN* is built upon graph neural networks and is designed to: 1) regularize the neighborhood aggregation of graph neural networks and 2) learn effective representation for the entire graph. Thus, we discuss related work to our method from these two perspectives below.

Neighborhood aggregation of Graph Neural Networks. The core idea shared across various graph neural networks is to perform neighborhood aggregation for each graph node to refine the node features iteratively [20, 21]. Typical convolutional neural networks [22] are successfully adapted to graph data either in the spectral domain [23, 24] or spatial domain [2, 25]. In particular, GCN [3] presents a scalable and efficient implementation of graph convolution. GraphSAGE [4] explores multiple potential aggregating functions while Graph Isomorphism Network (GIN) [5] builds upon GraphSAGE and presents a framework which is as theoretically powerful as the Weisfeiler-Lehman graph isomorphism test [12]. Another research line of modeling neighborhood aggregation is to learn the aggregating weights directly per edge around the center node. Edge-Conditioned Convolution (ECC) [13] learns aggregating weights conditioned on edge labels while GAT [6] employs an attention mechanism to measure compatibility between neighboring nodes for each edge.

All aforementioned methods focus on modeling the local dependencies between neighboring nodes when performing feature aggregation whilst the relevance of each node to the task of graph classification is ignored. While the node saliency can be modeled by the attention mechanism roughly, as ChebyGIN [8] and SAGPool [7] do, this type of methods has limited effectiveness in that it solely relies on the attention model without explicit modeling of saliency. Besides, the global node saliency has been also modeled based on 1) node degrees by Graphormer [9], or 2) the connectivities between nodes by DGM [10], or 3) the graph convolution output by DGCNN [11]. Nevertheless, all these methods do not explicitly model the saliency yet, namely the semantic similarities between each node and the global graph representation. In contrast, our *SAR-GNN* measures the node saliency by modeling the compatibility between the learned graph representation and the node features. The learned node saliency is then leveraged to regularize the modeling of neighborhood aggregation.

Representation learning for the entire graph. Most of the existing methods for graph classification learn the representation for the entire graph in an indirect manner: they allocate most of the model capacity to learning effective node features, and obtain the representation for the entire graph by simply aggregating node features via a readout function, mostly implemented as a variety of pooling functions. Typical examples include DiffPool [14] which designs an adaptive pooling method that collapses nodes hierarchically, ECC which pre-calculates a pooling map to coarsen graphs,

DGCNN [11] that sorts nodes first using SortPool before pooling, Eigen-Pooling [15] which is proposed based on graph Fourier transform to preserve graph structure during pooling, Graph U-Nets [16] which performs top-K pooling scheme, and SAGPool [7] which identifies the important nodes using the self-attention mechanism [26]. On the other hand, ESAN [27] represents a large graph as a set of distinguishable subgraphs based on the predefined policy. Unlike aforementioned prior work, we design a module termed as Graph Neural Memory specifically to distill a compact graph representation from node features progressively using a cross-attention mechanism. In particular, the compact graph representation and the node features are refined by two modules interdependently.

Although our model is mainly designed for effective graph representation learning, another potentially advantageous application of it, which is worth noting, is that the derived node saliency distribution by our *SAR-GNN* can be naturally used for interpretable explanation of GNN predictions. In this sense, our model is essentially consistent with GNNExplainer [28] which identifies compact subgraphs crucial for GNN’s predictions for interpretation, as well as PGExplainer [29] that is a generalizable GNN explainer by parameterizing the generation process of explanations.

3. Saliency-Aware Regularized Graph Neural Network

Built upon traditional graph neural networks, our *SAR-GNN* first captures the global node saliency w.r.t. the task of graph classification. The learned node saliency is then leveraged to regularize the neighborhood aggregation of graph neural networks and facilitate the message passing of features for salient nodes, thus our model can learn more effective representations for the entire graph.

3.1. Overview

Figure 1 illustrates the architecture of our *SAR-GNN*, which consists of two core modules: 1) a traditional graph neural network serving as the backbone for learning features of each node in a graph and 2) Graph Neural Memory which is specifically designed for distilling a compact graph representation for the entire graph from the node features of the backbone. Two modules work interdependently to refine both the compact graph representation and node features progressively through the stacked layers.

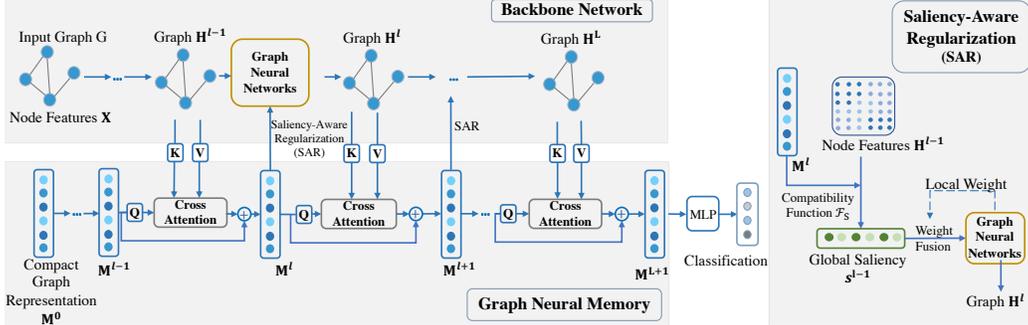


Figure 1: Architecture of the Saliency-Aware Regularized Graph Neural Network (SAR-GNN). It consists of two core modules: 1) a traditional graph neural network serving as the backbone network for learning node features and 2) the Graph Neural Memory for distilling a compact graph representation from node features of the backbone. The learned graph representation is leveraged to measure the global node saliency and regularize the backbone. Thus, two modules work interdependently to refine node features and the compact graph representation in an iterative manner. The operation of Saliency-Aware Regularization in the left panel is elaborated in the right panel.

Taken a graph G with N nodes as input, *SAR-GNN* first employs the Graph Neural Memory to perform feature distillation from the node features in the backbone network and obtains a compact graph representation for the whole graph:

$$\mathbf{M}^l = \mathcal{F}_M(\mathbf{H}^{l-1}, \mathbf{M}^{l-1}). \quad (1)$$

Here $\mathbf{M}^l \in \mathbb{R}^{d_M}$ denotes the compact graph representation with d_M dimensions at the l -th stacked layer, which is refined iteratively from the previous representation \mathbf{M}^{l-1} by incorporating the newly distilled information from the node feature matrix $\mathbf{H}^{l-1} \in \mathbb{R}^{N \times d_H}$ at the $(l-1)$ -th layer. Each of N nodes has d_H dimensions of features. \mathcal{F}_M is the transformation function by the Graph Neural Memory that will be explained concretely in Section 3.2. The obtained compact graph representation \mathbf{M}^l is further leveraged to measure the global saliency of each node for graph classification. In particular, we learn global saliency by modeling the compatibility between \mathbf{M}^l and node features via a function \mathcal{F}_S (explicated in Section 3.3):

$$\mathbf{s}^{l-1} = \mathcal{F}_S(\mathbf{M}^l, \mathbf{H}^{l-1}). \quad (2)$$

Herein, $\mathbf{s}^{l-1} \in \mathbb{R}^N$ is a vector, which sums to 1, denoting the saliency distribution for all nodes in the $(l-1)$ -th layer of the backbone network. It can be interpreted as the semantic similarities between each node and the graph

representation. The saliency \mathbf{s}^{l-1} is then used to regularize the backbone network and thereby refines the node features in the l -th layer:

$$\mathbf{H}^l = \mathcal{F}_B(\mathbf{H}^{l-1}, \mathbf{s}^{l-1}), \quad (3)$$

where \mathcal{F}_B is the transformation function of the backbone network. The refined node features \mathbf{H}^l are in turn used to distill the compact graph representation in the next layer \mathbf{M}^{l+1} as shown in Equation 1. Thus, the compact graph representation and the node features are refined interdependently. The compact graph representation in the last layer of the Graph Neural Memory \mathbf{M}^{L+1} (L is the number of feature-refining layers in the backbone network) is finally utilized for graph classification.

We will first elaborate on the design of the Graph Neural Memory, and then describe how to perform global regularization for the backbone network using the compact graph representation obtained from the Graph Neural Memory. Finally, we show that our *SAR-GNN* can be optimized in an end-to-end manner as a whole.

3.2. Graph Neural Memory

The Graph Neural Memory is designed to distill a compact latent representation for the entire graph from node features in the backbone network. As shown in Figure 1, the Graph Neural Memory learns a vectorial latent embedding $\mathbf{M} \in \mathbb{R}^{d_M}$ as the compact graph representation and refines it iteratively using cross-attention operation. Specifically, the initial latent embedding in the 0-th layer (denoted \mathbf{M}^0) is parameterized as a learnable vector and initialized randomly. In each layer of the Graph Neural Memory, the latent embedding is refined by attending to the node features of the previous layer in the backbone to perform feature distillation with cross-attention operation. Formally, the transformation function of the Graph Neural Memory \mathcal{F}_M in Equation 1 is modeled as:

$$\begin{aligned} \mathbf{q}_m &= \mathbf{M}^{l-1} \mathbf{W}^q, \quad \mathbf{K} = \mathbf{H}^{l-1} \mathbf{W}^k, \quad \mathbf{V} = \mathbf{H}^{l-1} \mathbf{W}^v, \\ \mathbf{M}^l &= \mathcal{F}_{\text{MLP}}(\mathbf{M}'), \quad \mathbf{M}' = \text{softmax}\left(\frac{\mathbf{q}_m \mathbf{K}^\top}{\sqrt{d}}\right) \mathbf{V}, \end{aligned} \quad (4)$$

where the latent embedding \mathbf{M}^{l-1} in the $l-1$ -th layer of the Graph Neural Memory serves as the query while the node feature matrix \mathbf{H}^{l-1} in the $(l-1)$ -th layer of the backbone network is used as both the key and value in

cross-attention operation. $\mathbf{W}^q \in \mathbb{R}^{d_M \times d}$, $\mathbf{W}^K \in \mathbb{R}^{d_H \times d}$, $\mathbf{W}^V \in \mathbb{R}^{d_H \times d}$ are learnable parameter matrices for linear transformations. \mathcal{F}_{MLP} is a Multilayer Perceptron (MLP) module consisting of two fully connected layers (with transformation matrices $\mathbf{W}_{\text{MLP}}^1 \in \mathbb{R}^{d_M \times d}$ and $\mathbf{W}_{\text{MLP}}^2 \in \mathbb{R}^{d \times d}$) with ReLU function. In our implementation, the transformation \mathcal{F}_M described in Equation 4 is typically performed k iterations (tuned as a hyper-parameter) in a single layer of the Graph Neural Memory. Note that residual connections are applied between adjacent layers of the Graph Neural Memory.

The learned latent embedding \mathbf{M}^l is in turn used to regularize the backbone network for refining the node feature matrix \mathbf{H}^l (Section 3.3). Thus the Graph Neural Memory is able to distill graph representations from newly refined node features in each layer instead of constant node information. Similar way of information distillation via cross attention has been previously adopted in Transformer-based models like BERT [30] and ViT [31] that learn a ‘CLS’ token for classification, DETR for object detection [32] and Perceiver for multi-modal feature learning [33]. Note that the initial latent embedding \mathbf{M}^0 is modeled as learnable parameters to learn an appropriate initial feature point that is compatible with the latent feature space in the backbone.

3.3. Saliency-Aware Regularization of the Backbone

The distilled compact graph representation from the Graph Neural Memory is used to perform global regularization on the backbone network, which enables the backbone network to perceive the global node saliency w.r.t. graph classification when learning node features. Specifically, we first utilize the compact graph representation to measure the global saliency for each node, and then regularize the backbone network using the node saliency distribution.

Measuring the global node saliency. The node saliency w.r.t. graph classification is measured by modeling the compatibility between the compact graph representation and features of each node in the graph. A straightforward way is to calculate the similarity between them by dot product in a projected latent space. The obtained similarity scores for all nodes are then normalized. Thus, the compatibility function \mathcal{F}_S in Equation 2 is formulated as:

$$\begin{aligned} \mathbf{q}_s &= \mathbf{M}^l \mathbf{W}_s^q, \quad \mathbf{K}_s = \mathbf{H}^{l-1} \mathbf{W}_s^K, \\ \mathbf{s}^{l-1} &= \text{softmax}\left(\frac{\mathbf{q}_s \mathbf{K}_s^\top}{\sqrt{d_s}}\right), \end{aligned} \tag{5}$$

where $\mathbf{W}_s^q \in \mathbb{R}^{d_M \times d}$ and $\mathbf{W}_s^K \in \mathbb{R}^{d_H \times d}$ are learnable parameter matrices to project the compact graph representation \mathbf{M}^l and the node feature matrix \mathbf{H}^{l-1} into the same latent space. Note that the scaling factor $\frac{1}{\sqrt{d}}$ is used to avoid the explosive growth of dot product between two vectors.

Regularizing the backbone network with node saliency. Typical graph neural networks refine features of a node by performing weighted aggregation over (transformed) features of its neighboring nodes, where the weights are either derived based on the node degrees like GCN or basic uniform distributions such as GraphSAGE or GIN. These methods solely model the local dependencies between neighboring nodes and fail to incorporate the global node saliency w.r.t. graph classification. Intuitively, the nodes with higher saliency score should receive more attention than those with lower saliency score during information propagation between adjacent nodes. Hence, we regularize the backbone network with the learned node saliency distribution.

We formulate the aggregating weights for the neighborhood of the i -th node (including itself) in the l -th layer as local weight distribution $\mathbf{a}_i^l \in \mathbb{R}^{|\mathcal{N}_i|+1}$, where \mathcal{N}_i denotes the set of neighboring nodes of the i -th node. The learned global saliency \mathbf{s}^l is viewed as the global weight distribution. Then we regularize the feature aggregation for each node in the backbone network by fusing the local and global weights together. We discuss two fusion mechanisms:

- **Weighted sum** over the local and global weights:

$$\begin{aligned} \mathbf{w}_i^l(j) &= \text{softmax}(\mathbf{a}_i^l(j) + \beta \mathbf{s}^l(j)), \\ \forall j \in \{i\} \cup \mathcal{N}_i, 1 \leq i \leq N, \end{aligned} \quad (6)$$

where $\mathbf{w}_i^l \in \mathbb{R}^{|\mathcal{N}_i|+1}$ is the regularized aggregating weights for the i -th node and its neighboring nodes, and j is a node index. β is a hyper-parameter tuned on a validation set. Softmax function is applied for normalization.

- **Scaling regularization** over the local weights by the global weights:

$$\begin{aligned} \mathbf{w}_i^l(j) &= \text{softmax}((1 + \mathbf{s}^l(j))^\gamma \mathbf{a}_i^l(j)), \\ \forall j \in \{i\} \cup \mathcal{N}_i, 1 \leq i \leq N, \end{aligned} \quad (7)$$

wherein, $\gamma > 0$ is a hyper-parameter.

The proposed saliency-aware regularization mechanism is readily applicable to most graph neural networks. In our implementation, the backbone network is instantiated with four classical types of graph neural networks: GCN, GraphSAGE, GIN and GAT. The transformation function \mathcal{F}_B in Equation 3 is formulated correspondingly as follows.

- **GCN**, which takes into account the degree distribution when performing neighborhood aggregation:

$$\begin{aligned} \mathbf{a}^l &= \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}, \\ \mathbf{H}^{l+1} &= \sigma(\mathbf{w}^l \mathbf{H}^l \Theta^l). \end{aligned} \quad (8)$$

Herein, we consider $\mathbf{a}^l \in \mathbb{R}^{N \times N}$ as the local aggregating weight matrix in the renormalized form [3] derived from the adjacency matrix $\tilde{\mathbf{A}}$ and degree matrix $\tilde{\mathbf{D}}$. $\mathbf{w}^l \in \mathbb{R}^{N \times N}$ is the regularized aggregating weight matrix for all nodes, calculated by Equation 6 or 7. Θ^l is the learnable parameter matrix for linear transformation in l -th layer and σ is the activation function.

- **GraphSAGE**. Taking the mean aggregator of GraphSAGE as an example, the local aggregating weight for all neighboring nodes is equal to 1. Thus, the features for i -th node $\mathbf{H}_i^{(l+1)}$ is refined as:

$$\mathbf{H}_i^{(l+1)} = \sigma\left(\Theta^l \text{MEAN}(\{\mathbf{w}_i^l(i) \mathbf{H}_i^l\} \cup \{\mathbf{w}_i^l(j) \mathbf{H}_j^l, \forall j \in \mathcal{N}_i\})\right). \quad (9)$$

Here \mathbf{w}_i^l is the regularized weights for the i -th node and its neighboring nodes by Equation 6 or 7.

- **GIN**, whose local aggregating weights are equal to 1 except for the center node:

$$\begin{aligned} \mathbf{a}^l(i) &= 1 + \epsilon^{l+1}, \quad \mathbf{a}^l(j) = 1, \forall j \in \mathcal{N}_i, \\ \mathbf{H}_i^{l+1} &= \mathcal{F}_{\text{MLP}}^l\left(\mathbf{w}_i^l(i) \mathbf{H}_i^l + \sum_{j \in \mathcal{N}_i} \mathbf{w}_i^l(j) \mathbf{H}_j^l\right), \end{aligned} \quad (10)$$

where $\mathcal{F}_{\text{MLP}}^l$ denotes the transformation function for a Multilayer Perceptron at the l -th layer of backbone network, and \mathbf{w}_i^l is the regularized weights for the i -th node and its neighboring nodes by Equation 6 or 7.

- **GAT**, which uses an attention mechanism to calculate the local weight distribution between each node and its neighboring nodes:

$$\begin{aligned} \mathbf{a}_i^l(j) &= \frac{\exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]))}{\sum_{j \in \{i\} \cup \mathcal{N}(i)} \exp(\text{LeakyReLU}(\mathbf{a}^\top [\mathbf{W}\mathbf{h}_i \parallel \mathbf{W}\mathbf{h}_j]))}, \\ \mathbf{h}_i^{l+1} &= \sigma\left(\sum_{j \in \{i\} \cup \mathcal{N}(i)} \mathbf{w}_i^l(j) \mathbf{W}\mathbf{h}_j\right). \end{aligned} \quad (11)$$

where \mathbf{h}_i are the feature vector for the i -th node and \mathbf{W} is a shared transformation matrix. \mathbf{a} is the parameters for a linear transformation to calculate a scalar matching score while \parallel denotes the concatenation operation. $\mathbf{w}_i^l(j)$ is the regularized weights between the i -th node and the j -th node.

3.4. End-to-End Parameter Learning

Since the compact graph representation in the 1-th layer of the Graph Neural Memory \mathbf{M}^1 is learned from the input node features \mathbf{X} , the Graph Neural Memory contains one more layer than the backbone network. Suppose the backbone network contains in total L feature-refining layers, the learned compact graph representation in the last layer of the Memory \mathbf{M}^{L+1} is utilized for predicting the graph label:

$$P(\mathbf{y}|\mathbf{M}^{L+1}) = \text{softmax}(\mathcal{F}_{\text{MLP}}(\mathbf{M}^{L+1})), \quad (12)$$

where $\mathbf{y} \in \mathbb{R}^K$ is predicted probabilities for all potential K categories. \mathcal{F}_{MLP} refers to the transformation function for a Multilayer Perceptron.

The proposed *SAR-GNN* can be optimized in an end-to-end manner by a cross-entropy loss:

$$\mathcal{L} = - \sum_{i=1}^{N_t} \log P(\mathbf{y}^i|\mathbf{X}^i). \quad (13)$$

Here \mathbf{X}^i and \mathbf{y}^i refer to the i -th training sample with its graph label while N_t is the training size.

4. Experiments

To evaluate our proposed *SAR-GNN*, we first conduct ablation study to investigate the effectiveness of each essential technique in our model, then

we conduct extensive experiments to compare our model with other state-of-the-art methods for graph classification on seven challenging graph datasets across various types of graph data. Moreover, we also perform qualitative evaluation to validate the effectiveness of the learned node saliency and the learned graph representation.

4.1. Experimental Setup

Datasets. We conduct experiments on seven datasets across various types of graph data for evaluation. 1) *MUTAG* [34] is a chemical dataset containing two categories of compounds: mutagenic aromatic and heteroaromatic compounds. 2) *ENZYMES* [35] is a protein dataset comprising 600 protein tertiary structures, which are grouped into 6 top-level enzyme classes. 3) *PROTEINS* [36] is also a protein dataset which contains two types of samples: enzymes or non-enzymes. 4-5) *IMDB-MULTI* [37] and *IMDB-BINARY* [37] both include ego-network graphs involving 1000 actors who play roles in movies appeared in IMDB. All graphs in *IMDB-MULTI* are categories into three genre classes including Comedy, Romance and Sci-Fi while *IMDB-BINARY* has two categories: Action and Romance. 6) *TRIANGLES* [8] is a large synthetic dataset consisting of 45,000 graph samples, and the task is to predict the number of triangles in a graph (ranging from 1 to 10). Note that we perform evaluation on the split of the test set, which has the same number of categories of graphs as the training set. 7) *Letter-high* [38] is also a synthetic dataset containing 2,250 graph samples and each sample represents a distorted letter drawings.

Implementation details. Following typical evaluation protocol for graph classification [21], we perform 10-fold cross validation for all datasets except for *TRIANGLES* dataset which has the official data division for training and test. For quantitative evaluation, we measure the accuracy of graph classification as the evaluation metric. Adam [39] is employed for gradient descent optimization.

4.2. Ablation Study

To investigate the effectiveness of each proposed component, we first perform ablation study with six variants of our *SAR-GNN* for all four types of backbones.

- **Base model**, which is exactly the backbone network, e.g., base model is GCN when it is used as the backbone. Following Errica et al. [21], we apply sum pooling to obtain graph representations for all *Base* models.

Table 1: Classification accuracy (%) of six variants of our model on *MUTAG* and *Letter-high* for ablation study. The backbone is instantiated with GCN, GIN, GraphSAGE and GAT, respectively. The standard deviation on 10-fold cross validation is also provided.

Variants \ Backbone	GCN		GIN		GraphSAGE		GAT	
	MUTAG	Letter-high	MUTAG	Letter-high	MUTAG	Letter-high	MUTAG	Letter-high
Base model	71.6 ±10.9	61.1±2.6	81.4±6.6	73.3±2.5	75.8±7.8	71.0±2.7	76.1 ±4.2	82.1±1.8
GNN-GNN	76.7±7.0	73.2±1.9	84.7±5.6	79.5±2.2	78.4±7.3	75.9±2.3	77.5 ±4.3	82.6±1.3
SAR-Pooling-W	76.2±6.0	70.7±1.8	83.8±4.0	75.3±1.8	82.8±6.8	76.8±2.5	79.9 ±2.7	83.4±4.6
SAR-Pooling-S	75.9±10.1	70.5±2.9	84.0±7.3	74.1±3.4	81.9±10.3	75.8±2.9	78.3 ±3.2	84.4±4.1
SAR-GNN-W	80.0±5.1	77.6±2.4	85.1±5.0	82.3±1.8	84.6±4.2	80.8±2.6	81.4 ±2.8	85.1±2.6
SAR-GNN-S	82.2±5.5	76.8±1.7	85.3±5.5	83.6±1.7	87.4±3.1	80.5±2.6	80.2 ±3.4	85.0±5.2

- **GNN-GNN**, which employs the Graph Neural Memory (GNN) to learn the compact graph representation for graph classification. Nevertheless, the backbone is not regularized with the node saliency.
- **SAR-Pooling**, which performs sum pooling over all node features in the last layer of the backbone for graph classification, instead of using the compact graph representation. Note that the backbone of this variant is regularized with node saliency. In particular, two proposed mechanisms for weight fusion (Equations 6 and 7), namely ‘weighted sum’ and ‘scaling regularization’, are evaluated separately. We denote them as **SAR-Pooling-W** and **SAR-Pooling-S** respectively.
- **SAR-GNN**, which is the intact version of our model. Both mechanisms for weight fusion are evaluated, denoted as **SAR-GNN-W** and **SAR-GNN-S** respectively.

Table 1 shows the experimental results of these variants on *MUTAG* and *Letter-high* datasets, from which we conduct following investigation.

Effect of regularization with learned node saliency. The large performance gap between *Base* model and *SAR-Pooling-W* as well as *SAR-Pooling-S* for all four backbones on both datasets reveals the effect of the global regularization with learned node saliency. Besides, the comparisons between *GNN-GNN* and *SAR-GNN* also demonstrate the advantages of such saliency-aware regularization.

Effect of the Graph Neural Memory. Compared with *Base* model, *GNM-GNN* improves the performance by a large margin for all types of backbones on both datasets. Additionally, *SAR-GNN* consistently outperforms *SAR-Pooling* using the same fusion mechanism in all cases (with different backbones and datasets). These results validate the superiority of the Graph Neural Memory over the pooling method in learning graph representation.

Comparison between ‘Weighted sum’ and ‘Scaling regularization’ for weight fusion. We further compare two different mechanisms for weight fusion in Equations 6 and 7. Comparing the performance between *SAR-GNN-W* and *SAR-GNN-S*, or between *SAR-Pooling-W* and *SAR-Pooling-S* in all cases, we observe that there is no clear winner between these two fusion mechanisms. In most cases, there is no big performance gap between them.

Investigation into the optimization policy. our *SAR-GNN* consists of two core modules: the backbone network for learning node features and Graph Neural Memory for distilling the compact graph representation. Two modules work interdependently to refine each other and the whole model can be optimized in an end-to-end manner. We investigate two optional optimization policies: 1) joint training of two modules and 2) alternating training of two modules.

- **Joint training**, which optimizes two modules jointly by gradient descent method. Adam is employed in our implementation.
- **Alternating training**, in which two modules are optimized in an alternating manner, i.e., the parameters of one module are frozen when optimizing the other one. Such policy is essentially analogous to the block coordinate descent method when considering the two modules as two blocks of parameters to be optimized. The only difference is that we still employ gradient descent method to optimize each block of parameters instead of calculating the local optimal solution per block directly. Similar to the block coordinate descent method, we optimize two modules by conducting such alternating training iteratively to reach the learning convergence.

Figure 2 presents the experimental comparison between two optimization policies. The alternating training policy reaches the optimization convergence after several iterations of alternating optimization. We observe that the joint training policy achieves better performance for graph classification than the alternating training policy. Similar to the coordinate descent

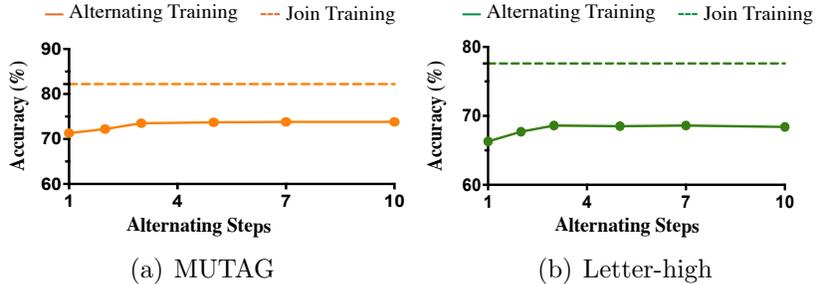


Figure 2: The performance of our *SAR-GNN* adopting two optimization policies, namely the joint training policy and the alternating training policy, on MUTAG (a) and Letter-high datasets (b), respectively.

method, the alternating training policy can theoretically achieve globally optimum performance if it can obtain the local optimum w.r.t. the optimized module in each iteration of per-module optimization for convex optimization problem. However, in our case both alternating training and joint training policies employ gradient descent for non-convex optimization of deep neural networks. Thus, the alternating training policy can hardly reach the optimum w.r.t. the optimized module in each alternating optimization step in practice. On the other hand, the joint training policy performs gradient back-propagation w.r.t. all parameters of the model whilst alternating training only propagates gradients w.r.t. the parameters of the optimized module of the model for optimization. Thus, we surmise that joint training policy is able to optimize the model along more optimized directions than alternating training, thereby yielding better performance.

4.3. Comparison with Other Methods

Next we compare our model with the state-of-the-art methods for graph classification. Apart from four backbone models (GCN, GIN, GraphSAGE and GAT), we also compare our model with 12 state-of-the-art methods for graph classification, including DGCNN [11], ECC [13], ChebyGIN [8], SAG-Pool [7], Graphormer [9], ESAN [27] (termed ‘DSS-GNN’ in the paper) based on the overall best policy ‘EGO+’, GSN [40], Graph U-Nets [16], DGM [10], UGformer [41], SLIM [42] and ADSF-RWR [43]. ADSF-RWR is initially designed for node classification and we adapt it for graph classification by applying sum pooling over all node embeddings to obtain the graph representation. We re-train and re-evaluate all methods under the same experimen-

Table 2: Classification accuracy (%) of different methods for graph classification on seven datasets. The standard deviations on 10-fold cross validation (different initializations for *TRIANGLES*) are provided. ‘SAR-GCN’ refers to our method using GCN as the backbone. The performance gains of our method over each of four *base* models are highlighted in *italic* in parentheses.

Method	MUTAG	ENZYMES	IMDB-MULTI	TRIANGLES	Letter-high	PROTEINS	IMDB-BINARY
DGCNN [11]	83.9±5.8	38.6±4.8	46.4±3.5	77.3±4.2	44.7±2.9	72.5±4.1	66.9±3.9
ECC [13]	81.4±7.6	29.5±8.2	43.5±3.1	80.4±0.5	80.8±1.5	72.3±4.0	66.2±2.1
ChebyGIN-unsup [8]	83.4±6.0	26.3±2.7	39.0±1.5	67.0±3.0	27.0±2.2	71.2±3.4	50.0±0.5
DSS-GNN (EGO+) [27]	84.1±5.9	62.5±5.7	50.1±2.4	93.3±1.1	83.6±2.0	69.2±5.6	71.5±3.6
GSN [40]	84.4±5.5	63.2±6.8	51.1±2.7	76.1±5.2	40.3±2.4	72.9±3.7	72.5±4.1
Graphormer [9]	83.1±6.7	58.7±5.2	48.9±2.7	56.1±2.1	76.3±2.4	76.3±3.4	72.9±3.3
SAGPOOL [7]	72.6±5.7	31.9±5.7	46.0±3.8	63.2±10.3	72.2±3.3	72.6±3.3	69.2±3.5
Graph U-Nets [16]	76.0±7.2	34.4±3.0	47.6±3.0	54.0±1.0	42.5±1.3	74.2±5.2	68.3±3.3
DGM [10]	86.2±6.7	40.3±2.2	49.2±2.5	53.4±2.3	56.4±3.4	68.8±6.9	73.1±3.8
UGformer [41]	78.2± 5.9	68.9±4.7	49.6± 2.1	81.0±4.6	80.6±2.7	70.1±3.5	70.1±2.3
SLIM [42]	78.8± 4.8	56.1±3.4	47.6± 2.1	68.9±3.1	71.3±1.8	74.1±4.7	70.8±5.7
ADSF-RWR [43]	76.4± 6.2	66.0±4.3	45.2± 4.9	62.4±5.2	80.1±2.4	73.9±6.8	69.4±2.7
GCN	71.6±10.9	68.1±4.0	48.2±3.8	46.0±1.0	61.1±2.6	73.3±3.1	68.4±6.7
GIN	81.4±6.6	62.3±5.0	48.5±3.3	51.6±2.7	73.3±2.5	73.5±3.4	70.2±2.8
GraphSAGE	75.8±7.8	62.1±6.2	47.6±3.5	52.5±0.4	71.0±2.7	73.0±4.6	70.7±4.4
GAT	76.1±4.2	65.4 ± 5.2	46.0±3.2	65.4±3.4	82.1±1.8	73.5±2.6	70.0±5.1
SAR-GCN (Ours)	82.2±5.5 (<i>+10.6</i>)	69.3±5.0 (<i>+1.2</i>)	49.1±3.1 (<i>+0.9</i>)	81.7±0.9 (<i>+35.7</i>)	77.6±2.4 (<i>+16.5</i>)	76.8±3.4 (<i>+3.5</i>)	69.2±4.8 (<i>+0.8</i>)
SAR-GIN (Ours)	85.3±5.5 (<i>+3.9</i>)	65.7±3.5 (<i>+3.4</i>)	51.6±3.7 (<i>+3.1</i>)	78.2±4.2 (<i>+26.6</i>)	84.8±1.3 (<i>+11.5</i>)	75.3±6.4 (<i>+1.8</i>)	73.7±3.9 (<i>+3.5</i>)
SAR-GraphSAGE (Ours)	87.4±3.1 (<i>+11.6</i>)	68.7±4.0 (<i>+6.6</i>)	48.7±3.2 (<i>+1.1</i>)	78.9±1.2 (<i>+26.4</i>)	80.5±2.6 (<i>+9.5</i>)	76.0±4.7 (<i>+3.0</i>)	70.9±4.1 (<i>+0.2</i>)
SAR-GAT (Ours)	81.4 ±2.8 (<i>+5.3</i>)	69.1±4.7 (<i>+3.7</i>)	47.7±2.6 (<i>+1.7</i>)	87.4±4.6 (<i>+22.0</i>)	85.1±2.6 (<i>+3.0</i>)	75.8 ±1.9 (<i>+2.3</i>)	72.5±2.8 (<i>+2.5</i>)

tal setting, using the same data split and evaluation metrics, to have a fair comparison. Note that we compare with ChebyGIN in both unsupervised (denoted as ChebyGIN-unsup) and supervised training settings for learning node saliencies to have a comprehensive comparison. In this set of experiments, we only report the performance of our model by selecting the better one between two fusion mechanisms based on the validation set.

Table 2 presents the experimental results of different methods on seven datasets. We make the following observations. First, our method achieves substantial performance gains compared to the corresponding base models used for backbone on all datasets except *IMDB-BINARY*, which validates the effectiveness and adaptability of our method across different types of backbones. The performance improvement of our method on *IMDB-BINARY* is minor presumably due to the ceiling effect, considering the fact that the performance of multiple models on this dataset is close to each other. Second, our method based on one of four backbones achieves the best performance on all seven datasets except *TRIANGLES*, which manifests the robustness of our method. In contrast, other methods, such as GSN, Graphormer and

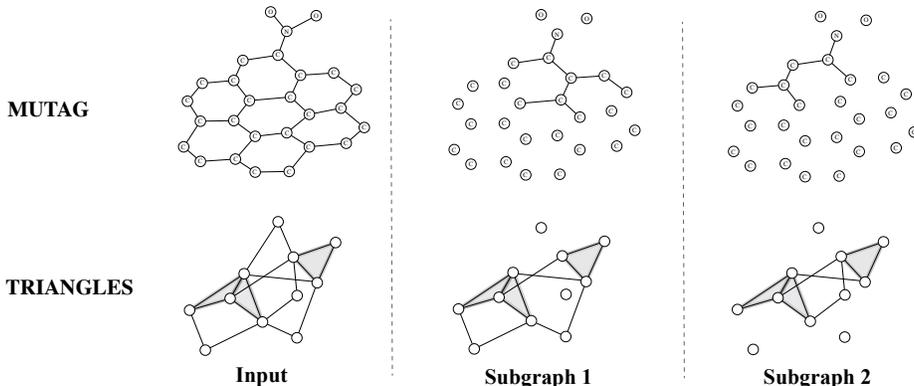


Figure 3: Visualization of two randomly selected subgraphs generated by DSS-GNN (EGO+) on samples from MUTAG and TRIANGLES datasets, respectively. The triangles are indicated in shaded regions for the TRIANGLES sample. While DSS-GNN can preserve the triangle structure in the generated subgraphs, its generated subgraphs on the MUTAG sample fail to preserve intact fused rings whose number is closely correlated to the recognition of mutagenic aromatic.

ECC, show unstable performance across different datasets. DSS-GNN performs substantially better than other methods on *TRIANGLES*, presumably because DSS-GNN can recognize the triangle structure as distinguishable subgraphs. To validate it, we visualize two randomly selected subgraphs learned by DSS-GNN in Figure 3, which shows that DSS-GNN can indeed generate subgraphs that preserve the triangles on the samples from *TRIANGLES* dataset. However, its generated subgraphs on the *MUTAG* sample fail to preserve intact fused rings whose number is closely correlated to the recognition of mutagenic aromatic [34], which accounts for its relatively inferior performance to our model on *MUTAG*. The third observation we make is that our method consistently performs well when built upon different backbones and the performance differences between them are minor. Interestingly, the performance differences among four base models are substantially reduced when they are integrated into our method.

Comparison with the implicit modeling of global saliency by attention mechanism. SAGPool and ChebyGIN use attention mechanism to model global saliency implicitly by learning a parameterized attention function \mathcal{F}_{att} taking the node features \mathbf{H} as input: $\mathbf{s}_{\text{att}} = \mathcal{F}_{\text{att}}(\mathbf{H})$. In contrast, our method models the global saliency explicitly by measuring the compatibility between the compact graph representation \mathbf{M} and node features \mathbf{H} :

Table 3: Classification accuracy (%) of our methods instantiated with different backbones as well as ChebyGIN in both unsupervised and supervised learning settings on TRIANGLES dataset.

Learning settings	ChebyGIN	Our method			
		SAR-GCN	SAR-GIN	SAR-GraphSAGE	SAR-GAT
Unsupervised	67.0±3.0	81.7±0.9	78.2±4.2	78.9±1.2	87.4±4.6
Supervised	88.0±1.0	91.2±1.7	89.8±2.1	89.3±2.0	91.8±1.5

$\mathbf{s}_{\text{our}} = \mathcal{F}_S(\mathbf{M}, \mathbf{H})$, as shown in Equation 2 and 5. The iterative interdependent refining mechanism enables our model to learn more effective node features and graph representation, which in turn lead to more accurate global saliency. The large performance superiority of our model over SAGPool and ChebyGIN in Table 2 in the paper demonstrates the advantage of our model. Besides, the qualitative comparison in Figure 4 and Figure 5 also reveal such advantage of our model over SAGPool and ChebyGIN.

Comparison with ‘ChebyGIN’ in the supervised setting for learning global node saliency. Our model, as well as other methods involved in comparison, is optimized for graph classification only using graph category labels but with no annotations of the global saliency. In this sense, the global saliency is learned in an unsupervised setting. To further investigate the performance of the our model in the supervised setting, we conduct experiments on TRIANGLES dataset by leveraging both the annotations of the global saliency and graph categories for supervision. We compare our method and the supervised version of ChebyGIN, which is presented in Table 3. The results show that all four instantiations of our method with different backbones consistently outperform ChebyGIN in the supervised setting, which reveals the effectiveness of our method in such setting. Besides, comparing the performance between unsupervised and supervised settings, both our method and ChebyGIN achieve substantial improvements. Such results indicate that the supervision on the global saliency yields more effective learning of global saliency and thereby more precise graph classification, which also implies the essential benefit of learning global node saliencies to the task of graph classification.

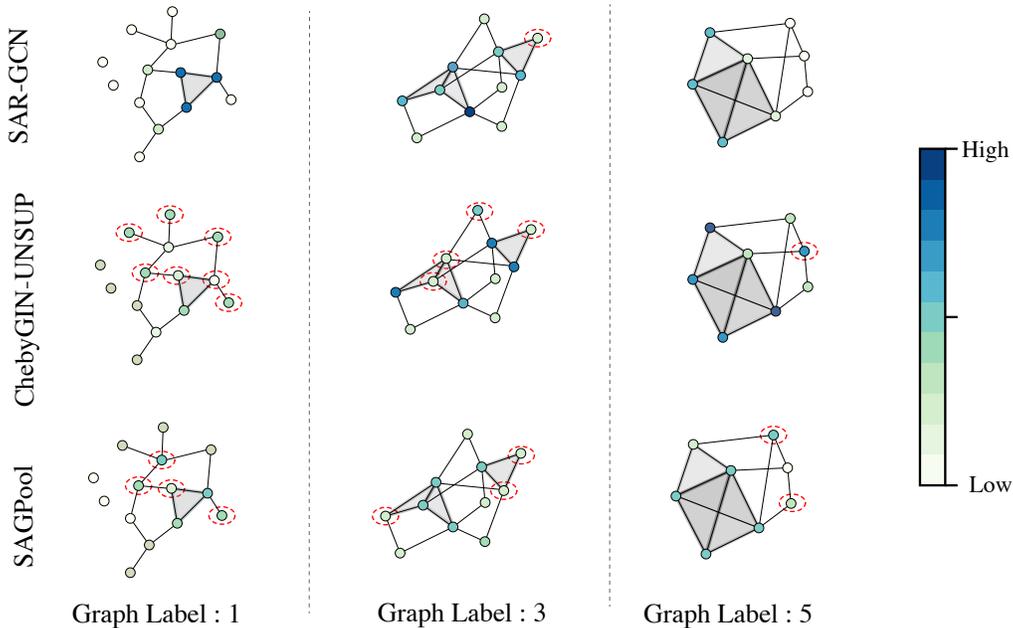


Figure 4: Visualization of learned node saliency by our SAR-GCN, ChebyGIN-unsup and SAGPool on three graphs randomly selected from TRIANGLES test data. The triangles are indicated in shaded regions and the number of triangles is given as the graph label. Our SAR-GCN tends to assign higher saliency scores to those nodes associated with more triangles, and thus shows higher accuracy than ChebyGIN-unsup and SAGPool in capturing the discriminative nodes associated with more triangles. The nodes that are assigned distinctly inaccurate saliency weights are marked with red dotted circles.

4.4. Qualitative Evaluation

In this set of experiments, we perform qualitative evaluation on the learned node saliency and the graph representation by our model, respectively.

Evaluation of the learned node saliency. To validate whether the learned node saliency by our method can accurately reflect the relevance of each node to the task of graph classification, we visualize the learned node saliency on *TRIANGLES* and *MUTAG* datasets in Figures 4 and 5, respectively. We also visualize the learned global attention scores by ChebyGIN-unsup and SAGPool which employs an attention mechanism to model the global saliency in an implicit manner.

As shown in Figure 4, our method tends to assign higher saliency scores

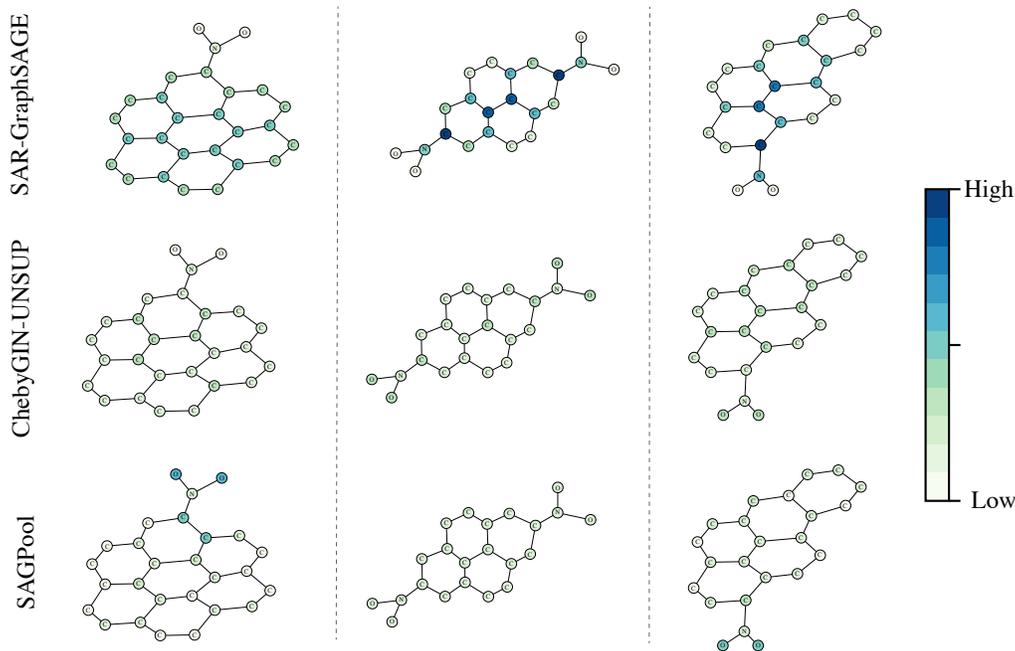


Figure 5: Visualization of learned node saliency by our SAR-GraphSAGE, ChebyGIN-unsup and SAGPool on three graphs, labeled as ‘mutagenic aromatic’, randomly selected from MUTAG. Our SAR-GraphSAGE can capture the key nodes located at the center of fused rings which are discriminative for ‘mutagenic aromatic’, whilst ChebyGIN-unsup and SAGPool either assign roughly uniform attention to all nodes or fail to capture the key nodes.

to those nodes associated with more triangles, which is reasonable since these nodes are more relevant to the graph category, namely the number of triangles in the graph. Although the learned attention scores by ChebyGIN-unsup and SAGPool also show a similar pattern, the results are less accurate than those of our method.

Figure 5 visualizes the learned saliency by our model as well as ChebyGIN-unsup and SAGPool on three positive samples labeled as ‘mutagenic aromatic’, randomly selected from *MUTAG*. Our model assigns higher saliency scores to the key nodes which locate at the center of fused rings, which is consistent with the chemical knowledge that ‘compounds with three or more fused rings are much more mutagenic, other factors being equal, than those with one or two’ [34]. In contrast, ChebyGIN-unsup and SAGPool either

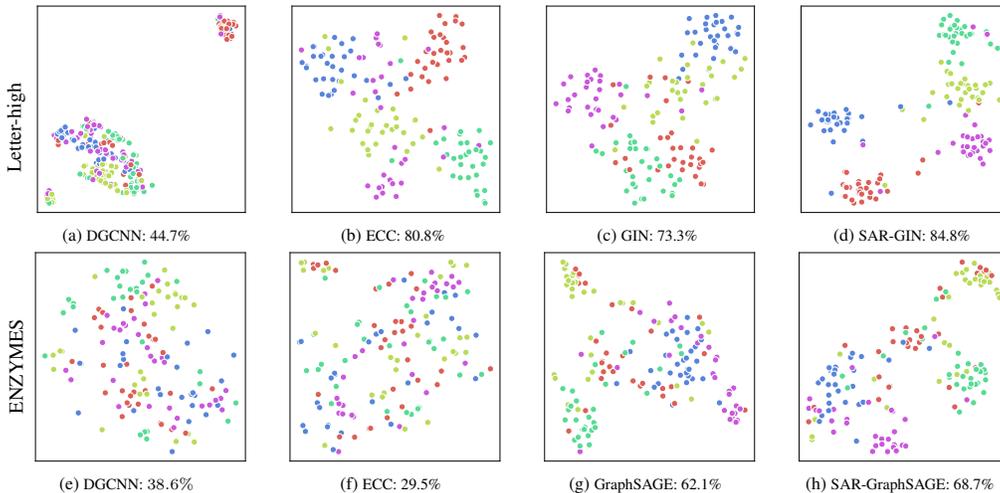


Figure 6: t-SNE maps of test data from five randomly selected classes of *Letter-high* and *ENZYMES*, constructed based on the learned graph representations by four different methods. We show results of our model with different backbones (GIN and GraphSAGE) on two datasets to show the robustness of our model. 30 randomly selected samples are visualized for each class. The classification accuracy for each method on each dataset is also presented.

learn a uniform distribution of attention weights for all nodes or fail to capture those key nodes. These visualizations clearly show the superiority of our method over ChebyGIN-unsup in learning the global node saliency for graph classification.

Evaluation of the learned graph representation. To gain further insight into the quality of the learned global graph representation by our method and other methods, we show t-SNE [44] maps of test data from *Letter-high* and *ENZYMES* datasets in Figure 6, constructed on the learned graph representations by four methods. We show results of our model with different backbones (GIN and GraphSAGE) on two datasets to show the robustness of our model. More effective graph representations typically lead to more separable clustering between classes in the t-SNE map. These maps reveal the consistent results with the quantitative evaluation among different methods in Section 4.3, which demonstrates that our method is able to learn more effective graph representation than other methods.

4.5. Model Complexity

We conduct extensive investigation into the model complexity of our model. In particular, we evaluate the model complexity in terms of model size and computational complexity, respectively.

Model complexity w.r.t. model size. We first conduct experiments to evaluate the model complexity of our model w.r.t. model size.

- **Is the performance gain of our model yielded from more learnable parameters compared to the backbone?** To investigate it, we compare the performance as well as the running time of both our *SAR-GCN* and the backbone GCN with the same amount of parameters by configuring the model structure of each module. Figure 7 (a) presents the results of both models on *TRIANGLES* dataset as a function of parameter amount, which show that our model consistently outperforms the backbone model significantly with the same amount of parameters, taking comparable running time. Furthermore, our model reaches performance saturation with a larger model capacity than the backbone, revealing the greater potential of our model. Both models exhibit overfitting to some degree when configured with excessive model capacity.
- **Cost performance in terms of the augmented model size of the Graph Neural Memory.** We measure the performance gain of increasing the model size of the Graph Neural Memory with a constant backbone GCN to evaluate the cost performance in terms of the augmented model size between our *SAR-GCN* and the backbone GCN. The results in Figure 7 (b) reveal that our model outperforms the backbone by a large margin even augmented with a relatively small size of Graph Neural Memory. Larger size of Graph Neural Memory yields more performance gain due to more modeling capacity.

Computational complexity. The computational cost of our model mainly consists of three parts, namely the cost incurred by Graph Neural Memory, the backbone and the saliency-aware regularization, respectively. We analyze the theoretical computational complexity of these three parts to calculate the overall complexity of our model. We take the instantiation of our model with GCN as an example.

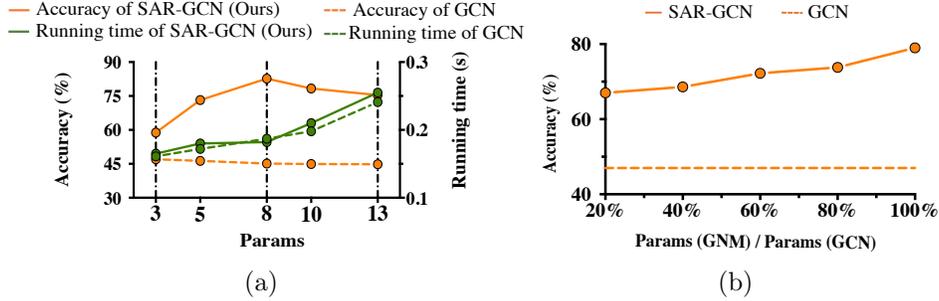


Figure 7: (a) Accuracy and average running time for one graph sample of our SAR-GCN and GCN on TRIANGLES dataset with increasing model size. Note that two models are compared with the same model size in this set of experiments. (b) Performance of our SAR-GCN on TRIANGLES as a function of increasing model size of the Graph Neural Memory (GNM) built upon a backbone (GCN) with constant model size. The model size of the Graph Neural Memory is expressed as a percentage of the model size of the backbone along the horizontal axis.

- **Complexity of Graph Neural Memory (GNM).** The computational complexity of Graph Neural Memory can be obtained by calculating the cost of each modeling step of \mathcal{F}_M shown in Equation 4. Specifically, the complexity for calculating \mathbf{m}_q , \mathbf{K} and \mathbf{V} is $O(d_M d)$, $O(N d_H d)$ and $O(N d_H d)$, respectively. Taking into account the complexity involved in the calculation of \mathbf{M}^l and \mathbf{M}' , the complexity for one layer of Graph Neural Memory is: $O(d_M d + N d_H d + N d + d^2)$.

In practice, d_M , d_H and d are approximately in the same order of magnitude, which results in the overall complexity of Graph Neural Memory consisting of L layers:

$$T_{\text{GNM}} = O(LkNd^2), \quad (14)$$

where $k \in \{1, 2\}$ is the iteration number of \mathcal{F}_M in Equation 4.

- **Complexity of the saliency-aware regularization.** Measuring the global node saliency formulated in Equation 5 involves three matrix multiplications, thus the complexity is $O(Ld_M d + LNd_H d + LNd)$ which approximately equals to $O(LNd^2)$. Besides, the complexity of regularization with the node saliency, performing weight fusion either by the weighted sum mechanism or the scaling mechanism, is $O(LN^2)$. Thus,

the total complexity for the saliency-aware regularization is:

$$T_{\text{SAR}} = O(LNd^2 + LN^2). \quad (15)$$

- **Complexity of the backbone (GCN).** The computational complexity of GCN can be derived from the modeling process of GCN in Equation 8. It has been analyzed [45] that the time complexity of GCN is:

$$T_{\text{GCN}} = O(LNd^2 + L|E|d), \quad (16)$$

where $|E|$ is the number of graph edges.

Combining the complexity of all three parts, the overall computational complexity of our *SAR-GCN* is:

$$T_{\text{SAR-GCN}} = O(LNd^2 + LN^2 + L|E|d). \quad (17)$$

Compared with the complexity of backbone (GCN) in Equation 16, our model has the same order of magnitude of complexity when satisfying $N = O(d^2)$.

To have a more comprehensive comparison between our model and the backbone, we conduct experiments to compare the performance of them with the same computational complexity. Figure 8 shows two sets of comparisons using different backbones on TRIANGLES dataset as a function of increasing computational complexity by adjusting the model scale. We observe that our model outperforms the backbones by a large margin consistently, while increasing the model complexity can potentially lead to more superiority. These results demonstrate that the performance gain of our model over the backbone does not result from the extra introduced computational complexity.

5. Conclusion

In this work we have presented the Saliency-Aware Regularized Graph Neural Network (*SAR-GNN*) for graph classification. Our *SAR-GNN* learns the global node saliency w.r.t. graph classification, and leverages it to regularize the neighborhood aggregation for feature learning. As a result, our method pays more attention to those nodes that are more relevant to graph classification and is able to learn more effective representations for the entire graph. We have shown the effectiveness of the proposed method built upon different types of graph neural networks by extensive experiments.

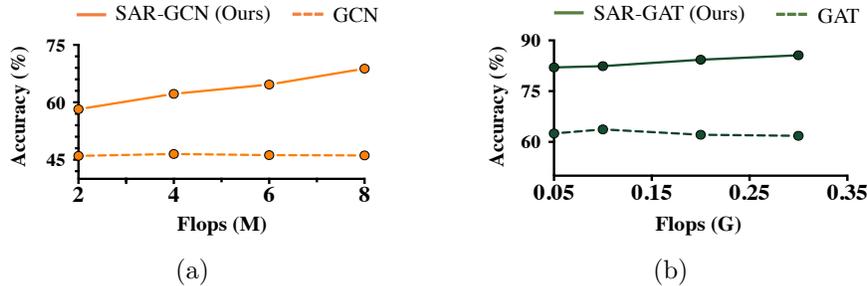


Figure 8: Accuracy of our model and the corresponding backbone on TRIANGLES dataset as a function of increasing computational complexity. GCN and GAT are used as backbones for comparisons in (a) and (b), respectively.

References

- [1] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, G. Monfardini, The graph neural network model, *IEEE transactions on neural networks* 20 (1) (2008) 61–80.
- [2] A. Micheli, Neural network for graphs: A contextual constructive approach, *IEEE Transactions on Neural Networks* 20 (3) (2009) 498–511.
- [3] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, in: *International Conference on Learning Representations*, 2017.
- [4] W. L. Hamilton, Z. Ying, J. Leskovec, Inductive representation learning on large graphs, in: *Advances in Neural Information Processing Systems*, 2017, pp. 1024–1034.
- [5] K. Xu, W. Hu, J. Leskovec, S. Jegelka, How powerful are graph neural networks?, in: *International Conference on Learning Representations*, 2019.
- [6] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio, Graph attention networks, in: *International Conference on Learning Representations*, 2018.
- [7] J. Lee, I. Lee, J. Kang, Self-attention graph pooling, in: *International Conference on Machine Learning*, 2019, pp. 3734–3743.

- [8] B. Knyazev, G. W. Taylor, M. Amer, Understanding attention and generalization in graph neural networks, in: *Advances in Neural Information Processing Systems*, 2019, pp. 4202–4212.
- [9] C. Ying, T. Cai, S. Luo, S. Zheng, G. Ke, D. He, Y. Shen, T.-Y. Liu, Do transformers really perform badly for graph representation?, in: *Advances in Neural Information Processing Systems*, 2021, pp. 28877–28888.
- [10] C. Bodnar, C. Cangea, P. Liò, Deep graph mapper: Seeing graphs through the neural lens, *Frontiers in big Data* 4 (2021).
- [11] M. Zhang, Z. Cui, M. Neumann, Y. Chen, An end-to-end deep learning architecture for graph classification, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [12] A. A. Leman, B. Weisfeiler, A reduction of a graph to a canonical form and an algebra arising during this reduction, *Nauchno-Technicheskaya Informatsiya* 2 (9) (1968) 12–16.
- [13] M. Simonovsky, N. Komodakis, Dynamic edge-conditioned filters in convolutional neural networks on graphs, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3693–3702.
- [14] R. Ying, J. You, C. Morris, X. Ren, W. L. Hamilton, J. Leskovec, Hierarchical graph representation learning with differentiable pooling, in: *Advances in Neural Information Processing Systems*, 2018, pp. 4800–4810.
- [15] Y. Ma, S. Wang, C. C. Aggarwal, J. Tang, Graph convolutional networks with EigenPooling, in: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019, pp. 723–731.
- [16] H. Gao, S. Ji, Graph U-Nets, in: *International Conference on Machine Learning*, 2019, pp. 2083–2092.
- [17] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, S. Y. Philip, A comprehensive survey on graph neural networks, *IEEE transactions on neural networks and learning systems* 32 (1) (2020) 4–24.

- [18] N. Shervashidze, P. Schweitzer, E. J. Van Leeuwen, K. Mehlhorn, K. M. Borgwardt, Weisfeiler-Lehman graph kernels., *Journal of Machine Learning Research* 12 (9) (2011).
- [19] T. Zhang, Y. Wang, Z. Cui, C. Zhou, B. Cui, H. Huang, J. Yang, Deep Wasserstein graph discriminant learning for graph classification, in: *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021, pp. 10914–10922.
- [20] T. Chen, S. Bian, Y. Sun, Are powerful graph neural nets necessary? A dissection on graph classification, in: *International Conference on Learning Representations*, 2019.
- [21] F. Errica, M. Podda, D. Bacciu, A. Micheli, A fair comparison of graph neural networks for graph classification, in: *International Conference on Learning Representations*, 2020.
- [22] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [23] J. Bruna, W. Zaremba, A. Szlam, Y. LeCun, Spectral networks and locally connected networks on graphs, in: *International Conference on Learning Representations*, 2014.
- [24] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, in: *Advances in Neural Information Processing Systems*, 2016, pp. 3844–3852.
- [25] M. Niepert, M. Ahmed, K. Kutzkov, Learning convolutional neural networks for graphs, in: *International Conference on Machine Learning*, 2016, pp. 2014–2023.
- [26] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [27] B. Bevilacqua, F. Frasca, D. Lim, B. Srinivasan, C. Cai, G. Balamurugan, M. M. Bronstein, H. Maron, Equivariant subgraph aggregation networks, in: *International Conference on Learning Representations*, 2022.

- [28] R. Ying, D. Bourgeois, J. You, M. Zitnik, J. Leskovec, Gnnexplainer: Generating explanations for graph neural networks, in: *Advances in Neural Information Processing Systems*, 2019.
- [29] D. Luo, W. Cheng, D. Xu, W. Yu, B. Zong, H. Chen, X. Zhang, Parameterized explainer for graph neural network, in: *Advances in Neural Information Processing Systems*, 2020.
- [30] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, in: *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.
- [31] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, S. Yan, Tokens-to-token vit: Training vision transformers from scratch on imagenet, in: *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 558–567.
- [32] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, S. Zagoruyko, End-to-end object detection with transformers, in: *European Conference on Computer Vision*, 2020, pp. 213–229.
- [33] A. Jaegle, F. Gimeno, A. Brock, A. Zisserman, O. Vinyals, J. Carreira, Perceiver: General perception with iterative attention, in: *International Conference on Machine Learning*, 2021.
- [34] A. K. Debnath, R. L. Lopez de Compadre, G. Debnath, A. J. Shusterman, C. Hansch, Structure-activity relationship of mutagenic aromatic and heteroaromatic Nitro compounds. Correlation with molecular orbital energies and hydrophobicity, *Journal of medicinal chemistry* 34 (2) (1991) 786–797.
- [35] I. Schomburg, A. Chang, C. Ebeling, M. Gremse, C. Heldt, G. Huhn, D. Schomburg, BRENDA, the enzyme database: updates and major new developments, *Nucleic acids research* 32 (2004) D431–D433.
- [36] K. M. Borgwardt, C. S. Ong, S. Schönauer, S. Vishwanathan, A. J. Smola, H.-P. Kriegel, Protein function prediction via graph kernels, *Bioinformatics* 21 (suppl_1) (2005) i47–i56.

- [37] P. Yanardag, S. V. N. Vishwanathan, Deep graph kernels, in: Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining, 2015, pp. 1365–1374.
- [38] K. Riesen, H. Bunke, IAM graph database repository for graph based pattern recognition and machine learning, in: Structural, Syntactic, and Statistical Pattern Recognition, 2008, pp. 287–297.
- [39] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization (2015).
- [40] G. Bouritsas, F. Frasca, S. P. Zafeiriou, M. Bronstein, Improving graph neural network expressivity via subgraph isomorphism counting, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [41] D. Q. Nguyen, T. D. Nguyen, D. Phung, Universal graph transformer self-attention networks, in: Companion Proceedings of the Web Conference 2022, 2022, pp. 193–196.
- [42] Y. Zhu, K. Zhang, J. Wang, H. Ling, J. Zhang, H. Zha, Structural landmarking and interaction modelling: A $\hat{\text{slim}}^{\hat{\text{a}}}$ network for graph classification, in: Proceedings of the AAAI Conference on Artificial Intelligence, 2022, pp. 9251–9259.
- [43] K. Zhang, Y. Zhu, J. Wang, J. Zhang, Adaptive structural fingerprints for graph attention networks, in: International Conference on Learning Representations, 2019.
- [44] L. van der Maaten, G. Hinton, Visualizing data using t-SNE, *Journal of Machine Learning Research* 9 (2008) 2579–2605.
- [45] Y. Q. Derrick Blakely, Jack Lanchantin, Time and space complexity of graph convolutional networks, in: Accessed on: Dec, 2021, 31., 2021.