

Genetic-Evolutionary Graph Neural Networks: A Paradigm for Improved Graph Representation Learning

Anonymous authors

Paper under double-blind review

Abstract

Message-passing graph neural networks have become the dominant framework for learning over graphs. However, empirical studies continually show that message-passing graph neural networks tend to generate over-smoothed representations for nodes after iteratively applying message passing. This over-smoothing problem is a core issue that limits the representational capacity of message-passing graph neural networks. We argue that the fundamental problem with over-smoothing is a lack of diversity in the generated embeddings, and the problem could be reduced by enhancing the embedding diversity in the embedding generation process. To this end, we propose genetic-evolutionary graph neural networks, a new paradigm for graph representation learning inspired by genetic algorithms. We view each layer of a graph neural network as an evolutionary process and develop operations based on crossover and mutation to prevent embeddings from becoming similar to one another, thus enabling the model to generate improved graph representations. The proposed framework is well-motivated, as it directly draws inspiration from genetic algorithms for preserving population diversity. We experimentally validate the proposed framework on six benchmark datasets on different tasks. The results show that our method significantly advances the performance of current graph neural networks, resulting in new state-of-the-art results for graph representation learning on these datasets.

1 Introduction

Graphs are a general data structure for representing and analyzing complex relationships among entities. Many real-world systems, such as social networks, molecular structures, communication networks, can be modeled using graphs. It is essential to develop intelligent models for uncovering the underlying patterns and interactions within these graph-structured systems. Recent years have seen an enormous body of studies on learning over graphs. The studies include graph foundation models, geometry processing and deep graph embedding. These advances have produced new state-of-the-art or human-level results in various domains, including recommender systems (Zhang et al., 2024), chemical synthesis (Xie et al., 2024), and 2D and 3D vision tasks (Chen et al., 2024; Kim et al., 2023).

Graph neural networks have emerged as a dominant framework for learning from graph-structured data. The development of graph neural network models can be motivated from different approaches. The fundamental graph neural networks have been derived as a generalization of convolutions to non-Euclidean data (Bruna et al., 2014), as well as by analogy to classic graph isomorphism tests (Hamilton et al., 2017). Regardless of the motivations, the defining feature of the graph neural network framework is that it utilizes a form of message passing wherein messages are exchanged between nodes and updated using neural networks (Hamilton, 2020). During each graph neural network layer, the model aggregates features from a node’s local neighbourhood and then updates the node’s representation according to the aggregated information.

Message passing is at the heart of current graph neural networks. However, this paradigm of message passing also has major limitations. Theoretically, it is connected to the Weisfeiler-Lehman (WL) isomorphism test as well as to simple graph convolutions. The representational capacity of message-passing graph neural networks is inherently bounded by the WL isomorphism test (Xu et al., 2019). Empirical studies continually

find that message-passing graph neural networks suffer from the problem of over-smoothing (Wu et al., 2023; Hamilton, 2020). That is, the representations for all nodes can become very similar to one another after too many message passing iterations. These core limitations prevent graph neural networks from learning more meaningful representations from graphs. In recent years, an increasing number of studies have been devoted to addressing the bottlenecks, such as normalization and regularization techniques (Zhao & Akoglu, 2020b), and combining the global self-attention mechanism (Rampasek et al., 2022), exploring generalized message passing (Barceló et al., 2020). Despite these advances, improving the capability of graph neural network models still remains a fundamental challenge in learning from graph-structured data.

To learn meaningful graph representations, it is crucial to generate embeddings for all nodes that depend on both the graph structure and node attributes. However, when the over-smoothing phenomenon occurs, the representations for all nodes begin to look identical to each other. The consequence is that node-specific feature information becomes lost. To prevent this issue, it is important to preserve the diversity of generated embeddings throughout their layer-wise generation process. *In this paper, we propose genetic-evolutionary graph neural networks, a new paradigm for graph representation learning that integrates the idea from genetic algorithms for maintaining population diversity into the message-passing graph neural network framework.*

Genetic algorithms, inspired by Charles Darwin’s theory of natural evolution, emulate the process of natural selection, wherein the fittest individuals are selected to reproduce and generate the next generation of offspring. Over successive generations, biological organisms evolve using biologically inspired operators, such as selection, crossover and mutation, towards “survival of the fittest”. This approach has become a powerful tool for generating high-quality solutions for complex optimization and search problems.

In genetic algorithms, crossover and mutation play a key role in generating diverse individuals for selection, preventing the algorithms from premature convergence (Gupta & Ghafir, 2012). Crossover introduces variety by combining genetic information from different parents, and mutation introduces small random changes in genetic information. In this work, we view the iterative node embedding process as an evolutionary process, in which each layer of message passing produces a new generation of embeddings. We introduce two crossover operations, i.e., cross-generation crossover and sibling crossover, and a mutation operation, and we develop two graph neural network building blocks based on the operations. At each layer of a graph neural network, we first use message passing to update node representations and then apply crossover and mutation to prevent embeddings from becoming similar to one another, thus enabling the model to learn improved graph representations.

Unlike previous methods, such as residual connections (He et al., 2016), SSFG (Zhang et al., 2022) and PairNorm (Zhao & Akoglu, 2020a), this work proposes operations by drawing inspiration from genetic algorithms for addressing the over-smoothing problem in graph neural networks. Our framework is well-motivated as it views the layer-wise node embedding process as analogous to the genetic evolutionary process. It is a general paradigm that can be integrated into different graph neural network models. We conduct experiments on six benchmark datasets on different graph tasks. We show that the use of our framework significantly improves the performance of the baseline graph neural networks, advancing the state-of-the-art results for graph representation learning on the datasets.

The main contributions of this paper can be summarized as follows. (1) This paper proposes a new framework named genetic-evolutionary graph neural networks for learning from graph-structured data. The core idea behind the proposed framework is to model each layer of a graph neural network as an evolutionary process. We develop three key operations inspired by crossover and mutation from genetic algorithms to enhance the diversity of generated embeddings at each layer. (2) The proposed framework is well-motivated, as it is directly inspired by genetics. It is a general paradigm which can be integrated into current message-passing graph neural networks. Empirical evaluations are conducted on six popular datasets on different graph tasks, and the results demonstrate that the proposed framework significantly improves the performance of the baseline graph neural networks.

2 Related Work

2.1 Graph Neural Networks

Most current graph neural networks can be categorized into spectral approaches and spatial approaches (Velićković et al., 2018). The spectral approaches are developed based on spectral graph theory. The key idea of spectral graph neural networks is that convolutions are defined in the spectral domain through an extension of the Fourier transform to graphs. In contrast, spatial graph neural networks define convolutions in spatially localized neighbourhoods. The behaviour of the convolutions is analogous to that of kernels in convolutional neural networks which aggregate features from spatially-defined patches in an image.

Gilmer et al. (2017) reformulated common models as a message-passing neural network framework, wherein node features are exchanged between nodes and updated using neural networks. A common issue with message-passing graph neural networks is known as the over-smoothing problem. This issue of over-smoothing was first identified by Li et al. (2018). It can also be viewed as a consequence of the neighbourhood aggregation operation in the message-passing update (Hamilton, 2020). The follow-up studies for limiting over-smoothing include graph normalization and regularization techniques (Zhao & Akoglu, 2020a; Chen et al., 2022), combining the global self-attention with local message passing (Rampasek et al., 2022), and improved graph attention approaches (Wu et al., 2024). Additionally, there have been studies on uncovering over-smoothing in basic graph neural network models from theoretical analysis (Oono & Suzuki, 2020). Luan et al. (2024) investigated the impact of the homophily principle, i.e., nodes with the same labels are more likely to be connected, on node distinguishability and showed that graph neural network is capable of generating meaningful representations regardless of homophily levels. Kelesis et al. (2023) introduced a simple modification to the slope of the ReLU activation function to mitigate the issue of oversmoothing. Rong et al. (2020) introduces a method for reducing oversmoothing by randomly removing edges during training for node classification tasks. Kelesis et al. (2025) proposed a weight initialization method that takes into account the graph topology for reducing oversmoothing in graph neural networks.

2.2 Genetic Algorithms

Genetic algorithms, first introduced by John Holland in the 1960s (Holland, 1992), are a heuristic optimization method inspired by the principles of nature selection and evolutionary biology. Genetic algorithms are a powerful method for solving complex constrained and unconstrained optimization problems. In the past decades, genetic algorithms have emerged as a widely adopted tool for addressing optimization and search problems across various fields including scheduling, mathematics, robotics and bioinformatics (Alhijawi & Awajan, 2023; Davidor, 1991; McCall, 2005).

In machine learning, genetic algorithms have been successfully applied for feature selection (Babatunde et al., 2014), neural network optimization (Miller et al., 1989) and hyperparameter tuning for models like neural networks and support vector machines (Alibrahim & Ludwig, 2021). In object detection, hyperparameter evolution which uses a genetic algorithm was applied for optimizing hyperparameters in YOLO models (Redmon, 2016). Sehgal et al. (2019) showed that evolving the weights of a deep neural network using a genetic algorithm was a competitive approach for training reinforcement learning models. Shi et al. (2022) used the genetic approach for graph neural network architecture search. Unlike previous work, to the best of our knowledge, this is the first work that introduces operations inspired by genetic algorithms for regularizing graph neural networks through enhancing embedding diversity.

3 Methodology

3.1 Graph Neural Networks

A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ can be defined through a set of nodes \mathcal{V} and a set of edges \mathcal{E} between pairs of these nodes. Each node $u \in \mathcal{V}$ is associated with a node-level feature \mathbf{x}_u . Graph neural networks are a general framework for representation learning over the graph \mathcal{G} and $\{\mathbf{x}_u, \forall u \in \mathcal{V}\}$. At its core, the graph neural network framework iteratively updates the representation for every node using a form of message passing.

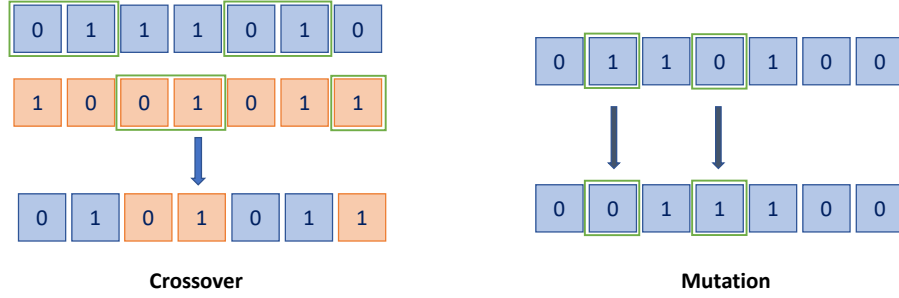


Figure 1: Crossover recombines of the genetic information of parents to produce an offspring. Mutation introduces small random changes in genetic information.

During each message-passing iteration, each node $u \in \mathcal{V}$ aggregates the representations of the nodes in its neighborhood, and the representation for node u is then updated according to the aggregated representation. Following Hamilton (2020), this message-passing framework can be expressed as follows:

$$\mathbf{h}_u^{(k)} = \text{Update}^{(k)} \left(\mathbf{h}_u^{(k-1)}, \text{Aggregate}^{(k)}(\{\mathbf{h}_v^{(k-1)}, \forall v \in \mathcal{N}(u)\}) \right), \quad (1)$$

where *Update* and *Aggregate* are neural networks, and $\mathcal{N}(u)$ is the set nodes in u 's neighbourhood. The superscripts are used for distinguishing the embeddings and functions at different iterations. At each iteration k , the *Aggregate* function takes the set of embeddings of nodes in $\mathcal{N}(u)$ as input and generates an aggregated message $\mathbf{m}_{\mathcal{N}(u)}^{(k)}$. The *Update* function then generates the updated embedding for node u based on the message $\mathbf{m}_{\mathcal{N}(u)}^{(k)}$ and u 's previous embedding $\mathbf{h}_u^{(k-1)}$. The embeddings at $k = 0$ are initialized to the node-level features, i.e., $\mathbf{h}_u^{(0)} = \mathbf{x}_u, \forall u \in V$. After K iterations of message passing, every node embedding contains information from its K -hop neighborhood.

This message passing formalism is currently the dominant framework for learning over graphs. However, a common issue with message-passing graph neural networks is over-smoothing. The idea of over-smoothing is that the embeddings for all nodes begin to become similar and are relatively uninformative after too many rounds of message passing. This issue of over-smoothing can be viewed as a consequence of the neighborhood aggregation operation. Li et al. (2018) showed that the graph convolution of the basic graph convolutional network model (Kipf & Welling, 2016) can be seen as a special form of Laplacian smoothing that generates the representation for every node using the weighted average of a node's itself and its neighbours' embeddings. But after applying too many rounds of Laplacian smoothing, the representations for all nodes will become indistinguishable from each other. From the graph signal processing perspective, multiplying a signal by high powers of the symmetric normalized adjacency matrix $\mathbf{A}_{\text{sym}} = \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{D}^{\frac{1}{2}}$, which corresponds to a convolutional filter the lowest eigenvalues, or frequencies, of the symmetric normalized Laplacian $\mathbf{L}_{\text{sym}} = \mathbf{I} - \mathbf{A}_{\text{sym}}$. Thus, the simple graph neural network that stacks multiple rounds of graph convolution converges all the node representations to constant values within connected components on the graph, i.e., the "zero-frequency" of the Laplacian (Hamilton, 2020).

3.2 Genetic-Revolutionary Graph Neural Networks

3.2.1 Motivation

In the above, we discussed the over-smoothing problem in message-passing graph neural network. We see that the fundamental issue is the loss of embedding diversity at each layer throughout the generation process. Thus, we can view the trade-off between model performance and depth of popular graph neural network models from this perspective. Graph neural networks need to model complex relationships and long-term dependencies using more layers to improve the performance. However, using too many layers will eliminate node-specific features, which leads to significantly reduced model performance.

Graph neural networks generate embeddings for nodes through an iterative message-passing process. At each message-passing iteration, the representation for every node is updated according to the information aggregated from the node’s graph neighbourhood. We can view this iterative process as an genetic evolutionary process, wherein graph nodes are individuals of a population, and the model is to evolve a population of nodes over multiple generations to obtain their expressive representations for graph tasks.

In genetic algorithms, a very homogeneous population, i.e., little population diversity, is considered as the major reason for premature convergence to suboptimal solutions (Whitley, 2001). Therefore, it is crucial to preserve the diversity of population during the evolutionary process. Similarly, we need to maintain the diversity of generated embedding in their generation to prevent the model from converging to a local optimum in optimization.

To preserve the population diversity, genetic algorithms use the operators of crossover and mutation to generate diverse individuals and select those best fit the environment to evolve over successive generations. The crossover operation recombines of the characteristics of each ancestor of an offspring, and the mutation operation randomly changes the genetic information to increase the variability (see Figure 1). In a similar manner, we can generalize the mechanisms to the embedding generation process. By integrating crossover and mutation methods within the message-passing framework, we can prevent generated embeddings from becoming too similar to each other. This ultimately would enhance the model representational capacity.

3.2.2 Improving Graph Neural Networks with Genetic Operations

We view each layer of a graph neural network as a genetic evolution process, in which the nodes represent individuals of a population and their embeddings represent chromosomes that store genetic information. During each graph neural network layer, we first use message passing to update the embeddings for all nodes and then use genetic operations to increase the diversity of generated embeddings. We propose three operations inspired by genetic algorithms: (1) cross-generation crossover, (2) sibling crossover, and (3) mutation.

Genetically, crossover is a process in which the genetic information of two parents is recombined to produce new offspring, resulting in the exchange of genetic material between parental chromosomes. This mechanism forms the basis for driving biological variation, shaping differences in traits within species and introducing novel traits previously unseen in a population. It basically helps promote the evolutionary process by enabling novel gene combinations to emerge and spread across generations. Fundamentally, this process creates diversity at the level of genes that reflects difference in chromosomes of different individuals.

Cross-generation crossover. Similar to crossover in genetics, the cross-generation operation in our framework recombines the embedding for a node generated by message-passing and the node’s previous layer embedding. For $\bar{\mathbf{h}}_u^{(k)} = (\bar{\mathbf{h}}_{u,1}^{(k)}, \dots, \bar{\mathbf{h}}_{u,d}^{(k)})$ and $\mathbf{h}_u^{(k-1)} = (\mathbf{h}_{u,1}^{(k-1)}, \dots, \mathbf{h}_{u,d}^{(k-1)})$ which represent the embedding for node u generated by message passing and u ’s previous layer embedding, cross-generation crossover can be expressed as follows:

$$\begin{aligned} \mathbf{h}_u^{(k)} &= \text{Crossover}(\bar{\mathbf{h}}_u^{(k)}, \mathbf{h}_u^{(k-1)}) \\ \text{where } \mathbf{h}_{u,i}^{(k)} &= \begin{cases} \mathbf{h}_{u,i}^{(k-1)} & \text{if } \lambda_i < p \\ \bar{\mathbf{h}}_{u,i}^{(k)} & \text{else} \end{cases}, i = 1, \dots, d, \end{aligned} \quad (2)$$

and $\lambda_i \sim U(0, 1)$ and p is a probability indicating information from the previous layer embedding. At each dimension, the feature is randomly selected from the embedding generated using message passing or from the embedding inputted to this layer. Because each round of message passing generates a smoothed version of the input, recombining information from a node’s previous layer embedding reduces the smoothness of the generated embeddings. This operation is a parameter-free method and can be integrated into current graph neural networks.

Sibling crossover is an operation that randomly selects information from siblings. In our implementation, we generate multi-head outputs using message passing as siblings and update the embedding for a node by randomly selecting information from the multi-head outputs.

Algorithm 1 Pseudocode for cross-generation crossover in a PyTorch-like style.

```

# h, h_in: representaton generated by message passing and the previous layer embedding
# f_prob: probability of recombining information from parent
# self.dist: a Bernoulli distribution defined by torch.distributions.Bernoulli(torch.tensor(self.f_prob)):

def forward(self, h, h_in):
    if self.training == True:
        crossover_mask = self.dist.sample(h.shape) # generate crossover mask

        # crossover from h and h_in
        h = h_in * crossover_mask + h * (1 - crossover_mask)
    else:
        h = h_in * self.f_prob + h * (1 - self.f_prob)

    return h

```

Algorithm 2 Pseudocode for mutation in a PyTorch-like style.

```

# self.running_mean: the mean of h over the training set
# self.running_var: the variance of h over the training set
# self.mutation_prob: probability of mutation

def forward(self, h):
    if self.training == True:
        mean = h.mean([0])
        var = h.var([0])
        n = h.numel() / h.size(1)

        with torch.no_grad():
            # momentum update of running_mean and running_var
            self.running_mean = self.momentum * mean + (1 - self.momentum) * self.running_mean
            self.running_var = self.momentum * var * n / (n - 1) + (1 - self.momentum) * self.running_var

        # generate mutatioin noise
        gaussian_noise = torch.randn(h.shape)

    if self.training == True:
        mutation_mask = Bernoulli.sample(h.shape) # generate mutation mask
        h = (gaussian_noise * self.running_var + self.running_mean) * mutation_mask + h * (1 - mutation_mask)
    else:
        h = self.running_mean * self.mutation_prob + h * (1 - self.mutation_prob)

    return h

```

$$\mathbf{h}_u^{(k)} = \text{Crossover}(\bar{\mathbf{h}}_u^{(k, head_1)}, \dots, \bar{\mathbf{h}}_u^{(k, head_z)})$$

$$\text{where } \mathbf{h}_{u,i}^{(k)} = \bar{\mathbf{h}}_{u,i}^{(k, h_{ij})},$$
(3)

$h_{ij} \sim \text{Categorical}(\frac{1}{z}, \dots, \frac{1}{z})$, and z is the number of heads. Each $\bar{\mathbf{h}}_u^{(k, head_h)}$ in the multi-head outputs represents a sibling generated using the same input. This operation also increases individual diversity by randomly combining information from different siblings.

Mutation is the process in which some genes of individuals are randomly changed. In our framework, the feature at each dimension is randomly replaced by a value sampled from a Gaussian distribution, wherein the statistics are calculated using batches. For a batch of m vectors $\mathcal{B} = \{\mathbf{h}_u^1, \mathbf{h}_u^2, \dots, \mathbf{h}_u^m\}$, we calculate the mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\delta}$ of the feature over the training set as follows.

$$\boldsymbol{\mu} \leftarrow \mathbb{E}_{\mathcal{B}}(\boldsymbol{\mu}_{\mathcal{B}})$$

$$\boldsymbol{\delta} \leftarrow \frac{m}{m-1} \mathbb{E}_{\mathcal{B}}(\boldsymbol{\delta}_{\mathcal{B}}^2)$$
(4)

where $\boldsymbol{\mu}_{\mathcal{B}}$ and $\boldsymbol{\delta}_{\mathcal{B}}^2$ are the mean and variance of the batch \mathcal{B} . Here we use the unbiased variance estimate. Then we randomly sample a vector $\boldsymbol{\gamma}$ from a multivariate Gaussian distribution $N(\mathbf{0}, \mathbf{I})$ and update the feature as follows:

$$\tilde{\mathbf{h}}_u^i = (\boldsymbol{\gamma} \boldsymbol{\delta} + \boldsymbol{\mu}) \text{mask} + \mathbf{h}_u^i (1 - \text{mask})$$
(5)

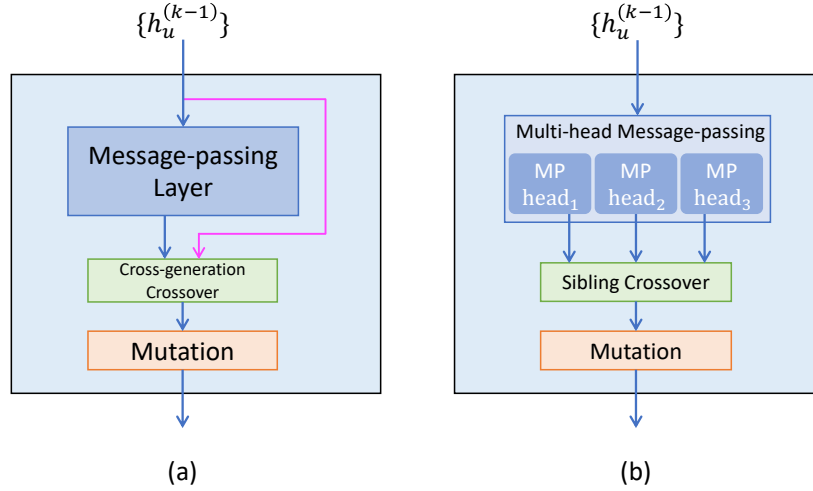


Figure 2: Building block architectures: Block (a) applies cross-generation to a node’s embedding generated using message passing and the node’s previous layer embedding, and Block (b) applies sibling crossover to a set of outputs generated using multi-head message passing.

where the **mask** $\sim \text{Bernoulli}(\text{mutation_rate})$. The mutation operation is also a parameter-free method. It basically introduces randomness to features as a regularization method, enabling the model to explore new space for optimization.

3.3 Model Architecture

Algorithm 1 and Algorithm 2 show our Pytorch-style pseudo-code for the cross-generation crossover operation and mutation operation respectively. The code for sibling crossover can be easily adapted from Algorithm 1. We design two building blocks based on the cross-generation crossover operation and sibling crossover operation (see Figure 2). The first building block applies the cross-generation crossover after message passing, followed by the mutation operation. Note that this building block is compatible with different graph neural network models and it does not introduce additional trainable parameters. The other building block applies sibling crossover to a set of multi-head outputs, followed by the mutation operation. This method requires the model to generate multiple siblings using a multi-head message passing.

The embedding generation process takes the graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and features for all nodes $\mathbf{x}_u, \forall u \in \mathcal{V}$, as input. This is followed by K building blocks that generate hidden embeddings. Finally, a readout function is applied to the output of the last block to generate the graph representation. For node-level tasks, the embeddings generated by the last block are directly used.

4 Empirical Evaluation

4.1 Datasets and Setup

The experiments are conducted on six benchmark datasets, i.e., MNIST, CIFAR10, PascalVOC-SP, COCO-SP, Peptides-func and Peptides-struct (Dwivedi et al., 2020; 2022) on three graph tasks, graph classification, node classification, and graph regression. We closely follow the setup as Dwivedi et al. (2020; 2022) for training and evaluating the models. The details of the datasets and evaluation metrics are provided in the appendix section.

Table 1: Classification accuracy (%) on MNIST and CIFAR10 on the superpixel graph classification task. The cross-generation crossover and mutation operations are applied to the base graph neural network model.

Model	MNIST	CIFAR10
MoNet (Monti et al., 2017)	90.805±0.032	54.655±0.518
GraphSAGE (Hamilton et al., 2017)	97.312±0.097	65.767±0.308
GIN (Xu et al., 2019)	96.485±0.252	55.255±1.527
GCNII (Chen et al., 2020)	90.667±0.143	56.081±0.198
PNA (Corso et al., 2020)	97.94±0.12	70.35±0.63
DGN (Beaini et al., 2021)	–	72.838±0.417
CRaWl (Toenshoff et al., 2021)	97.944±0.050	69.013±0.259
GIN-AK+ (Zhao et al., 2021)	–	72.19±0.13
3WLGNN (Maron et al., 2019)	95.075±0.961	59.175±1.593
EGT (Hussain et al., 2022)	98.173±0.087	68.702±0.409
GatedGCN + SSFG (Zhang et al., 2022)	97.985±0.032	71.938±0.190
EdgeGCN (Zhang et al., 2023)	98.432±0.059	76.127±0.402
Expformer (Shirzad et al., 2023)	98.550±0.039	74.754±0.194
TIGT (Choi et al., 2024)	98.230±0.133	73.955±0.360
RandAlign + GatedGCN (Zhang & Xu, 2024)	98.512±0.033	76.395±0.186
GCN (Kipf & Welling, 2016)	90.705±0.218	55.710±0.381
GCN + Ours	95.926±0.031	59.157±0.130
GAT (Veličković et al., 2018)	95.535±0.205	64.223±0.455
GAT + Ours	97.643±0.108	69.247±0.248
GatedGCN (Bresson & Laurent, 2017)	97.340±0.143	67.312±0.311
GatedGCN + Ours	98.526±0.031	78.242±0.239
GPS (Rampasek et al., 2022)	98.051±0.126	72.298±0.356
Finetuned GPS (Tönshoff et al., 2024)	98.186±0.107	75.680±0.188
Finetuned GPS + Ours	98.685±0.029	80.636±0.195

4.2 Results

CIFAR10 and MNIST. Table 1 reports the results on the two datasets on the superpixel classification task. We use the GPS (Rampasek et al., 2022) as the base model. The GPS model is a hybrid of local aggregation and global aggregation architecture. It uses GatedGCN for local aggregation and uses Transformer for global aggregation. We apply cross-generation and mutation (i.e., block (a) in Figure 2) to the base GatedGCN model. The crossover rate is set to 0.5 and mutation rate is set to 0.1. We see from Table 1 that our method improves the performance of the base Finetuned GPS by a large margin, with a relative improvement of 0.44% and 6.55% on MNIST and CIFAR10 respectively. It simultaneously outperforms both Expformer (Shirzad et al., 2023) and RandAlign (Zhang & Xu, 2024), which previously achieved the best performance on MNIST and CIFAR10 respectively.

PascalVOC-SP and COCO-SP. The two datasets are long-range prediction datasets compared to MNIST and CIFAR10. The task is to predict if a node corresponds to a region of an image which belongs to a particular class. We use Finetuned GPS (Tönshoff et al., 2023) as the base model. The Finetuned GPS is also a hybrid of GatedGCN and Transformer architecture. We apply cross-generation and mutation to the base GatedGCN model. The crossover rate is set to 0.9 and mutation rate is set to 0.05. The results are reported in Table 2. Previously, Finetuned GPS achieved the best performance among the baseline models on the two datasets. As compared to Finetuned GPS, the use of our method results in a relative improvement

Table 2: Results on PascalVOC-SP and COCO-SP on the node classification task. The cross-generation crossover and mutation operations are applied to the base graph neural network model.

Model	PascalVOC-SP (F1 \uparrow)	COCO-SP (F1 \uparrow)
GCN (Kipf & Welling, 2016)	0.1268 \pm 0.0060	0.0841 \pm 0.0010
GINE (Hu et al., 2019)	0.1265 \pm 0.0076	0.1339 \pm 0.0044
GCNII (Chen et al., 2020)	0.1698 \pm 0.0080	0.1404 \pm 0.0011
GatedGCN (Bresson & Laurent, 2017)	0.2873 \pm 0.0219	0.2641 \pm 0.0045
GatedGCN + RWSE (Rampasek et al., 2022)	0.2860 \pm 0.0085	0.2574 \pm 0.0034
Transformer + LapPE Dwivedi et al. (2022)	0.2694 \pm 0.0098	0.2618 \pm 0.0031
SAN + LapPE Dwivedi et al. (2022)	0.3230 \pm 0.0039	0.2592 \pm 0.0158
SAN + RWSE Dwivedi et al. (2022)	0.3216 \pm 0.0027	0.2434 \pm 0.0156
Expformer Shirzad et al. (2023)	0.3975 \pm 0.0037	0.3455 \pm 0.0009
RandAlign + GPS (Zhang & Xu, 2024)	0.4242 \pm 0.0011	0.3567 \pm 0.0026
Finetuned GCN (Tönshoff et al., 2024)	0.2078 \pm 0.0031	0.1338 \pm 0.0007
Finetuned GCN + Ours	0.2445\pm0.0008	0.1485\pm0.0004
Finetuned GatedGCN (Tönshoff et al., 2024)	0.3880 \pm 0.0040	0.2922 \pm 0.0018
Finetuned GatedGCN + Ours	0.4176\pm0.0013	0.3226\pm0.0010
GPS (Rampasek et al., 2022)	0.3748 \pm 0.0109	0.3412 \pm 0.0044
Finetuned GPS (Tönshoff et al., 2024)	0.4440 \pm 0.0065	0.3884 \pm 0.0055
Finetuned GPS + Ours	0.4832\pm0.0031	0.4002\pm0.0019

Table 3: Results on Pepti-func and Pepti-struct. The sibling crossover and mutation operations are applied to the base graph neural network model.

Model	Peptides-func (AP \uparrow)	Peptides-struct (MAE \downarrow)
GCN	0.5930 \pm 0.0023	0.3496 \pm 0.0013
GINE	0.5498 \pm 0.0079	0.3547 \pm 0.0045
GCNII (Chen et al., 2020)	0.5543 \pm 0.0078	–
GatedGCN	0.5864 \pm 0.0077	0.3420 \pm 0.0013
Gated + RWSE	0.6069 \pm 0.0035	0.3357 \pm 0.0006
Transformer+LapPE	0.6326 \pm 0.0126	0.2529 \pm 0.0016
SAN+LapPE	0.6384 \pm 0.0121	0.2683 \pm 0.0043
SAN+RWSE	0.6439 \pm 0.0075	0.2545 \pm 0.0012
Expformer (Shirzad et al., 2023)	0.6527 \pm 0.0043	0.2481 \pm 0.0007
GPS (Rampasek et al., 2022)	0.6535 \pm 0.0041	0.2500 \pm 0.0005
Finetuned GPS (Tönshoff et al., 2023)	0.6534 \pm 0.0091	0.2509 \pm 0.0014
Finetuned GPS + Ours	0.6730\pm0.0053	0.2486\pm0.0009
Finetuned GCN (Tönshoff et al., 2023)	0.6860 \pm 0.0050	0.2460 \pm 0.0007
Finetuned GCN + Ours	0.7021\pm0.0034	0.2426\pm0.0014

of 8.83% and 3.04% respectively without using additional model parameters. Once again, our framework achieves new state-of-the-art performance on the two datasets.

Peptides-func and Peptides-struct. We use Finetuned GCN (Tönshoff et al., 2023) as the base model on the two datasets. We use sibling crossover and mutation to the base model. The number of siblings is set to 2 and mutation rate is set to 0.1. The results are reported in Table 6. Finetuned GCN is a strong

Table 4: Ablation study: Importance of crossover and mutation on the model performance on CIFAR10 and PascalVOC-SP.

Base Model	Crossover	Mutation	CIFAR10	PascalVOC-SP
Finetuned GPS (Tönshoff et al., 2023)	×	×	75.680±0.188	0.4440±0.0065
	✓	×	79.434±0.228	0.4952±0.0098
	×	✓	77.029±0.203	0.4554±0.0077
	✓	✓	80.636±0.195	0.4832±0.0031

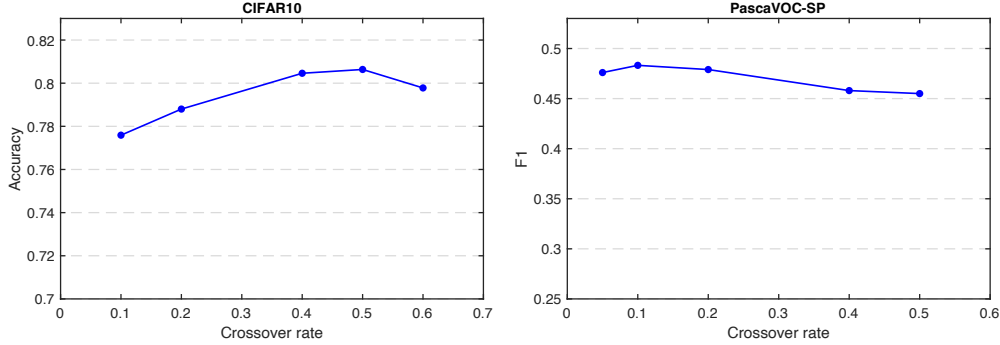
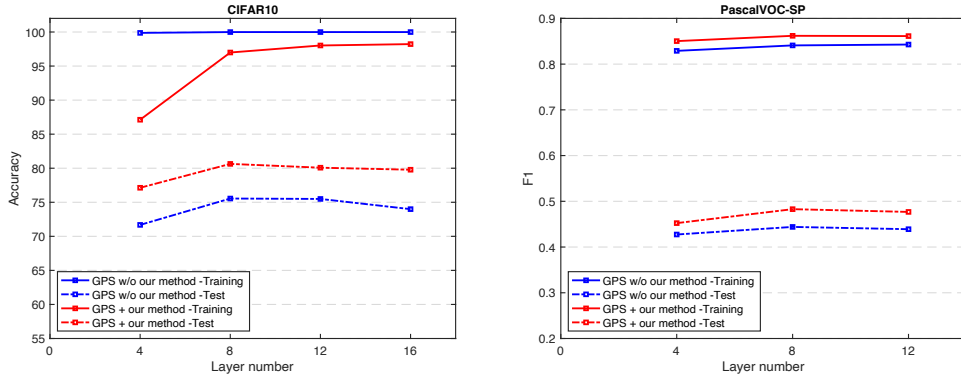
Figure 3: Impact of the crossover rate p on the performance of Finetuned GPS on CIFAR10 and PascalVOC-SP.

Figure 4: Results of our method on the base Finetuned GPS model with different layers on CIFAR10 and PascalVOC-SP.

baseline model in previous work. We see from Table 6 that the use of framework further improve the model performance.

Ablation Study. We conduct an ablation study on CIFAR10 and PascalVOC-SP to analyse the importance of crossover and mutation on the model performance. Table 5 shows the ablation study results. It can be seen from Table 5 that the crossover operation plays a major role in improving the model performance. The mutation operation helps further improve the model performance as a regularization method.

We further analyzed the impact of the crossover rate p on model performance on CIFAR10 and PascalVOC-SP. Figure 3 shows the experimental results. We see that the best performance is achieved when p is set to different values on the two datasets. When p is set to 0, it is equivalent to not using crossover. A recommended strategy for tuning p is starting from a small value, e.g., 0.05, and then gradually increasing it to find the optimal value. The strategy is also used for tuning the mutation rate in our evaluations.

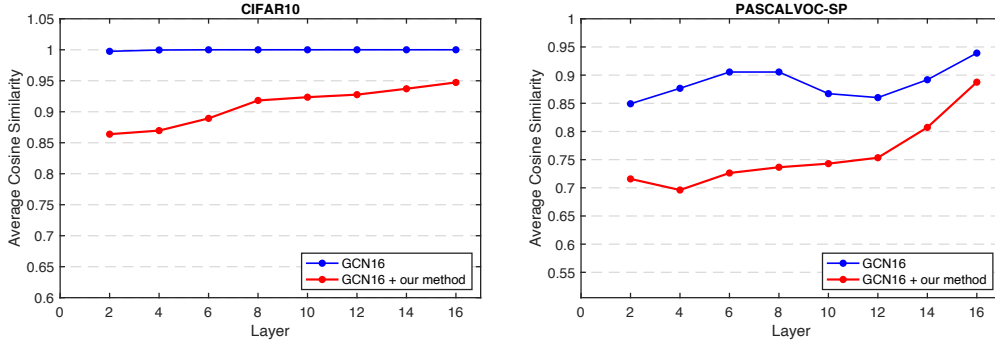


Figure 5: Average cosine similarity at different model layers on CIFAR10 and PascalVOC-SP.

We conducted experiments to analyze the performance of our method on the base Finetuned GPS model with different layers on CIFAR10 and PascalVOC-SP. The results are shown in Figure 4. We also analyzed the performance of our method on the base Finetuned GPS model with different layers on CIFAR10, and the results are reported in Figure 7 in the appendix section. It can be seen from Figure 4 and Figure 7 that the use of our method improves the model generalization performance on the base models.

We compared our method with the basic GCN in which residual connections and batch normalizations are not used on MNIST and CIFAR10. The results are shown in Table 5. We see that the model performance drops without using these techniques and that the joint use of our method with residual connections and batch normalizations yields the best task performance.

We validated our method for few-shot learning on the Omniglot dataset (Lake et al., 2015). The results are reported in Table 6. The base GNN (Satorras & Estrach, 2018) model achieves competitive results while remaining simpler than other methods. The use of our method further improves the performance of GNN (Satorras & Estrach, 2018) for 1-shot 20-way and 5-shot 20-way experiments.

Figure 5 visualizes the average cosine similarity at different layers of a GCN model on CIFAR10 and PascalVOC-SP. On the CIFAR10 dataset, we see that the average cosine similarity converges to 1 at an early stage for the base GCN model, which indicates that the oversmoothing problem occurs. With our method, it can be seen that the average cosine similarity is reduced. On the PascalVOC dataset, the average cosine similarity is also reduced when our method is applied to the base GCN model. Figure 6 shows the results of our method on vanilla GCN with increased layers on MNIST and CIFAR10. We see that the results for the vanilla GCN drop significantly on the two datasets. In contrast, our method enables the model to maintain stable performance, exhibiting only a slight drop with deeper layers. The results show that our method effectively prevents the model performance from a severe drop with increased layers.

Discussion. The idea of genetic operations is to enhance embedding diversity during training to prevent the model from premature convergence. The genetic operations fundamentally work as a regularization approach by introducing randomness. The crossover operation can be understood as a random pooling method, which randomly recombines features of different input embeddings according to a predefined probability value. The mutation operation can be seen as a generalized Dropout method that randomly replaces the feature point at a dimension with a value sampled from a Gaussian distribution, wherein the mean and variance of the Gaussian distribution are calculated using batch data. The Dropout method, in contrast, randomly replace the feature point at a dimension with the value of 0. The fundamental idea of genetic operations provides new perspective on machine learning tricks. It could be potentially applied for regularizing other machine learning models, such as multilayer perceptrons and Transformers.

In this work, the hyperparameters, i.e., crossover rate and mutation rate, are determined by empirical tuning through experiments. We found that of the crossover rate and the mutation rate were set to different values to achieve the best performance on different tasks. In future work, automatic methods could be explored for searching the optimal crossover and mutation rates. This would help save the extensive hyperparameter tuning procedure.

Table 5: Comparison of our method with the basic GCN, wherein residual connections and batch normalizations (BN) are not used.

Model	MNIST	CIFAR10
GCN (w/o residual connections and BN)	87.590 \pm 0.336	48.810 \pm 1.045
GCN (with residual connections and BN)	90.705 \pm 0.218	55.710 \pm 0.381
GCN (with residual connections and BN) + Ours	95.926 \pm 0.031	59.157 \pm 0.130

Table 6: Results on the base GNN model (Satorras & Estrach, 2018) for few-shot learning on Omniglot.

Model	20-Way	
	1-shot	5-shot
Siamese Net (Koch et al., 2015)	88.2%	97.0%
Matching Networks (Vinyals et al., 2016)	93.8%	98.5%
Prototypical Networks (Snell et al., 2017)	95.4%	98.8%
ConvNet with Memory (Kaiser et al., 2017)	95.0%	98.6%
Meta Networks (Munkhdalai & Yu, 2017)	97.0%	—
R2-D2 (Bertinetto et al., 2019)	96.24%	99.20%
GNN (Satorras & Estrach, 2018)	97.4%	99.0%
GNN + Ours	98.03%	99.47%

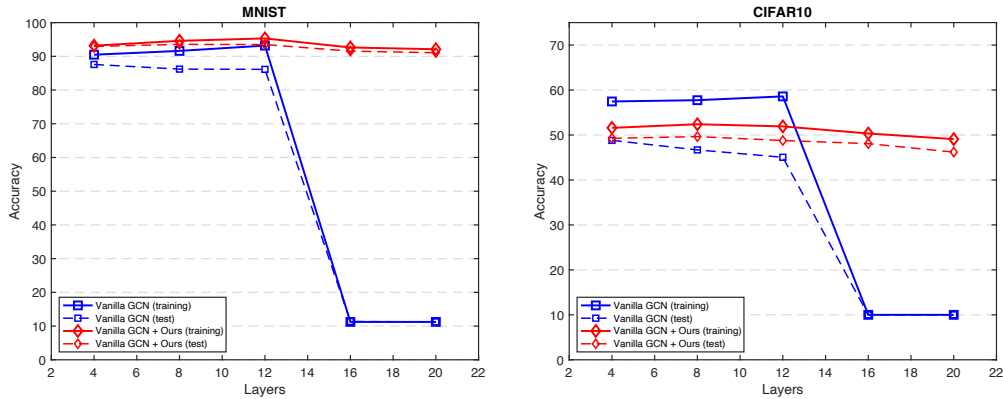


Figure 6: Performance of our method on the vanilla GCN with increased layers.

The proposed genetic operations introduce randomness in the training procedure. Therefore, the model requires more training epochs to optimize, which increases the training time. Additionally, the optimal combination of genetic operations and the hyperparameters need to be tuned experimentally for different tasks and different base models, which further increases the training time.

5 Conclusions

This paper presents a new framework called genetic-evolutionary graph neural networks for graph representation learning. The key idea of our approach is to view each layer of a graph neural network as a genetic evolutionary process and use biogenetics-inspired operations to prevent the over-smoothing problem in graph neural networks. We developed three operations, i.e., cross-generation crossover, sibling crossover and mutation, inspired by genetic algorithms and presented two building blocks based on the the operations for graph

representation learning. An important advantage of the proposed framework lies in its interpretability, as it frames layerwisely graph representation learning as an evolutionary process. The experimental evaluations were conducted on six popular datasets on different graph tasks. The results showed that the use of our framework significantly improves the performance of the base graph neural networks, achieving new state-of-the-art performance for graph representation learning on these datasets. We also presented ablations of our framework, showing the importance of each operation on the overall model performance.

References

- Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. Slic superpixels compared to state-of-the-art superpixel methods. *IEEE transactions on pattern analysis and machine intelligence*, 34(11):2274–2282, 2012.
- B Alhijawi and A Awajan. Genetic algorithms: theory, genetic operators, solutions, and applications, evol. intel., 2023, 2023.
- Hussain Alibrahim and Simone A Ludwig. Hyperparameter optimization: Comparing genetic algorithm against grid search and bayesian optimization. In *2021 IEEE Congress on Evolutionary Computation (CEC)*, pp. 1551–1559. IEEE, 2021.
- Oluleye H Babatunde, Leisa Armstrong, Jinsong Leng, and Dean Diepeveen. A genetic algorithm-based feature selection. 2014.
- Pablo Barceló, Egor V Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, and Juan-Pablo Silva. The logical expressiveness of graph neural networks. In *8th International Conference on Learning Representations (ICLR 2020)*, 2020.
- Dominique Beaini, Saro Passaro, Vincent Létourneau, Will Hamilton, Gabriele Corso, and Pietro Liò. Directional graph networks. In *International Conference on Machine Learning*, pp. 748–758. PMLR, 2021.
- Luca Bertinetto, Joao F Henriques, Philip HS Torr, and Andrea Vedaldi. Meta-learning with differentiable closed-form solvers. In *International conference on learning representations*, 2019.
- Xavier Bresson and Thomas Laurent. Residual gated graph convnets. *arXiv preprint arXiv:1711.07553*, 2017.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *International Conference on Learning Representations*, 2014.
- Chaoqi Chen, Yushuang Wu, Qiyuan Dai, Hong-Yu Zhou, Mutian Xu, Sibe Yang, Xiaoguang Han, and Yizhou Yu. A survey on graph neural networks and graph transformers in computer vision: A task-oriented perspective. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In *International Conference on Machine Learning*, pp. 1725–1735. PMLR, 2020.
- Yihao Chen, Xin Tang, Xianbiao Qi, Chun-Guang Li, and Rong Xiao. Learning graph normalization for graph neural networks. *Neurocomputing*, 493:613–625, 2022.
- Yun Young Choi, Sun Woo Park, Minhoo Lee, and Youngho Woo. Topology-informed graph transformer. *arXiv preprint arXiv:2402.02005*, 2024.
- Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*, 33:13260–13271, 2020.
- Yuval Davidor. *Genetic Algorithms and Robotics: A heuristic strategy for optimization*, volume 1. World Scientific Publishing Company, 1991.
- Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson. Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.

- Vijay Prakash Dwivedi, Ladislav Rampasek, Mikhail Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu, and Dominique Beaini. Long range graph benchmark. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pp. 1263–1272. PMLR, 2017.
- Deepti Gupta and Shabina Ghafir. An overview of methods maintaining diversity in genetic algorithms. *International journal of emerging technology and advanced engineering*, 2(5):56–60, 2012.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pp. 1024–1034, 2017.
- William L Hamilton. Graph representation learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 14(3):1–159, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- John H Holland. *Adaptation in natural and artificial systems: an introductory analysis with applications to biology, control, and artificial intelligence*. MIT press, 1992.
- Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure Leskovec. Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*, 2019.
- Md Shamim Hussain, Mohammed J Zaki, and Dharmashankar Subramanian. Global self-attention as a replacement for graph convolution. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 655–665, 2022.
- Łukasz Kaiser, Ofir Nachum, Aurko Roy, and Samy Bengio. Learning to remember rare events. *arXiv preprint arXiv:1703.03129*, 2017.
- Dimitrios Kelesis, Dimitrios Vogiatzis, Georgios Katsimpras, Dimitris Fotakis, and Georgios Paliouras. Reducing oversmoothing in graph neural networks by changing the activation function. In *ECAI 2023*, pp. 1231–1238. IOS Press, 2023.
- Dimitrios Kelesis, Dimitris Fotakis, and Georgios Paliouras. Reducing oversmoothing through informed weight initialization in graph neural networks. *Applied Intelligence*, 55(7):632, 2025.
- Sangwon Kim, Dasom Ahn, and Byoung Chul Ko. Cross-modal learning with 3d deformable attention for action recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 10265–10275, 2023.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- Gregory Koch, Richard Zemel, Ruslan Salakhutdinov, et al. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2, pp. 1–30. Lille, 2015.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Brenden M Lake, Ruslan Salakhutdinov, and Joshua B Tenenbaum. Human-level concept learning through probabilistic program induction. *Science*, 350(6266):1332–1338, 2015.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. *AAAI Conference on Artificial Intelligence*, 2018.

- Sitao Luan, Chenqing Hua, Minkai Xu, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Jie Fu, Jure Leskovec, and Doina Precup. When do graph neural networks help with node classification? investigating the homophily principle on node distinguishability. *Advances in Neural Information Processing Systems*, 36, 2024.
- Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph networks. *arXiv preprint arXiv:1905.11136*, 2019.
- John McCall. Genetic algorithms for modelling and optimisation. *Journal of computational and Applied Mathematics*, 184(1):205–222, 2005.
- Geoffrey F Miller, Peter M Todd, and Shailesh U Hegde. Designing neural networks using genetic algorithms. In *ICGA*, volume 89, pp. 379–384, 1989.
- Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5115–5124, 2017.
- Tsendsuren Munkhdalai and Hong Yu. Meta networks. In *International conference on machine learning*, pp. 2554–2563. PMLR, 2017.
- Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node classification. *International Conference on Learning Representation (ICLR)*, 2020.
- Ladislav Rampasek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Dominique Beaini. Recipe for a general, powerful, scalable graph transformer. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- J Redmon. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*, 2020.
- Victor Garcia Satorras and Joan Bruna Estrach. Few-shot learning with graph neural networks. In *International conference on learning representations*, 2018.
- Adarsh Sehgal, Hung La, Sushil Louis, and Hai Nguyen. Deep reinforcement learning using genetic algorithm for parameter optimization. In *2019 Third IEEE International Conference on Robotic Computing (IRC)*, pp. 596–601. IEEE, 2019.
- Min Shi, Yufei Tang, Xingquan Zhu, Yu Huang, David Wilson, Yuan Zhuang, and Jianxun Liu. Genetic-gnn: Evolutionary architecture search for graph neural networks. *Knowledge-based systems*, 247:108752, 2022.
- Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J Sutherland, and Ali Kemal Sinop. Exphormer: Sparse transformers for graphs. In *International Conference on Machine Learning*, 2023.
- Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. *Advances in neural information processing systems*, 30, 2017.
- Jan Toenshoff, Martin Ritzert, Hinrikus Wolf, and Martin Grohe. Graph learning with 1d convolutions on random walks. *arXiv preprint arXiv:2102.08786*, 2021.
- Jan Tönshoff, Martin Ritzert, Eran Rosenbluth, and Martin Grohe. Where did the gap go? reassessing the long-range graph benchmark. *The Second Learning on Graphs Conference (LoG 2023)*, 2023.
- Jan Tönshoff, Martin Ritzert, Eran Rosenbluth, and Martin Grohe. Where did the gap go? reassessing the long-range graph benchmark. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL <https://openreview.net/forum?id=Nm0WX86sKv>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph Attention Networks. *International Conference on Learning Representations*, 2018.

- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, Daan Wierstra, et al. Matching networks for one shot learning. *Advances in neural information processing systems*, 29, 2016.
- Darrell Whitley. An overview of evolutionary algorithms: practical issues and common pitfalls. *Information and software technology*, 43(14):817–831, 2001.
- Xinyi Wu, Amir Ajorlou, Zihui Wu, and Ali Jadbabaie. Demystifying oversmoothing in attention-based graph neural networks. *Advances in Neural Information Processing Systems*, 36:35084–35106, 2023.
- Xinyi Wu, Amir Ajorlou, Zihui Wu, and Ali Jadbabaie. Demystifying oversmoothing in attention-based graph neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jiancong Xie, Yi Wang, Jiahua Rao, Shuangjia Zheng, and Yuedong Yang. Self-supervised contrastive molecular representation learning with a chemical synthesis knowledge graph. *Journal of Chemical Information and Modeling*, 64(6):1945–1954, 2024.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? 2019.
- Haimin Zhang and Min Xu. Randalalign: A parameter-free method for regularizing graph convolutional networks. *arXiv preprint arXiv:2404.09774*, 2024.
- Haimin Zhang, Min Xu, Guoqiang Zhang, and Kenta Niwa. Ssfg: Stochastically scaling features and gradients for regularizing graph convolutional networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- Haimin Zhang, Jiahao Xia, Guoqiang Zhang, and Min Xu. Learning graph representations through learning and propagating edge features. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- Jiahao Zhang, Rui Xue, Wenqi Fan, Xin Xu, Qing Li, Jian Pei, and Xiaorui Liu. Linear-time graph neural networks for scalable recommendations. In *Proceedings of the ACM on Web Conference 2024*, pp. 3533–3544, 2024.
- Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. In *International Conference on Learning Representations*, 2020a.
- Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. *International Conference on Learning Representations*, 2020b.
- Lingxiao Zhao, Wei Jin, Leman Akoglu, and Neil Shah. From stars to subgraphs: Uplifting any gnn with local structure awareness. *arXiv preprint arXiv:2110.03753*, 2021.

A Appendix

Datasets. The experiments were conducted on the following six benchmark datasets.

- **MNIST and CIFAR10** are two datasets for superpixel graph classification (Dwivedi et al., 2020). The superpixels are converted from original images in MNIST (LeCun et al., 1998) and CIFAR10 (Krizhevsky et al., 2009) using the SLIC algorithm (Achanta et al., 2012).
- **PascalVOC-SP and COCO-SP** are two datasets of superpixels (Dwivedi et al., 2022), which are converted from images in original PascalVOC and COCO datasets. The task on the two datasets is to predict if a node corresponds to a region of an image which belongs to a particular class.
- **Peptides-func and Peptides-Struct** (Dwivedi et al., 2022) are two datasets of peptides molecular graphs. The nodes in the graphs represent heavy (non-hydrogen) atoms of the peptides, and the edges represent the bonds between these atoms. The graphs are categorized into 10 classes based on the peptide functions, e.g., antibacterial, antiviral, cell-cell communication. The two datasets are used for evaluating the model’s performance for multi-label graph classification and multi-label graph regression.

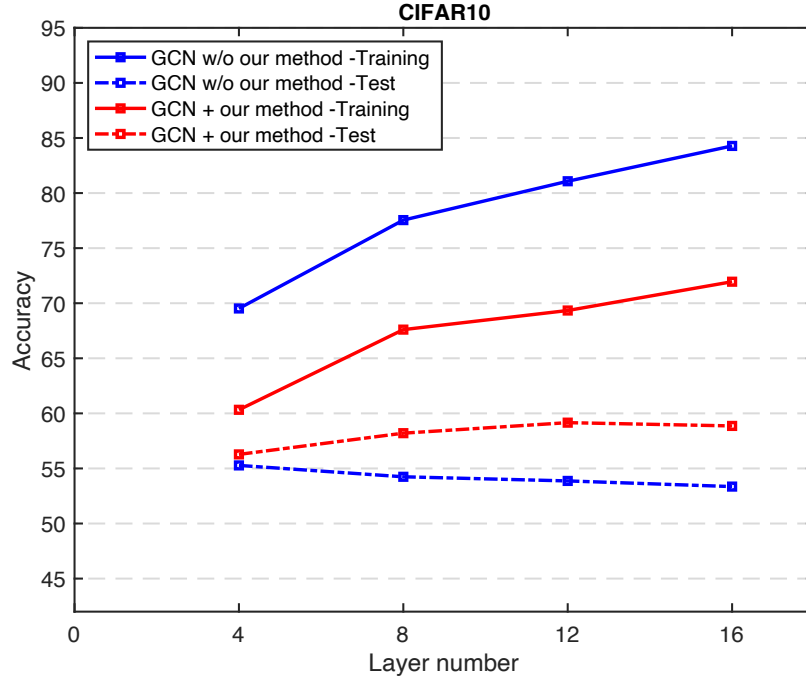


Figure 7: Results of our method on the base Finetuned GCN model with different layers on CIFAR10.

The statistics of the benchmark datasets used in the experiments are shown in below Table 7.

Table 7: Statistics of the six benchmark datasets used in the experiments.

Dataset	Graphs	Nodes	Avg. nodes/graph	#Training	#Validation	#Test	#Categories
MNIST	70K	—	40-75	55,000	5000	10,000	10
CIFAR10	60K	—	85-150	45,000	5000	10,000	10
PascalVOC-SP	11,355	5,443,545	479.40	8,489	1,428	1,429	20
COCO-SP	123,286	58,793,216	476.88	113,286	5,000	5,000	81
Peptides-func	15,535	2,344,859	150.94	70%	15%	15%	10
Peptides-struct	15,535	2,344,859	150.94	70%	15%	15	—

Evaluation Metrics. Following Dwivedi et al. (2020) and Rampasek et al. (2022), the following metrics are used evaluation on different tasks. The performance on MNIST and CIFAR10 on graph classification is evaluated using the classification accuracy. The performance on PascalVOC-SP and COCO-SP on node classification is evaluated using the macro weighted F1 score. The performance on Peptides-func on multi-label graph classification is evaluated using average precision (AP) across the categories. The performance on Peptides-struct on multi-label graph regression is evaluated using mean absolute error (MAE).

Algorithm 3 Cross-generation Crossover, Sibling Crossover, and Mutation

```

1: function CROSSGENERATIONCROSSOVER(h, h_in, crossover_prob)
2:   // current layer embedeing h, previous layer embedding h_in, crossover probability crossover_prob
3:   if model.training == True then // training phase
4:     crossover_mask = Bernoulli.sample(prob=crossover_prob) // each value in
      crossover_mask is sampled from the Bernoulli distribution with probability p
5:     h_crossover = h_in * crossover_mask + h * (1 - crossover_mask)
6:   else // inference phase
7:     h_crossover = h_in * p + h * (1 - p)
8:   end if
9:   return h_crossover
10: end function

11: function SIBLINGCROSSOVER(h_sibling_list)
12:   // list of sibling representations h_sibling_list.
13:   num_siblings = len(h_sibling_list)
14:   if model.training == True then // training phase
15:     sibling_mask = one_hot(Multinomial.sample(prob =  $\frac{1}{num\_siblings}$ , size = num_siblings))
16:     h_crossover = 0
17:     for i ∈ {1, ..., num_siblings} do
18:       h_crossover = h_sibling_list[i] * sibling_mask[i]
19:     end for
20:   else // inference phase
21:     h_crossover = mean(h_sibling_list)
22:   end if
23:   return h_crossover
24: end function

25: function MUTATION(h_batch, mutation_prob)
26:   mean = h_batch.mean()
27:   var = h_batch.var()
28:   batch_size = h_batch.batch_size()
29:   running_mean, running_var = UPDATE_RUNNING_MEAN_VAR(mean, var)
30:   if model.training == True then // training phase
31:     gaussian_noise = Gaussian.sample(h_batch.shape)
32:     mutation_mask = Bernoulli.sample(h_batch.shape)
33:     h = (gaussian_noise * running_var + running_mean) * mutation_mask + h_batch * (1 -
      mutation_mask)
34:   else // inference phase
35:     h = running_mean * mutation_prob + h_batch * (1 - mutation_prob)
36:   end if
37:   return h
38: end function

```

Algorithm 4 The embedding generation process with cross-generation crossover and mutation.

Input: Graph $G = (V, E)$; number of graph neural layers K ; input node features $\{\mathbf{x}_v, \forall v \in V\}$; crossover probability p

Output: Node embeddings $\mathbf{h}_u^{(K)}$ for all $u \in V$

```

1:  $\mathbf{h}_u^{(0)} \leftarrow \mathbf{x}_u, \forall u \in V$ 
2: for  $k = 1, \dots, K$  do
3:   for  $u \in V$  do
4:      $\mathbf{h}_u^{(k)} = \text{Update}^{(k)} \left( \mathbf{h}_u^{(k-1)}, \text{Aggregate}^{(k)}(\{\mathbf{h}_v^{(k-1)}, \forall v \in \mathcal{N}(u)\}) \right)$  // generate an embedding for  $u$ 
        using Equation (1)).
5:      $\mathbf{h}_u^{(k)} \leftarrow \text{CROSSGENERATIONCROSSOVER}(\mathbf{h}_u^{(k)}, \mathbf{h}_u^{(k-1)}, p)$ 
6:      $\mathbf{h}_u^{(k)} \leftarrow \text{MUTATION}(\mathbf{h}_u^{(k)}, \text{mutation\_prob})$ 
7:   end for
8: end for
```

Algorithm 5 The embedding generation process with sibling crossover and mutation.

Input: Graph $G = (V, E)$; number of graph neural layers K ; input node features $\{\mathbf{x}_v, \forall v \in V\}$; number of siblings num_siblings ; mutation probability mutation_prob

Output: Node embeddings $\mathbf{h}_u^{(K)}$ for all $u \in V$

```

1:  $\mathbf{h}_u^{(0)} \leftarrow \mathbf{x}_u, \forall u \in V$ 
2: for  $k = 1, \dots, K$  do
3:   for  $u \in V$  do
4:     for  $s \in \text{num\_siblings}$  do // generate multiple sibling representations using multi-head message
        passing.
5:        $\mathbf{h}_{u,s}^{(k)} = \text{Update}_s^{(k)} \left( \mathbf{h}_u^{(k-1)}, \text{Aggregate}_s^{(k)}(\{\mathbf{h}_v^{(k-1)}, \forall v \in \mathcal{N}(u)\}) \right)$ 
6:     end for
7:      $\mathbf{h}_u^{(k)} \leftarrow \text{SIBLINGCROSSOVER}([\mathbf{h}_{u,1}^{(k)}, \dots, \mathbf{h}_{u,\text{num\_siblings}}^{(k)}])$ 
8:      $\mathbf{h}_u^{(k)} \leftarrow \text{MUTATION}(\mathbf{h}_u^{(k)}, \mathbf{h}_u^{(k-1)}, \text{mutation\_prob})$ 
9:   end for
10: end for
```
