# NENN: Incorporate Node and Edge Features in Graph Neural Networks

**Yulei Yang** YANGYULEI18@NUDT.EDU.CN
**Dongsheng Li** LDS1201@163.COM
*National University of Defense Technology, Chashang, China*

## Abstract

Graph neural networks (GNNs) have attracted an increasing attention in recent years. However, most existing state-of-the-art graph learning methods only focus on node features and largely ignore the edge features that contain rich information about graphs in modern applications. In this paper, we propose a novel model to incorporate **N**ode and **E**dge features in graph **N**eural **N**etworks (NENN) based on a hierarchical dual-level attention mechanism. NENN consists of node-level attention layer and edge-level attention layer. The two types of layers of NENN are alternately stacked to learn and aggregate embeddings for nodes and edges. Specifically, the node-level attention layer aims to learn the importance of the node based neighbors and edge based neighbors for each node, while the edge-level attention layer is able to learn the importance of the node based neighbors and edge based neighbors for each edge. Leveraging the proposed NENN, the node and edge embeddings can be mutually reinforced. Extensive experiments on academic citation and molecular networks have verified the effectiveness of our proposed graph embedding model.

**Keywords:** Graph Neural Network, Attention Mechanism, Graph Convolutional Network

## 1. Introduction

Convolutional Neural Networks (CNNs) have become very useful and successful techniques to process various data with regular grid-like structure Krizhevsky et al. (2012); Simonyan and Zisserman (2014); Redmon et al. (2016); He et al. (2018); Jégou et al. (2017). However, the non-Euclidean graphs containing all kinds of nodes and edges are ubiquitous in the real world, such as social networks, bioprotein networks and citation networks, which can not be easily represented due to the complex and irregular structure.

In recent years, there is a growing interest in graph presentation learning methods. Representation learning is to map high-dimensional features from a graph to a low-dimensional vectors, so that many downstream problems can be easily solved, including node classification Ribeiro et al. (2017); Jacob et al. (2014); Donnat et al. (2018), link prediction Berg et al. (2017); Zhang and Chen (2018); Hasanzadeh et al. (2019) and graph classification Ying et al. (2018); Murphy et al. (2019); Lee et al. (2019).

Inspired by spectral graph theory, Bruna et al. (2013) generalizes CNNs to the graph domain based on the feature decomposition of graph Laplace matrix. In order to reduce the overhead of the decomposition, ChebNet Defferrard et al. (2016) is proposed to approximate convolution kernel with Chebyshev polynomials. As a pilot work, Kipf and Welling (2017)

proposes a graph convolutional network (GCN) as the first order approximation of ChebNet, which greatly simplifies the convolution filters by limiting the receptive field to the 1-hop neighbors for each node. Finally, the GCN model is successfully applied to semi-supervised node classification and achieves state-of-the-art performance. The basic idea behind GCN is to map a high-dimensional node representation to a low-dimensional vector by transforming, propagating, aggregating and updating node features across edges in a graph. Nevertheless, GCN model is essentially a spectral approach working on transductive learning tasks. As a result, GCN can not run on large and dynamic graphs effectively. To address the limitations of GCN, GraphSAGE Hamilton et al. (2017) extends GCN from a transductive approach to an inductive one using a spatial-based method to train embeddings for previously unseen nodes. GraphSAGE restricts neighborhood sampling to learn how to aggregate node features rather than train fixed node embeddings. In addition, GraphSAGE also proposes a mini batch training algorithm, which solves the problem that GCN cannot be applied to large graphs. Graph Attention Network (GAT) Veličković et al. (2018), a newfangled attention-based graph neural network, trains weight coefficients associated with neighbors for each node to learn node embeddings. It has demonstrated the effectiveness in graph embedding and shown the superiority over the previous methods.

Despite the success of existing graph neural networks, there are two enormous challenges. On the one hand, almost all previous literatures only leverage the node features and completely ignore the edge features that are completely likely to contain important information. For example, in molecule networks, a node represents an atom while an edge represents a bond connecting two atoms. A bond usually has some simple edge features (e.g., bond type, atom pair type, bond order, conjugated, ring status, aromaticity), which are closely related to atom features. On the other hand, how to measure the importance of neighborhood as well as the connecting edges or nodes is not fully considered.

In order to address the aforementioned challenges, we propose a novel graph neural network, named NENN, which incorporates node and edge features based on a dual-level attention mechanism, including node-level and edge-level attentions. Specifically, we aim to to learn the importance of node based neighbors and edge based neighbors and aggregate embeddings for each node in the node-level attention layer. Similarly, the embedding of each edge is generated in the edge-level attention layer.

We conduct extensive experiments on node classification, graph classification and graph regression to verify the effectiveness of the proposed NENN. For node classification, we use the benchmark citation network datasets: Cora, Citeseer Sen et al. (2008) and Pubmed Galileo Mark Namata and Huang (2012). For graph classification and graph regression, we demonstrate the proposed NENN is able to effectively generate node and edge embeddings by incorporating node and edge features on multiple molecular datasets: Tox21 Wu et al. (2018), HIV Wu et al. (2018), Freesolv Mobley and Guthrie (2014), and Lipophilicity Wu et al. (2018). The results show that the proposed NENN outperforms relevant baselines by a significant margin.

In a nutshell, our main contributions are summarized as follows:

- We propose a novel graph neural network (NENN) that incorporates both node and edge features, which can learn node embeddings as well as edge embeddigns simultaneously.

- To the best of our knowledge, this is the first attempt to take the influences of neighbors for nodes and edges into consideration based on a hierarchical dual-level attention mechanism , including node-level and edge-level attentions.

- We transform the roles of nodes and edges and extend them to the neighbor set of each other, which strengthens the connection between edges and nodes.

- Various graph-related tasks, including graph classification, graph regression, and node classification, are used to verify the scalability and generality for NENN.

- We conduct extensive experiments on benchmark data sets.The citation networks Sen et al. (2008); Galileo Mark Namata and Huang (2012) and multiple molecular networks Wu et al. (2018); Mobley and Guthrie (2014) are applied to demonstrate the effectiveness of our model. The results show the superiority of the proposed NENN compared with the state-of-the-art baselines.

## 2. Related Work

The real-world data usually appears in the form of non-Euclidea graphs. On of the most severe challenges in graph representations is to efficiently exploit the node and topology information. At present, graph representation learning methods can be roughly divided into three parts: matrix factorization, random walks and graph neural networks.

**Factorization-based approaches**. The key idea of matrix factorization is that the relation matrix (e.g. adjacency matrix and Laplace matrix) is decomposed to yield the low-dimensional representations. For example, Grarep Cao et al. (2015) reduces the dimension of the relation matrix by SVD decomposition to get the k-step network vertex representation for weighted graphs. HOPE Ou et al. (2016) preserves high-order proximities and captures the asymmetric transitivity for directed graphs. However, these methods can not efficiently process large-scale graphs since they have huge performance overheads for matrix factorization.

**Random walk**. Just as its name implies, a random walk on a graph starts with a node and recursively connects with a randomly selected neighbor until a threshold. DeepWalk Perozzi et al. (2014) is the first attempt to learn latent representations leveraging truncated random walks. node2vec Grover and Leskovec (2016) further designs a biased random walk to efficiently explore diverse neighborhoods based on DFS and BFS strategies. By recursively compressing the input graph to smaller but structurally similar graphs, HARP Chen et al. (2018) captures the global topological information about the input graph and generates presentations on this smaller graph during a random walk. Nevertheless, random walk methods are not the most successful methods so far. The shallow embedding methods use unshareable parameters and functions, which makes it impossible to embed nodes of a large-scale graphs to a low-dimensional space. Besides, they are only applied to learn embeddings for fixed graphs in the transductive setting, consequently, do not naturally generalize to unseen nodes.

**Graph neural networks**. In order to address the limitations of the previous methods, in recent years, graph neural networks are proposed to learn node embeddings. Graph neural network has many branches so far, but we mainly focus on the methods based on

graph convolution. Inspired by the successes of CNNs in image recognition, GCN Kipf and Welling (2017) stacks graph convolutional layers to aggregate local information from neighbors and encodes nodes into vectors. GraphSAGE Hamilton et al. (2017) is a novel inductive approach that generates latent embeddings for unseen nodes. Attention mechanisms have been widely applied to many tasks in deep learning. GAT Veličković et al. (2018) is introduced to learn the importance coefficients for nodes and its neighbors rather than treats the neighborhood information equally. However, the above graph neural networks not only ignore the essential edge features but also fail to differentiate the influences of their connecting edges.

**Edge-related Work**. To address the above issues, some models are proposed. Message passing neural network (MPNN) Justin Gilmer (2017), a generalized model, consisting of two phases: multiple message passing phases and readout phase, is proposed to predict molecular properties. Although MPNN adds edge information to the message passing phase, its passing mechanism can not recognize the correlation between nodes and edges. However, the node based neighbors and the edge based neighbors defined in two types of convolution layers of NENN can integrate the adjacency relationship between edges and nodes, which captures the local structure information about the graph. RGCN Schlichtkrull et al. (2018) uses a simple forward-pass rule:

$$h_i^{(l+1)} = \sigma \left( \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_i^r} \frac{1}{c_{i,r}} W_r^{(l)} h_j^{(l)} + W_0^{(l)} h_i^{(l)} \right) \tag{1}$$

where $W_r$ represents the weight matrix of relation $r \in \mathcal{R}$. This simple aggregation is too hard to calculate and performs poorly. EGNN Gong and Cheng (2019) introduces attention mechanism to explore edge features. In EGNN, except for the original edge features used in the first layer, the attention coefficient between the two nodes of in $l$-th layer is used as the edge features in $(l + 1)$-th layer, which greatly leads to the loss of edge information. Consequently, EGNN inherently reinforces the embedding of nodes with edge features, while edge embeddings are not enhanced by node embeddings. As a result, it does not suitable to edge embedding learning and link prediction. Compared with EGNN, NENN adopts a hierarchical dual-level attention mechanism where the the role of edge and node are alternated, to keep the edge features as a vector rather than a one-dimensional attention coefficient. In addition, NENN considers the relationship between edges and nodes and the graph structure fairly and comprehensively, which enables NENN to learn node, edge and even graph representations. CensNet Jiang et al. (2019) embeds both nodes and edges to a latent feature space by using line graph of the original undirected graph. However, the CensNet uses approximated spectral graph convolution in the layer-wise propagation, which makes the CensNet can not process large graphs and directed graphs. On the contrary, the basic idea behind of the proposed NENN is based on spatial domain, which enables various graphs to be processed. Although both CensNet and NENN adopt alternate convolution, NENN and CensNet use two completely different ideas. CensNet uses the original graph and its line graph Harary and Norman (1960) to perform alternate convolution, while NENN extends the neighboring nodes to the neighbors of an edge and the neighboring edges to the neighbors of a node. In addition, NENN also adopts the attention mechanism to learn
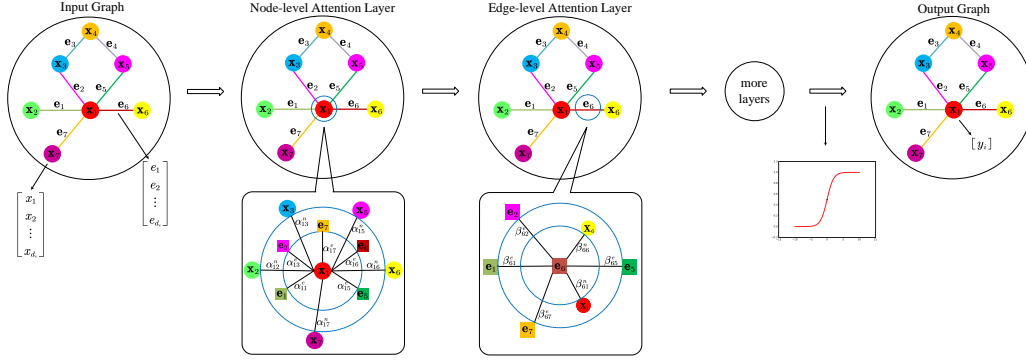
Figure 1: The overall framework of the proposed NENN for node embedding generation. The dual-level attention layers (i.e., node-level attention layer and edge-level attention layer) of NENN are alternately stacked to learn node and edge embeddings.

an importance coefficient from neighbors for each node or edge , which enables NENN to learn more efficient embeddings.

## 3. The Proposed Method: NENN

### 3.1. Preliminary

In this section, we propose a novel model to incorporate **N**ode and **E**dge features in graph **N**eural **N**etworks (NENN) based on a hierarchical dual-level attention mechanism. For a graph $G(V, E)$ with node features and edge features, where $V$ defines a set of $N_v = |V|$ nodes, $E$ is a set of $N_e = |E|$ edges. Let $\mathbf{X} = \{\mathbf{x}_i | i \in N_v\} \in \mathbb{R}^{N_v \times d_v}$ be node feature matrix, where $\mathbf{x}_i \in \mathbb{R}^{d_v}$ represents $d_v$-dimensional feature vector of node $i$. Let $\mathbf{E} = \{\mathbf{e}_i | i \in N_e\} \in \mathbb{R}^{N_e \times d_e}$ be edge feature matrix, where $\mathbf{e}_i \in \mathbb{R}^{d_e}$ denotes $d_e$-dimensional feature vector of edge $i$.

**Definition 1 (Node Based Neighbors)** *Given a graph $G(V, E)$, the node based neighbors $\mathcal{N}_i$ of node $i$ are defined as the set of nodes which connect with node $i$. In the same way, the node set $\mathcal{N}_j$ connected by edge $j$ represents the node based neighbors of edge $j$. Specially, the node based neighbors of node $i$ include itself.*

**Example 1** *As shown in Figure 1, nodes are represented by circles while edges are represented by squares. In node-level attention layer, the node based neighbors $\mathcal{N}_1$ of the red node whose feature vector is $\mathbf{x}_1$ denote the nodes $\{1, 2, 3, 5, 6, 7\}$ which consists of one-hop neighboring nodes and node 1.*

**Definition 2 (Edge Based Neighbors)** *Given a graph $G(V, E)$, the edge based neighbors $\mathcal{E}_i$ of node $i$ are defined as the set of edges connecting with node $i$. Similarly, the edge based neighbors $\mathcal{E}_j$ of edge $j$ are defined as the set of edges connecting with edge $j$. Specially, the edge based neighbors of edge $j$ include itself.*

**Example 2** *As shown in Figure 1, in edge-level attention layer, the edge based neighbors of the brown edge (i.e., $\mathbf{e}_6$) denote the neighboring edges $\{1, 2, 5, 6, 7\}$. In addition, the node based neighbors of the brown edge (i.e., $\mathbf{e}_6$) denote the one-hop neighboring node set $\{1, 6\}$.*

Figure 1 shows the overall process of the proposed NENN for node embeddings generation. The proposed NENN consists of two types of attention layers, node-level attention layer and edge-level attention layer. The dual-level attention layers are alternately stacked to learn node and edge embeddings. In the node-level attention layer, we aim to learn the node based neighbors importance $\alpha_{ij}^n$ and edge based neighbors importance $\alpha_{ij}^e$ for each node. In the edge-level attention layer, we aim to learn the node based neighbors importance $\beta_{ij}^n$ and edge based neighbors importance $\beta_{ij}^e$ for each edge. With the learned importance coefficients, we can aggregate and update node and edge embeddings in order.

## 3.2. Node-level Attention Layer

In the node-level attention layer, NENN focuses to learn node embeddings with the help of edge features that contain significant information. It is clear to observe that different neighbors of each node play a different role and show different importance in generating node embedding. For this reason, we introduce a node-level attention to learn the importance coefficients of node based neighbors and edge based neighbors for node $i$.

In the $l$-th layer, suppose the input features consist of node features, $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_{N_v}\}$, $\mathbf{x}_i \in \mathbb{R}^{d_v^{(l)}}$, and edge features, $\mathbf{E} = \{\mathbf{e}_1, \mathbf{e}_2, ..., \mathbf{e}_{N_e}\}$, $\mathbf{e}_i \in \mathbb{R}^{d_e^{(l)}}$. The importance of node $j$ or edge $k$ ($j \in \mathcal{N}_i, k \in \mathcal{E}_i$) to node $i$ can be reformulated as follows:

$$\mathbf{e}_{ij}^n = Att_{node}^n(W_n \mathbf{x}_i, W_n \mathbf{x}_j) \tag{2}$$

$$\mathbf{e}_{ik}^e = Att_{node}^e(W_n \mathbf{x}_i, W_e \mathbf{e}_k) \tag{3}$$

Here, $Att_{node}^n$ and $Att_{node}^e$ denote the deep neural networks, which perform node-level attention for node $i$. $W_n \in \mathbb{R}^{d_v^l \times d_v^{l+1}}$ and $W_e \in \mathbb{R}^{d_e^l \times d_e^{l+1}}$ are the learnable weight matrices that linearly transform the input features into high-level features. The importance coefficient $\mathbf{e}_{ij}^n$ means how important node $j$ is to node $i$, while $\mathbf{e}_{ij}^e$ represents the influence of edge $j$ to node $i$.

Then the structure information is integrated into the proposed NENN via masked attention, which means the embedding of node $i$ depends only on neighboring nodes j or edges k. Next, the importance coefficient of node $j$ to node $i$ is normalized via the softmax function:

$$\alpha_{ij}^n = softmax_j(\mathbf{e}_{ij}^n) = \frac{\exp\left(\sigma\left(a_n^T\left[W_n \mathbf{x}_i || W_n \mathbf{x}_j\right]\right)\right)}{\sum_{k \in \mathcal{N}_i} \exp\left(\sigma\left(a_n^T\left[W_n \mathbf{x}_i || W_n \mathbf{x}_k\right]\right)\right)} \tag{4}$$

where $||$ is the concatenation operation, $a_n \in \mathbb{R}^{2d_v^{l+1}}$ is the parameter vector of a single-layer feed-forward network and $\sigma$ denotes the activation function (e.g. LeakyReLU).

Then, the importance coefficient $\alpha_{ik}^e$ of node $i$ and edge $k$ can be derived as :

$$\alpha_{ik}^e = softmax_k(\mathbf{e}_{ik}^e) = \frac{\exp\left(\sigma\left(a_e^T\left[W_n \mathbf{x}_i || W_e \mathbf{e}_k\right]\right)\right)}{\sum_{j \in \mathcal{E}_i} \exp\left(\sigma\left(a_e^T\left[W_n \mathbf{x}_i || W_e \mathbf{e}_j\right]\right)\right)} \tag{5}$$

where $a_e \in \mathbb{R}^{d_v^{l+1}+d_e^{l+1}}$ is the parameter vector of a single-layer feed-forward network.

After learning the importance $\alpha_{ij}^n$ and $\alpha_{ik}^e$, the embedding $\mathbf{x}_{\mathcal{N}_i}$ of node $i$'s node based neighbors can be aggregated with the corresponding importance coefficients:

$$\mathbf{x}_{\mathcal{N}_i} = \sigma\Big(W_n \cdot \text{MEAN}(\{\alpha_{ij}^n \mathbf{x}_j, \forall j \in \mathcal{N}_i\})\Big) \tag{6}$$

where MEAN is a mean aggregator. Then, the embedding $\mathbf{x}_{\mathcal{E}_i}$ of node $i$'s edge based neighbors can be aggregated as follows:

$$\mathbf{x}_{\mathcal{E}_i} = \sigma\Big(W_e \cdot \text{MEAN}(\{\alpha_{ik}^e \mathbf{e}_k, \forall k \in \mathcal{E}_i\})\Big) \tag{7}$$

Finally, with the edge based neighbors' embedding $\mathbf{x}_{\mathcal{E}_i}$ and node based neighbors' embedding $\mathbf{x}_{\mathcal{N}_i}$, the embedding of node $i$ in the $(l+1)$-th layer can be combined:

$$\mathbf{x}_i^{(l+1)} = \text{CONCAT}(\mathbf{x}_{\mathcal{N}_i}^{(l)}, \mathbf{x}_{\mathcal{E}_i}^{(l)}) \tag{8}$$

where CONCAT represents concatenation operation. $\mathbf{x}_i^{(l+1)}$ is the returned embedding for node $i$ in the $(l)$-th layer.

### 3.3. Edge-level Attention Layer

We use edge features to enhance the node embeddings in node-level attention layer while the node features are fused to learn edge embeddings in edge-level attention layer. To update the edge embeddings, we first learn the importance of node based neighbors and edge based neighbors for each edge.

The importance coefficient $\beta_{ij}^e$ of edge $k$ ($k \in \mathcal{E}_i$) to edge $i$ is normalized via the softmax function:

$$\beta_{ik}^e = \frac{\exp\Big(\sigma\Big(q_e^T\Big[W_e\mathbf{e}_i || W_e\mathbf{e}_k\Big]\Big)\Big)}{\sum_{j \in \mathcal{E}_i} \exp\Big(\sigma\Big(q_e^T\Big[W_e\mathbf{e}_i || W_e\mathbf{e}_j\Big]\Big)\Big)} \tag{9}$$

where $q_e \in \mathbb{R}^{2d_e^{l+1}}$ is an attention vector. The importance coefficient $\beta_{ij}^n$ of node $j$ ($j \in \mathcal{N}_i$) to edge $i$ is normalized via the softmax function:

$$\beta_{ij}^n = \frac{\exp\Big(\sigma\Big(q_n^T\Big[W_e\mathbf{e}_i || W_n\mathbf{x}_j\Big]\Big)\Big)}{\sum_{k \in \mathcal{N}_i} \exp\Big(\sigma\Big(q_n^T\Big[W_e\mathbf{e}_i || W_n\mathbf{x}_k\Big]\Big)\Big)} \tag{10}$$

where $q_n \in \mathbb{R}^{d_v^{l+1}+d_e^{l+1}}$ is an attention vector.

Then, edge $i$'s middle node based neighbors embedding $\mathbf{e}_{\mathcal{E}_i}$ and edge based neighbors $\mathbf{e}_{\mathcal{N}_i}$ can be generated by a mean aggregator:

$$\mathbf{e}_{\mathcal{E}_i} = \sigma\Big(W_e \cdot \text{MEAN}(\{\beta_{ik}^e \mathbf{e}_k, \forall k \in \mathcal{E}_i\})\Big) \tag{11}$$

$$\mathbf{e}_{\mathcal{N}_i} = \sigma\Big(W_n \cdot \text{MEAN}(\{\beta_{ij}^n \mathbf{x}_j, \forall j \in \mathcal{N}_i\})\Big) \tag{12}$$

---

**Algorithm 1** Mini-batch NENN node embeddings generation algorithm

---

**Input:** Subgraph $G'(V', E')$;
    node features $\{\mathbf{x}_i, \forall_i \in V'\}$;
    edge features $\{\mathbf{e}_i, \forall_i \in E'\}$;
    network depth $L$;

**Output:** final node embeddings $\{\mathbf{x}_i^{(L)}, \forall_i \in V'\}$;

**1** $\mathbf{x}_i^{(0)} \leftarrow \mathbf{x}_i, \forall i \in V'^{(0)}$ **for** $l = 0 \cdots L$ **do**

**2**   find the node based neighbors $\mathcal{N}_i$ and edge based neighbors $\mathcal{E}_i$ **if** *layer l is a node-level attention layer or l = L* **then**

**3**    **for** *each node $i \in V'^{(l)}$* **do**

**4**     Calculate the importance coefficient $\alpha_{ij}^{n\,(l)}$ and $\alpha_{ik}^{e\,(l)}$   Calculate the embedding of node based neighbors $\mathbf{x}_{\mathcal{N}_i}^{(l)}$ and edge based neighbors $\mathbf{x}_{\mathcal{E}_i}^{(l)}$

     $\mathbf{x}_i^{(l+1)} \leftarrow \text{CONCAT}(\mathbf{x}_{\mathcal{N}_i}^{(l)}, \mathbf{x}_{\mathcal{E}_i}^{(l)})$

**5**    **end**

**6**   **end**

**7**   **if** *layer l is an edge-level attention layer* **then**

**8**    **for** *each edge $i \in E'^{(l)}$* **do**

**9**     Calculate the importance coefficient $\beta_{ij}^{n\,(l)}$ and $\beta_{ik}^{e\,(l)}$   Calculate the embedding of node based neighbors $\mathbf{e}_{\mathcal{N}_i}^{(l)}$ and edge based neighbors $\mathbf{e}_{\mathcal{E}_i}^{(l)}$

     $\mathbf{e}_i^{(l+1)} \leftarrow \text{CONCAT}(\mathbf{e}_{\mathcal{N}_i}^{(l)}, \mathbf{e}_{\mathcal{E}_i}^{(l)})$

**10**    **end**

**11**   **end**

**12** **end**

**13** $\mathbf{x}_i^{(L)} \leftarrow \text{CONCAT}(\mathbf{x}_{\mathcal{N}_i}^{(L-1)}, \mathbf{x}_{\mathcal{E}_i}^{(L-1)})$

---

Similarly, the embedding of edge $i$ in the $(l+1) - th$ layer can be derived as follows:

$$\mathbf{e}_i^{(l+1)} = \text{CONCAT}(\mathbf{e}_{\mathcal{N}_i}^{(l)}, \mathbf{e}_{\mathcal{E}_i}^{(l)}) \tag{13}$$

With the learned importance, the proposed NENN can pay more attention to some meaningful nodes or edges for the specific task. Note that the importance coefficients are asymmetric which means the influence between different roles in a graph can be quite different (i.g., $\alpha_{ij}^n \neq \alpha_{ji}^n$, $\beta_{ij}^e \neq \beta_{ij}^e$).

The computation of attention can be easily parallelized across all edges and nodes, which means the hierarchical dual-level attention of the proposed NENN is highly efficient. The overall time complexity is linear to the number of nodes and edges.

To process large-scale graphs in the real world, we construct some subgraphs $G' \in G$ according to Hamilton et al. (2017). The mini-batch training process of the proposed NENN for node embedding generation is shown in Algorithm 1.

|  | Cora | Citeseer | Pubmed |
|---|---|---|---|
| # Nodes | 2,708 | 3,327 | 19,717 |
| # Edges | 5,429 | 4,732 | 44,338 |
| # Node Features | 1,433 | 3,703 | 500 |
| # Edge Features | 2 | 2 | 2 |
| # Node Classes | 7 | 6 | 3 |

Table 1: Dataset statistics of citation networks for semi-supervised node classification.

|  | Tox21 | HIV | Lipophilicity | Freesolv |
|---|---|---|---|---|
| # Graphs | 7,831 | 41127 | 4,200 | 642 |
| # Node Features | 25 | 25 | 25 | 25 |
| # Edge Features | 55 | 80 | 34 | 21 |
| # Graph Classes | 12 | 3 | - | - |

Table 2: Dataset statistics of molecular networks for graph classification and regression.

### 3.4. Variants of NENN

In order to verify the validity of the proposed NENN in more detail, we implement our model with different settings according to how to aggregate and update the node and edge features. Specifically, there are four different variants: NENN-NCEC, NENN-NCEA, NENN-NAEC and NENN-NAEA (abbreviated as NENN). Specifically, N, E, C, A represent nodes, edges, convolution operation and attention mechanism, respectively. For example, NENN-NCEA represents node features aggregated by convolution operation and edge features aggregated by attention mechanism.

## 4. Experiments

We evaluate the proposed NENN on three benchmark tasks: (i) semi-supervised node classification on citation networks Cora, Citeseer and Pubmed; (ii) multi-task graph classification on multiple molecular datasets Tox21 and HIV; (iii) graph regression on molecular datasets Lipophilicity and Freesolv.

### 4.1. Benchmark Datasets

We conduct extensive experiments on citation networks and molecular networks to demonstrate the effectiveness of the proposed NENN. Generally, each citation network corresponds to a graph where nodes represent documents and edges represent citation relationships between documents. Different from citation networks, each kind of molecular dataset consists of multiple graphs. Specifically, compounds, atoms, bonds represent graphs, nodes and edges, respectively. More detailed dataset statistics are shown in Table 1 and Table 2.

**Cora, Citeseer, and Pubmed**. Cora, Citeseer Sen et al. (2008) and Pubmed Galileo Mark Namata and Huang (2012), are citation networks and widely used as benchmark datasets for semi-supervised node classification in GCN, GraphSAGE, GAT, EGNN, CensNet, etc.

**Tox21 and HIV**. Tox21 Wu et al. (2018) is a public dataset of toxicity measurements, which comprises 7831 compounds on 12 different quantitative toxicity measurements in-

| Dataset | Cora | | | Citeseer | | | Pubmed | | |
|---|---|---|---|---|---|---|---|---|---|
| Label rate | 0.5% | 1% | 3% | 0.3% | 0.5% | 1% | 0.03% | 0.05% | 0.1% |
| GCN | $53.66 \pm 0.03$ | $62.50 \pm 0.02$ | $75.77 \pm 0.01$ | $39.83 \pm 0.02$ | $48.53 \pm 0.03$ | $59.47 \pm 0.03$ | $58.84 \pm 0.04$ | $66.23 \pm 0.03$ | $74.20 \pm 0.03$ |
| GraphSAGE | $38.48 \pm 0.04$ | $50.51 \pm 0.03$ | $66.33 \pm 0.03$ | $28.70 \pm 0.04$ | $35.80 \pm 0.02$ | $53.36 \pm 0.02$ | $46.00 \pm 0.02$ | $55.42 \pm 0.03$ | $60.54 \pm 0.10$ |
| GAT | $43.45 \pm 0.02$ | $49.66 \pm 0.02$ | $57.85 \pm 0.04$ | $32.53 \pm 0.03$ | $39.46 \pm 0.05$ | $47.57 \pm 0.06$ | $51.25 \pm 0.03$ | $52.56 \pm 0.02$ | $61.55 \pm 0.03$ |
| MPNN | $49.34 \pm 0.03$ | $51.35 \pm 0.03$ | $56.55 \pm 0.02$ | $34.45 \pm 0.02$ | $41.24 \pm 0.02$ | $49.34 \pm 0.02$ | $51.57 \pm 0.02$ | $53.57 \pm 0.02$ | $62.44 \pm 0.02$ |
| CensNet | $59.26 \pm 0.02$ | $69.56 \pm 0.04$ | $80.56 \pm 0.02$ | $51.54 \pm 0.02$ | $59.64 \pm 0.02$ | $64.34 \pm 0.02$ | $63.66 \pm 0.02$ | $67.87 \pm 0.02$ | $71.78 \pm 0.02$ |
| EGNN | $61.65 \pm 0.02$ | $\mathbf{68.86 \pm 0.02}$ | $79.69 \pm 0.03$ | $\mathbf{51.64 \pm 0.02}$ | $57.56 \pm 0.03$ | $66.62 \pm 0.02$ | $62.58 \pm 0.04$ | $64.50 \pm 0.02$ | $74.50 \pm 0.02$ |
| NENN-NAEA | $\mathbf{65.53 \pm 0.02}$ | $69.44 \pm 0.02$ | $\mathbf{82.57 \pm 0.02}$ | $50.66 \pm 0.02$ | $\mathbf{60.64 \pm 0.02}$ | $\mathbf{68.23 \pm 0.02}$ | $\mathbf{65.56 \pm 0.02}$ | $\mathbf{69.50 \pm 0.02}$ | $\mathbf{77.70 \pm 0.03}$ |
| NENN-NCEC | $60.46 \pm 0.03$ | $67.55 \pm 0.02$ | $80.58 \pm 0.02$ | $49.56 \pm 0.03$ | $59.55 \pm 0.02$ | $66.65 \pm 0.02$ | $64.55 \pm 0.02$ | $67.55 \pm 0.02$ | $75.56 \pm 0.03$ |
| NENN-NCEA | $60.03 \pm 0.02$ | $68.30 \pm 0.02$ | $80.40 \pm 0.04$ | $50.03 \pm 0.02$ | $59.50 \pm 0.02$ | $65.50 \pm 0.02$ | $65.50 \pm 0.02$ | $68.50 \pm 0.02$ | $76.40 \pm 0.02$ |
| NENN-NAEC | $60.56 \pm 0.02$ | $67.38 \pm 0.02$ | $80.88 \pm 0.02$ | $49.66 \pm 0.01$ | $59.66 \pm 0.02$ | $67.52 \pm 0.02$ | $65.45 \pm 0.02$ | $66.48 \pm 0.02$ | $73.56 \pm 0.02$ |

Table 3: Classification accuracies on citation networks.

cluding AR, AhR, AR-LBD, etc. The HIV Wu et al. (2018) dataset originates from the Drug Therapeutics Program AIDS Antiviral Screen, which measures the ability of HIV replication for 41127 compounds. The two datasets are usded to activity prediction (i.e. binary graph classification) that labels compounds as either "active" or "inactive".

**Lipophilicity and Freesolv**. Lipophilicity Wu et al. (2018) is a public dataset used to measure the affects both membrane permeability and solubility, which provides experimental results of octanol/water distribution coefficient (logD at pH 7.4) of 4,200 compounds. Also, the Free Solvation Database (Freesolv) Mobley and Guthrie (2014) is a common dataset providing experimental and calculated results of hydration free energies for 642 small molecules in water. According to the characteristics of the two datasets, we conduct extensive graph regression experiments to predict solvation energies or solubility.

## 4.2. Baselines

We compare with some state-of-the-art baselines to verify the effectiveness of our node-level attention and edge-level attention of the proposed NENN. In addition, we also used four variants of NENN to verify the effect of convolution or attention mechanism on the edge-level layer or node-level layer.

- **GCN** Kipf and Welling (2017): A transductive graph convolutional network which aggregates local information from neighbors and encodes nodes into low-dimensional vectors.

- **GraphSAGE** Hamilton et al. (2017): A novel inductive approach that generates latent embeddings for unseen nodes. Instead of learning a fixed representation for each node, GraphSAGE learns a function that aggregate features from its neighbors by sampling its neighbors.

- **GAT** Veličković et al. (2018): A graph convolution network based on attention mechanism. GAT learns an importance coefficient for each neighbor of a node instead of ignoring the differences of neighbors like GCN.

- **MPNN** Justin Gilmer (2017): A message passing neural network based on message passing and readout. MPNN is proposed to predict molecular properties.

- **CensNet** Jiang et al. (2019): A spectral domain based graph convolutional network. CensNet learns the node and edge feature embeddings simultaneously based on original graph and its line graph.
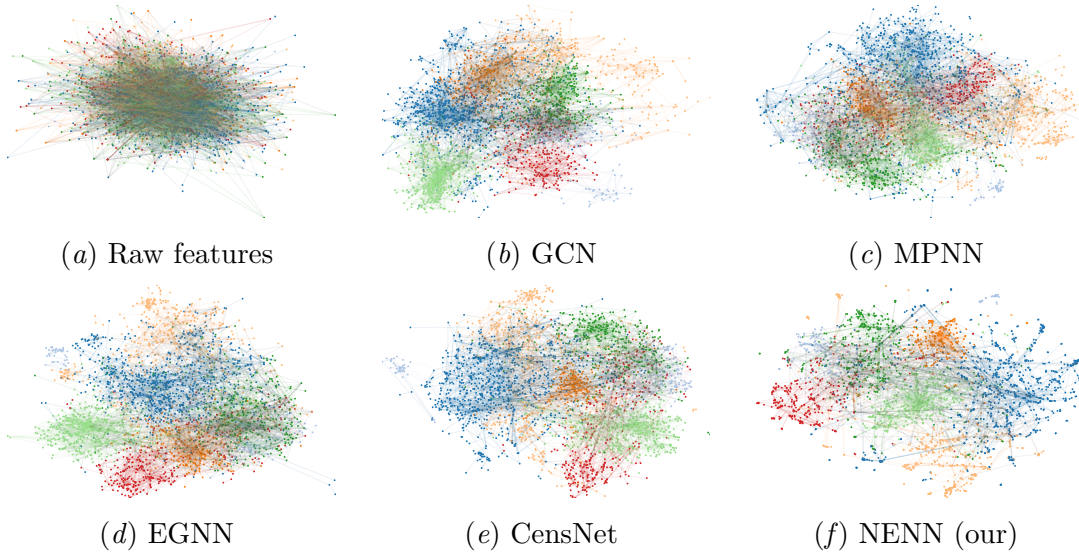
Figure 2: t-SNE visualization of semi-supervised node classification on Cora dataset.

- **EGNN** Gong and Cheng (2019): A graph neural network based on attention mechanism. EGNN explores edge features, but in the later layer, the edge feature vectors are converted to attention coefficients, which leads to the loss of edge information.

- **NENN-NAEA**: A variant of NENN. NENN-NAEA represents that both node and edge features are aggregated by attention mechanism.

- **NENN-NCEC**: A variant of NENN. NENN-NCEC represents that both node and edge features are aggregated by convolution operation.

- **NENN-NCEA**: A variant of NENN. NENN-NCEA represents that node features are aggregated by convolution operation and edge features are aggregated by attention mechanism.

- **NENN-NAEC**: A variant of NENN. NENN-NAEC represents that node features are aggregated by attention mechanism and edge features are aggregated by convolution operation.

### 4.3. Experimental Setup

Our experiments are all implemented by TensorFlow Abadi et al. (2016) and run on Ubuntu Linux 16.04 with NVIDIA RTX 2080 Ti. We use the Adam Kingma and Ba (2014) algorithm for training the models with the learning rate 0.0001 and batch size 128, number of epochs 200. The window size of an early stopping strategy is 200. We implement three layers NENN (i.e. node-level attention layer, edge-level attention layer and node-level attention layer). For all baselines, we split exactly the same training set, validation set and test set to ensure fairness. We follow the same splitting strategy in Jiang et al. (2019) and conduct experiments on citation networks with different label rate. For all molecular networks,
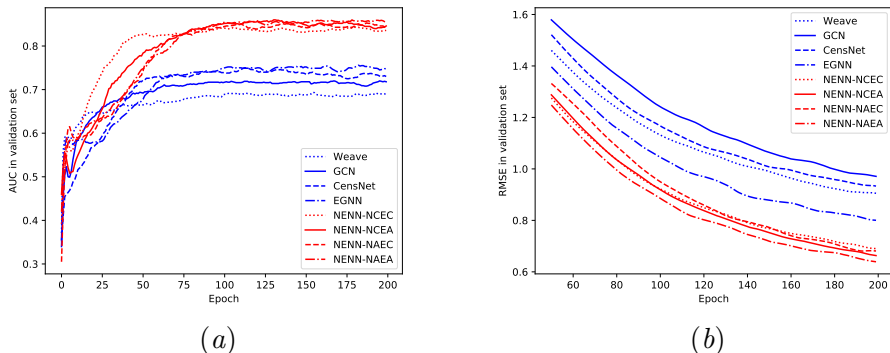
Figure 3: AUC in validation set for HIV networks (left) and RMSE in validation set for Lipophilicity networks (right).

training, validation and test dataset are split into with a ratio of 8:1:1. We run our models 5 times and report mean performance for each experiment.

### 4.4. Semi-supervised Node Classification

For semi-supervised node classification, we evaluate and report classification accuracies of the proposed NENN. We compared the proposed NENN with the representative models of the following five methods: GCN, GraphSAGE, GAT, CensNet, EGNN.

Table 3 reports classification accuracies of five baselines and the proposed NENN on citation networks. In all cases, the proposed NENN performs consistently much better than all baselines in 7 out of 9 experiments. From t-SNE Maaten and Hinton (2008) visualization in Figure 2, we can find that the proposed NENN can achieve more separated clusters than GCN, MPNN, EGNN, CensNet. It demonstrates that via incorporating node and edge features, the proposed NENN can learn a more meaningful node embedding.

### 4.5. Graph Classification

For graph classification, we predict molecular activity on the Tox21 and HIV datasets. We report the Area Under Curve (AUC) Hanley and McNeil (1982) scores and compare with some state-of-art baselines, including Random Forest (RF) Ali et al. (2012), Weave Kearnes et al. (2016), GCN, CensNet, EGNN.

We also compare some variants of our model NENN to validate the effectiveness. Table 4 reports the performance of all baselines and our models on four molecular networks. Remarkably, we can find that the embedding methods based on graph neural network have higher performance than the traditional methods based on graph kernel, which proves the superiority of graph representation learning methods. It is clear to observe that five baselines have rooms to improve on Tox21 and HIV datasets compared with our NENN models. Specifically, NENN improves upon the state-of-the-art GCN by a margin of 15.0% for the validation and 7.0% for the test on HIV dataset. Morever, since the feature dimensions of molecular networks are higher than those of citation networks, NENN-NCEC can improve

604

| Evaluation | AUC (Classification) | | | | RMSE(Regression) | | | |
|---|---|---|---|---|---|---|---|---|
| Dataset | Tox21 | | HIV | | Lipophilicity | | Freesolv | |
| Data Split | Validation | Test | Validation | Test | Validation | Test | Validation | Test |
| RF | $0.78 \pm 0.01$ | $0.75 \pm 0.03$ | $0.63 \pm 0.02$ | $0.62 \pm 0.02$ | $0.67 \pm 0.02$ | $0.66 \pm 0.04$ | $1.98 \pm 0.07$ | $1.62 \pm 0.14$ |
| Weave | $0.79 \pm 0.02$ | $0.80 \pm 0.02$ | $0.68 \pm 0.03$ | $0.71 \pm 0.05$ | $0.88 \pm 0.06$ | $0.89 \pm 0.04$ | $1.35 \pm 0.22$ | $1.37 \pm 0.14$ |
| GCN | $0.82 \pm 0.02$ | $0.84 \pm 0.01$ | $0.70 \pm 0.05$ | $0.77 \pm 0.02$ | $0.96 \pm 0.05$ | $0.98 \pm 0.03$ | $1.30 \pm 0.09$ | $1.35 \pm 0.26$ |
| EGNN | $0.82 \pm 0.01$ | $0.82 \pm 0.01$ | $0.73 \pm 0.06$ | $0.71 \pm 0.05$ | $0.79 \pm 0.02$ | $0.75 \pm 0.01$ | $1.07 \pm 0.08$ | $1.01 \pm 0.12$ |
| CensNet | $0.78 \pm 0.00$ | $0.79 \pm 0.00$ | $0.74 \pm 0.01$ | $0.73 \pm 0.02$ | $0.94 \pm 0.02$ | $0.83 \pm 0.02$ | $1.22 \pm 0.02$ | $1.46 \pm 0.01$ |
| NENN-NAEA | $\mathbf{0.86 \pm 0.02}$ | $0.85 \pm 0.01$ | $0.84 \pm 0.05$ | $\mathbf{0.84 \pm 0.01}$ | $0.67 \pm 0.06$ | $\mathbf{0.67 \pm 0.03}$ | $\mathbf{1.02 \pm 0.04}$ | $\mathbf{1.01 \pm 0.01}$ |
| NENN-NCEC | $0.81 \pm 0.02$ | $0.80 \pm 0.02$ | $0.80 \pm 0.01$ | $0.81 \pm 0.01$ | $0.77 \pm 0.04$ | $0.76 \pm 0.05$ | $1.25 \pm 0.04$ | $1.22 \pm 0.02$ |
| NENN-NCEA | $0.85 \pm 0.02$ | $0.84 \pm 0.01$ | $\mathbf{0.85 \pm 0.02}$ | $0.83 \pm 0.01$ | $\mathbf{0.66 \pm 0.02}$ | $0.70 \pm 0.08$ | $1.13 \pm 0.08$ | $1.09 \pm 0.04$ |
| NENN-NAEC | $0.85 \pm 0.01$ | $\mathbf{0.86 \pm 0.02}$ | $0.84 \pm 0.01$ | $0.83 \pm 0.01$ | $0.69 \pm 0.02$ | $0.71 \pm 0.01$ | $1.14 \pm 0.06$ | $1.11 \pm 0.06$ |

Table 4: Prediction results for the four molecular networks.

more in graph classification than in node classification. This also proves the validity of NENN that integrates higher-dimensional edge features into representation learning produce.

## 4.6. Graph Regression

For graph regression, we predict the solvation energies or solubility on Lipophilicity and Freesolv datasets. Root mean square error (RMSE) Chai and Draxler (2014) are adopted as the evaluation metric. The baselines of graph regression are the same as graph classification.

Table 4 shows the experimental results of RMSE. We highlight the best performance for all molecular networks. It's obvious that all of the variants of NENN significantly outperform the compared methods in RMSE. Figure 3(a) shows that the NENN obtains the best AUC score in validation set for HIV networks after around 50 epochs compared with other baselines. Figure 3(b) shows that the NENN greatly reduces RMSE. From Table 4, we can find that the performance of NENN-NAEA is much better than that of NENN-NCEC in both graph classification and graph regression, which proves the dual-level attention mechanism considering both the importance of neighbors based on nodes and edges can learn more efficient representations.

## 5. Conclusion

In this paper, we introduce an efficient embedding architecture, named NENN, which incorporates node and features to enhance the node and edge embeddings across neural network layer. The proposed NENN alternately stacks node-level attention layer and edge-level attention layer to learn the importance of node based neighbors and edge based neighbors. Leveraging the proposed NENN, the node and edge embeddings can be mutually reinforced. Extensive experiments on semi-supervised node classification, graph classification and graph regression demonstrate the effectiveness of NENN.

# References

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensor-flow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.

Jehad Ali, Rehanullah Khan, Nasir Ahmad, and Imran Maqsood. Random forests and decision trees. *International Journal of Computer Science Issues (IJCSI)*, 9(5):272, 2012.

Rianne van den Berg, Thomas N Kipf, and Max Welling. Graph convolutional matrix completion. *arXiv preprint arXiv:1706.02263*, 2017.

Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.

Shaosheng Cao, Wei Lu, and Qiongkai Xu. Grarep: Learning graph representations with global structural information. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 891–900. ACM, 2015.

Tianfeng Chai and Roland R Draxler. Root mean square error (rmse) or mean absolute error (mae)?–arguments against avoiding rmse in the literature. *Geoscientific model development*, 7(3):1247–1250, 2014.

Haochen Chen, Bryan Perozzi, Yifan Hu, and Steven Skiena. Harp: Hierarchical representation learning for networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

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.

Claire Donnat, Marinka Zitnik, David Hallac, and Jure Leskovec. Learning structural node embeddings via diffusion wavelets. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1320–1329, 2018.

Lise Getoor Galileo Mark Namata, Ben London and Bert Huang. Query-driven active survey-ing for collective classification. 2012.

Liyu Gong and Qiang Cheng. Exploiting edge features for graph neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9211–9219, 2019.

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.

William L. Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *NIPS*, 2017.

James A Hanley and Barbara J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 1982.

Frank Harary and Robert Z Norman. Some properties of line digraphs. *Rendiconti del Circolo Matematico di Palermo*, 9(2):161–168, 1960.

Arman Hasanzadeh, Ehsan Hajiramezanali, Krishna Narayanan, Nick Duffield, Mingyuan Zhou, and Xiaoning Qian. Semi-implicit graph variational auto-encoders. In *Advances in Neural Information Processing Systems*, pages 10711–10722, 2019.

Anfeng He, Chong Luo, Xinmei Tian, and Wenjun Zeng. A twofold siamese network for real-time object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4834–4843, 2018.

Yann Jacob, Ludovic Denoyer, and Patrick Gallinari. Learning latent representations of nodes for classifying in heterogeneous social networks. In *Proceedings of the 7th ACM international conference on Web search and data mining*, pages 373–382, 2014.

Simon Jégou, Michal Drozdzal, David Vazquez, Adriana Romero, and Yoshua Bengio. The one hundred layers tiramisu: Fully convolutional densenets for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 11–19, 2017.

Xiaodong Jiang, Pengsheng Ji, and Sheng Li. Censnet: convolution with edge-node switching in graph neural networks. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 2656–2662. AAAI Press, 2019.

Patrick F. Riley Oriol Vinyals George E. Dahl Justin Gilmer, Samuel S. Schoenholz. Neural message passing for quantum chemistry. In *International Conference on Machine Learning (ICML)*, 2017.

Steven Kearnes, Kevin McCloskey, Marc Berndl, Vijay Pande, and Patrick Riley. Molecular graph convolutions: moving beyond fingerprints. *Journal of computer-aided molecular design*, 30(8):595–608, 2016.

Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.

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.

Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. *arXiv preprint arXiv:1904.08082*, 2019.

Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

David L Mobley and J Peter Guthrie. Freesolv: a database of experimental and calculated hydration free energies, with input files. *Journal of computer-aided molecular design*, 28 (7):711–720, 2014.

Ryan L Murphy, Balasubramaniam Srinivasan, Vinayak Rao, and Bruno Ribeiro. Relational pooling for graph representations. *arXiv preprint arXiv:1903.02541*, 2019.

Mingdong Ou, Peng Cui, Jian Pei, Ziwei Zhang, and Wenwu Zhu. Asymmetric transitivity preserving graph embedding. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1105–1114. ACM, 2016.

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.

Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.

Leonardo FR Ribeiro, Pedro HP Saverese, and Daniel R Figueiredo. struc2vec: Learning node representations from structural identity. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 385–394, 2017.

Michael Schlichtkrull, Thomas N. Kipf, Peter Bloem, Rianne Vanden Berg, and Max Welling. Modeling relational data with graph convolutional networks. In *European Semantic Web Conference*, 2018.

Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi-Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.

Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=rJXMpikCZ. accepted as poster.

Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.

Zhitao Ying, Jiaxuan You, Christopher Morris, Xiang Ren, Will Hamilton, and Jure Leskovec. Hierarchical graph representation learning with differentiable pooling. In *Advances in neural information processing systems*, pages 4800–4810, 2018.

Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. In *Advances in Neural Information Processing Systems*, pages 5165–5175, 2018.