

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 has good interpretability, 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 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, chemical synthesis, and 2D and 3D vision tasks (Zhang et al., 2024; Xie et al., 2024; 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 network was 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

054 message-passing graph neural networks is inherently bounded by the WL isomorphism test. Empirical
055 studies continually find that message-passing graph neural networks suffer from the problem of
056 over-smoothing. That is, the representations for all nodes can become very similar to one another
057 after too many message passing iterations. These core limitations prevent graph neural networks
058 from more meaningful representations from graphs. In recent years, increasing studies have been
059 devoted to addressing the bottlenecks, such as normalization and regularization techniques (?), and
060 combining the global self-attention mechanism (Rampasek et al., 2022), exploring generalized mes-
061 sage passing (Barceló et al., 2020). Regardless of these advances, improving the capability of graph
062 neural network models still remains a fundamental challenge in learning from graph-structured data.

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

071 Genetic algorithms, inspired by the Charles Darwin’s theory of natural evolution, emulate the pro-
072 cess of natural selection, wherein the fittest individuals are selected to reproduce and generate the
073 next generation of offspring. Genetic algorithms employ a set of evolution-inspired operations, in-
074 cluding mutation, crossover, and selection (Mitchell, 1998). Over multiple generations, biological
075 organisms evolve based on the principle of natural selection, or “survival of the fittest”, enabling
076 them to accomplish target tasks. Genetic algorithms have been successfully applied in solving com-
077 plex optimization and search problems. In machine learning, genetic algorithms have also been
078 used for feature selection (Babatunde et al., 2014) and hyperparameter tuning for models like neural
079 networks and support vector machines (Alibrahim & Ludwig, 2021).

080 In genetic algorithms, the crossover and mutation operations play a key role in generating diverse
081 individuals for selection, preventing the algorithms from premature convergence (Gupta & Ghafir,
082 2012). Crossover introduces variety by combining genetic information from different parents, and
083 mutation introduces small random changes in genetic information. In this work, we view the itera-
084 tive node embedding process as an evolutionary process, in which each layer of message passing
085 produces a new generation of embeddings. We introduce two crossover operations, i.e., cross-
086 generation crossover and sibling crossover, and a mutation operation, and we develop two graph
087 neural network building blocks based on the operations. At each layer of a graph neural network,
088 we first use message passing to update node representations and then apply crossover and muta-
089 tion to prevent embeddings from becoming similar to one another, thus enabling the model to learn
improved graph representations.

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

099 The main contributions of this paper can be summarized as follows. (1) This paper proposes a new
100 framework named genetic-evolutionary graph neural networks for learning from graph-structured
101 data. The core idea behind the proposed framework is to model each layer of a graph neural network
102 as an evolutionary process. We develop three key operations inspired by crossover and mutation
103 from genetic algorithms to enhance the diversity of generated embeddings at each layer. (2) The
104 proposed framework offers good interpretability, as it is directly inspired by biogenetics. It is a
105 general paradigm which can be integrated into current message-passing graph neural networks. Em-
106 pirical evaluations are conducted on six popular datasets on different graph tasks, and the results
107 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 (Veličković 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.

Both spectral and spatial graph neural networks are essentially message-passing neural networks that employ a paradigm of message passing wherein embeddings are exchanged between nodes and updated using neural networks (Gilmer et al., 2017). 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, 2020; Chen et al., 2022), combing 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. \(2024\) analyzed homophily by studying intra- and inter-class node distinguishability and showed that graph neural network is capable of generating meaningful representations regardless of homophily levels.](#)

2.2 GENETIC ALGORITHMS

Genetic algorithm methods are inspired by the mechanisms of evolution and natural genetics (Srinivas & Patnaik, 1994). Genetic algorithms were first introduced by Holland (1992) as a heuristic method based on the principle of nature selection. Over the past years, genetic algorithms have emerged as a powerful tool for solving complex optimization and search problems across numerous fields such as scheduling, mathematics and networks (Alhijawi & Awajan, 2023).

In machine learning, genetic algorithms have been applied for optimizing neural networks (Miller et al., 1989) and designing neural network architectures (Jones, 1993). Researchers have also used genetic algorithms for optimizing hyperparameters in neural networks and support vector machines (Alibrahim & Ludwig, 2021; Shanthy & Chethan, 2022). 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.

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 reorientation 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. 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

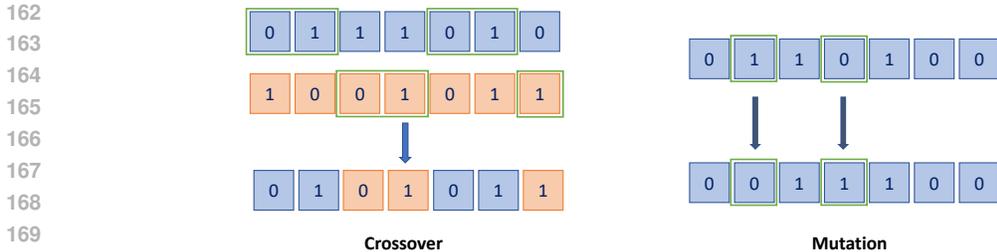


Figure 1: Crossover recombines of the genetic information of parents to produce an offspring. Mutation introduces small random changes in genetic information.

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 diversity of embeddings 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 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 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 converging 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:

$$\mathbf{h}_u^{(k)} = \text{Crossover}(\bar{\mathbf{h}}_u^{(k)}, \mathbf{h}_u^{(k-1)})$$

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

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.

$$\mathbf{h}_u^{(k)} = \text{Crossover}(\bar{\mathbf{h}}_u^{(k, \text{head}_1)}, \dots, \bar{\mathbf{h}}_u^{(k, \text{head}_z)})$$

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

$h_{ij} \sim \text{Categorical}(\frac{1}{z}, \dots, \frac{1}{z})$, and z is the number of heads. Each $\bar{\mathbf{h}}_u^{(k, \text{head}_h)}$ 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 Gaus-

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

```

270
271
272 # h, h_in: representaton generated by message passing and the previous layer embedding
273 # f_prob: probabily of recombining information from parent
274 # self.dist: a Bernoulli distribution defined by torch.distributions.Bernoulli(torch.tensor(
275     self.f_prob)):
276
277 def forward(self, h, h_in):
278     if self.training == True:
279         crossover_mask = self.dist.sample(h.shape) # generate crossover mask
280
281         # crossover from h and h_in
282         h = h_in * crossover_mask + h * (1 - crossover_mask)
283     else:
284         h = h_in * self.f_prob + h * (1 - self.f_prob)
285
286     return h

```

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

```

285
286 # self.running_mean: the mean of h over the training set
287 # self.running_var: the variance of h over the training set
288 # self.mutation_prob: probability of mutation
289
290 def forward(self, h):
291     if self.training == True:
292         mean = h.mean([0])
293         var = h.var([0])
294         n = h.numel() / h.size(1)
295
296         with torch.no_grad():
297             # momentum update of running_mean and running_var
298             self.running_mean = self.momentum * mean + (1 - self.momentum) * self.running_mean
299             self.running_var = self.momentum * var * n / (n - 1) + (1 - self.momentum) * self.
300                 running_var
301
302         # generate mutatioin noise
303         gaussian_noise = torch.randn(h.shape)
304
305         if self.training == True:
306             mutation_mask = Bernoulli.sample(h.shape) # generate mutation mask
307             h = (gaussian_noise * self.running_var + self.running_mean) * mutation_mask + h * (1 -
308                 mutation_mask)
309         else:
310             h = self.running_mean * self.mutation_prob + h * (1 - self.mutation_prob)
311
312     return h

```

sian 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.

$$\begin{aligned}
 \boldsymbol{\mu} &\leftarrow \mathbb{E}_{\mathcal{B}}(\boldsymbol{\mu}_{\mathcal{B}}) \\
 \boldsymbol{\delta} &\leftarrow \frac{m}{m-1} \mathbb{E}_{\mathcal{B}}(\boldsymbol{\delta}_{\mathcal{B}}^2)
 \end{aligned} \tag{4}$$

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:

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

where the $\mathbf{mask} \sim \text{Bernoulli}(\text{mutation_rate})$. 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

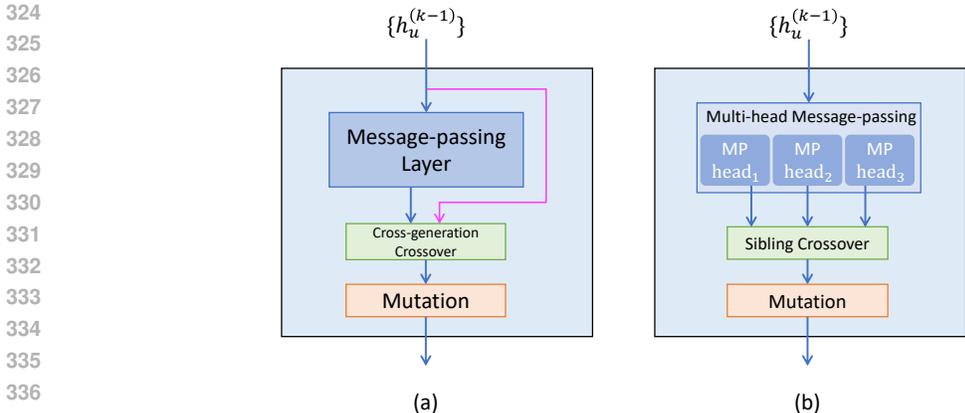


Figure 2: Building block architectures: Block (a) applies cross-generation to a node’s embedding generated using message passing and the node’s previous layer embedding, and Block (b) applies sibling crossover to a set of outputs generated using multi-head message passing.

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 GPS model.

Model	MNIST	CIFAR10
GCN (Kipf & Welling, 2016)	90.705±0.218	55.710±0.381
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 (Rampasek et al., 2022)	90.705±0.218	55.710±0.381
Ours + GCN	95.926±0.031	59.157±0.130
GPS (Rampasek et al., 2022)	98.051±0.126	72.298±0.356
Finetuned GPS	98.186±0.107	75.680±0.188
Ours + Finetuned GPS	98.685±0.029	80.636±0.195

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.

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 GPS model.

Model	PascalVOC-SP (F1)	COCO-SP (F1)
GCN Kipf & Welling (2016)	0.1268±0.0060	0.0841±0.0010
GINE Hu et al. (2019)	0.1265±0.0076	0.1339±0.0044
GCNII Chen et al. (2020)	0.1698±0.0080	0.1404±0.0011
GatedGCN Bresson & Laurent (2017)	0.2873±0.0219	0.2641±0.0045
GatedGCN + RWSE (Rampasek et al., 2022)	0.2860±0.0085	0.2574±0.0034
Transformer + LapPE Dwivedi et al. (2022)	0.2694±0.0098	0.2618±0.0031
SAN + LapPE Dwivedi et al. (2022)	0.3230±0.0039	0.2592±0.0158
SAN + RWSE Dwivedi et al. (2022)	0.3216±0.0027	0.2434±0.0156
Exphormer Shirzad et al. (2023)	0.3975±0.0037	0.3455±0.0009
RandAlign + GPS (Zhang & Xu, 2024)	0.4242±0.0011	0.3567±0.0026
Fine-tuned GCN (Tönshoff et al., 2023)	0.2078±0.0031	–
Ours + Finetuned GCN	0.2241±0.0020	–
GPS (Rampasek et al., 2022)	0.3748±0.0109	0.3412±0.0044
Fine-tuned GPS (Tönshoff et al., 2023)	0.4440±0.0065	0.3884±0.0055
Ours + Finetuned GPS	0.4832±0.0031	0.4002±0.0019

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

Model	Peptides-func (AP ↑)	Peptides-struct (MAE ↓)
GCN	0.5930±0.0023	0.3496±0.0013
GINE	0.5498±0.0079	0.3547±0.0045
GCNII (Chen et al., 2020)	0.5543±0.0078	–
GatedGCN	0.5864±0.0077	0.3420±0.0013
Gated + RWSE	0.6069±0.0035	0.3357±0.0006
Transformer+LapPE	0.6326±0.0126	0.2529±0.0016
SAN+LapPE	0.6384±0.0121	0.2683±0.0043
SAN+RWSE	0.6439±0.0075	0.2545±0.0012
Exphormer (Shirzad et al., 2023)	0.6527±0.0043	0.2481±0.0007
GPS (Rampasek et al., 2022)	0.6535±0.0041	0.2500±0.0005
Finetuned GPS (Tönshoff et al., 2023)	0.6534±0.0091	0.2509±0.0014
Finetuned GCN (Tönshoff et al., 2023)	0.6860±0.0050	0.2460±0.0007
Ours + Finetuned GCN	0.7021±0.0034	0.2426±0.0014

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.

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

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

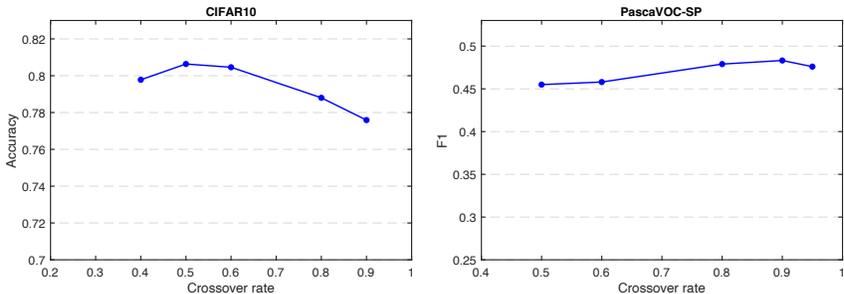


Figure 3: Impact of the crossover rate p on the model performance on CIFAR10 and PascalVOC-SP.

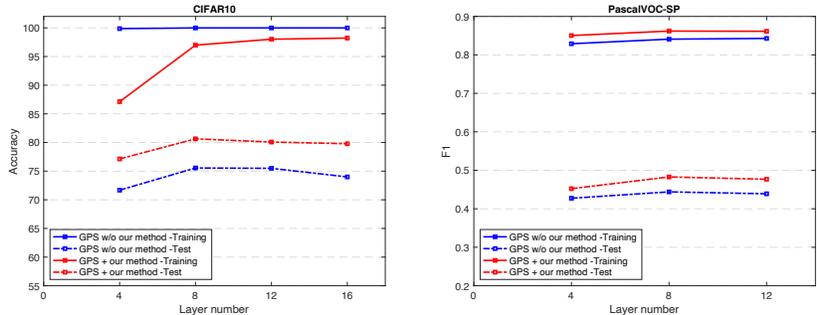


Figure 4: Results of our method on the base Finetuned GPS model with different layers on CIFAR10 and PascalVOC-SP.

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 model by a large margin, with a relative improvement of 0.648% and 11.53% on MNIST and CIFAR10 respectively. It simultaneously outperforms both Exphormer (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 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

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±0.336	48.810±1.045
GCN (with residual connections and BN)	90.705±0.218	55.710±0.381
GCN (with residual connections and BN) + Ours	95.926±0.031	59.157±0.130

of siblings is set to 2 and mutation rate is set to 0.1. The results are reported in Table 3. Finetuned GCN is a strong baseline model in previous work. We see from Table 3 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 0.9 or 0.95 and then gradually decreasing it to find the optimal value.

We conducted experiments to analysis 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 the performance of our method on the base Finetuned GPS model with different layers on CIFAR10, and the results are reported in Figure 5 in the appendix section. It can be seen from Figure 4 and Figure 5 that the use of our method improves the model generalization performance on the base models.

We further 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.

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.

- 540 B Alhijawi and A Awajan. Genetic algorithms: theory, genetic operators, solutions, and applications,
541 evol. intel., 2023, 2023.
- 542
- 543 Hussain Alibrahim and Simone A Ludwig. Hyperparameter optimization: Comparing genetic al-
544 gorithm against grid search and bayesian optimization. In *2021 IEEE Congress on Evolutionary*
545 *Computation (CEC)*, pp. 1551–1559. IEEE, 2021.
- 546 Oluleye H Babatunde, Leisa Armstrong, Jinsong Leng, and Dean Diepeveen. A genetic algorithm-
547 based feature selection. 2014.
- 548
- 549 Pablo Barceló, Egor V Kostylev, Mikael Monet, Jorge Pérez, Juan Reutter, and Juan-Pablo Silva.
550 The logical expressiveness of graph neural networks. In *8th International Conference on Learning*
551 *Representations (ICLR 2020)*, 2020.
- 552 Dominique Beaini, Saro Passaro, Vincent Létourneau, Will Hamilton, Gabriele Corso, and Pietro
553 Liò. Directional graph networks. In *International Conference on Machine Learning*, pp. 748–
554 758. PMLR, 2021.
- 555
- 556 Xavier Bresson and Thomas Laurent. Residual gated graph convnets. *arXiv preprint*
557 *arXiv:1711.07553*, 2017.
- 558
- 559 Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally
560 connected networks on graphs. *International Conference on Learning Representations*, 2014.
- 561 Chaoqi Chen, Yushuang Wu, Qiyuan Dai, Hong-Yu Zhou, Mutian Xu, Sibe Yang, Xiaoguang Han,
562 and Yizhou Yu. A survey on graph neural networks and graph transformers in computer vision:
563 A task-oriented perspective. *IEEE Transactions on Pattern Analysis and Machine Intelligence*,
564 2024.
- 565
- 566 Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph con-
567 volutional networks. In *International Conference on Machine Learning*, pp. 1725–1735. PMLR,
568 2020.
- 569
- 570 Yihao Chen, Xin Tang, Xianbiao Qi, Chun-Guang Li, and Rong Xiao. Learning graph normalization
571 for graph neural networks. *Neurocomputing*, 493:613–625, 2022.
- 572
- 573 Yun Young Choi, Sun Woo Park, Minhoo Lee, and Youngho Woo. Topology-informed graph trans-
574 former. *arXiv preprint arXiv:2402.02005*, 2024.
- 575
- 576 Gabriele Corso, Luca Cavalleri, Dominique Beaini, Pietro Liò, and Petar Veličković. Principal
577 neighbourhood aggregation for graph nets. *Advances in Neural Information Processing Systems*,
578 33:13260–13271, 2020.
- 579
- 580 Vijay Prakash Dwivedi, Chaitanya K Joshi, Thomas Laurent, Yoshua Bengio, and Xavier Bresson.
581 Benchmarking graph neural networks. *arXiv preprint arXiv:2003.00982*, 2020.
- 582
- 583 Vijay Prakash Dwivedi, Ladislav Rampasek, Mikhail Galkin, Ali Parviz, Guy Wolf, Anh Tuan Luu,
584 and Dominique Beaini. Long range graph benchmark. In *Thirty-sixth Conference on Neural*
585 *Information Processing Systems Datasets and Benchmarks Track*, 2022.
- 586
- 587 Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
588 message passing for quantum chemistry. In *International conference on machine learning*, pp.
589 1263–1272. PMLR, 2017.
- 590
- 591 Deepti Gupta and Shabina Ghafir. An overview of methods maintaining diversity in genetic al-
592 gorithms. *International journal of emerging technology and advanced engineering*, 2(5):56–60,
593 2012.
- 594
- 595 Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
596 In *Advances in neural information processing systems*, pp. 1024–1034, 2017.
- 597
- 598 William L Hamilton. Graph representation learning. *Synthesis Lectures on Artificial Intelligence and*
599 *Machine Learning*, 14(3):1–159, 2020.

- 594 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
595 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
596 770–778, 2016.
- 597 John H Holland. *Adaptation in natural and artificial systems: an introductory analysis with appli-*
598 *cations to biology, control, and artificial intelligence*. MIT press, 1992.
- 600 Weihua Hu, Bowen Liu, Joseph Gomes, Marinka Zitnik, Percy Liang, Vijay Pande, and Jure
601 Leskovec. Strategies for pre-training graph neural networks. *arXiv preprint arXiv:1905.12265*,
602 2019.
- 604 Md Shamim Hussain, Mohammed J Zaki, and Dharmashankar Subramanian. Global self-attention
605 as a replacement for graph convolution. In *Proceedings of the 28th ACM SIGKDD Conference on*
606 *Knowledge Discovery and Data Mining*, pp. 655–665, 2022.
- 607 Antonia J Jones. Genetic algorithms and their applications to the design of neural networks. *Neural*
608 *Computing and Applications*, 1(1):32–45, 1993.
- 610 Sangwon Kim, Dasom Ahn, and Byoung Chul Ko. Cross-modal learning with 3d deformable atten-
611 tion for action recognition. In *Proceedings of the IEEE/CVF international conference on com-*
612 *puter vision*, pp. 10265–10275, 2023.
- 613 Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional net-
614 works. *arXiv preprint arXiv:1609.02907*, 2016.
- 616 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
617 2009.
- 618 Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to
619 document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 621 Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for
622 semi-supervised learning. *AAAI Conference on Artificial Intelligence*, 2018.
- 623 Sitao Luan, Chenqing Hua, Minkai Xu, Qincheng Lu, Jiaqi Zhu, Xiao-Wen Chang, Jie Fu, Jure
624 Leskovec, and Doina Precup. When do graph neural networks help with node classification?
625 investigating the homophily principle on node distinguishability. *Advances in Neural Information*
626 *Processing Systems*, 36, 2024.
- 628 Haggai Maron, Heli Ben-Hamu, Hadar Serviansky, and Yaron Lipman. Provably powerful graph
629 networks. *arXiv preprint arXiv:1905.11136*, 2019.
- 631 Geoffrey F Miller, Peter M Todd, and Shailesh U Hegde. Designing neural networks using genetic
632 algorithms. In *ICGA*, volume 89, pp. 379–384, 1989.
- 633 Melanie Mitchell. *An introduction to genetic algorithms*. MIT press, 1998.
- 634 Federico Monti, Davide Boscaini, Jonathan Masci, Emanuele Rodola, Jan Svoboda, and Michael M
635 Bronstein. Geometric deep learning on graphs and manifolds using mixture model cnns. In
636 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5115–5124,
637 2017.
- 639 Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node
640 classification. *International Conference on Learning Representation (ICLR)*, 2020.
- 642 Ladislav Rampasek, Mikhail Galkin, Vijay Prakash Dwivedi, Anh Tuan Luu, Guy Wolf, and Do-
643 minique Beaini. Recipe for a general, powerful, scalable graph transformer. In Alice H. Oh,
644 Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information*
645 *Processing Systems*, 2022.
- 646 J Redmon. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE*
647 *conference on computer vision and pattern recognition*, 2016.

- 648 Adarsh Sehgal, Hung La, Sushil Louis, and Hai Nguyen. Deep reinforcement learning using genetic
649 algorithm for parameter optimization. In *2019 Third IEEE International Conference on Robotic*
650 *Computing (IRC)*, pp. 596–601. IEEE, 2019.
- 651 DL Shanthi and N Chethan. Genetic algorithm based hyper-parameter tuning to improve the perfor-
652 mance of machine learning models. *SN Computer Science*, 4(2):119, 2022.
- 653 Hamed Shirzad, Ameya Velingker, Balaji Venkatachalam, Danica J Sutherland, and Ali Kemal
654 Sinop. Exphormer: Sparse transformers for graphs. In *International Conference on Machine*
655 *Learning*, 2023.
- 656 Mandavilli Srinivas and Lalit M Patnaik. Genetic algorithms: A survey. *computer*, 27(6):17–26,
657 1994.
- 658 Jan Toenshoff, Martin Ritzert, Hinrikus Wolf, and Martin Grohe. Graph learning with 1d convolu-
659 tions on random walks. *arXiv preprint arXiv:2102.08786*, 2021.
- 660 Jan Tönshoff, Martin Ritzert, Eran Rosenbluth, and Martin Grohe. Where did the gap go? reassess-
661 ing the long-range graph benchmark. *The Second Learning on Graphs Conference (LoG 2023)*,
662 2023.
- 663 Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua
664 Bengio. Graph Attention Networks. *International Conference on Learning Representations*,
665 2018.
- 666 Darrell Whitley. An overview of evolutionary algorithms: practical issues and common pitfalls.
667 *Information and software technology*, 43(14):817–831, 2001.
- 668 Xinyi Wu, Amir Ajorlou, Zihui Wu, and Ali Jadbabaie. Demystifying oversmoothing in attention-
669 based graph neural networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- 670 Jiancong Xie, Yi Wang, Jiahua Rao, Shuangjia Zheng, and Yuedong Yang. Self-supervised con-
671 trastive molecular representation learning with a chemical synthesis knowledge graph. *Journal of*
672 *Chemical Information and Modeling*, 64(6):1945–1954, 2024.
- 673 Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
674 networks? 2019.
- 675 Haimin Zhang and Min Xu. Randalign: A parameter-free method for regularizing graph convolu-
676 tional networks. *arXiv preprint arXiv:2404.09774*, 2024.
- 677 Haimin Zhang, Min Xu, Guoqiang Zhang, and Kenta Niwa. Ssfg: Stochastically scaling features and
678 gradients for regularizing graph convolutional networks. *IEEE Transactions on Neural Networks*
679 *and Learning Systems*, 2022.
- 680 Haimin Zhang, Jiahao Xia, Guoqiang Zhang, and Min Xu. Learning graph representations through
681 learning and propagating edge features. *IEEE Transactions on Neural Networks and Learning*
682 *Systems*, 2023.
- 683 Jiahao Zhang, Rui Xue, Wenqi Fan, Xin Xu, Qing Li, Jian Pei, and Xiaorui Liu. Linear-time graph
684 neural networks for scalable recommendations. In *Proceedings of the ACM on Web Conference*
685 *2024*, pp. 3533–3544, 2024.
- 686 Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. In *International*
687 *Conference on Learning Representations*, 2020.
- 688 Lingxiao Zhao, Wei Jin, Leman Akoglu, and Neil Shah. From stars to subgraphs: Uplifting any gnn
689 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.

The statistics of the benchmark datasets used in the experiments are shown in below Table 6.

Table 6: 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).

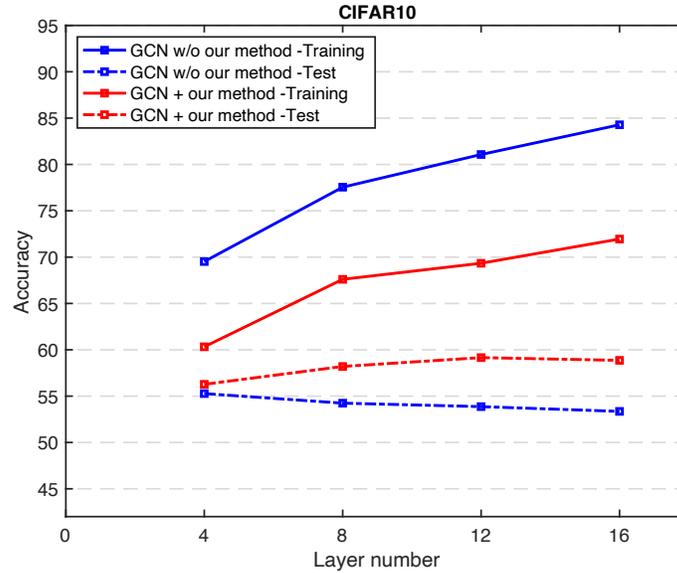


Figure 5: Results of our method on the base Finetuned GCN model with different layers on CIFAR10.

```

7 class Crossover(nn.Module):
8
9     def __init__(self, f_prob=0):
10         super(Crossover, self).__init__()
11         self.f_prob = f_prob #prob of sampling 1. for Bern
12         self.mm = torch.distributions.bernoulli.Bernoulli(torch.tensor([self.f_prob],
device='cuda'))
13
14     def forward(self, h, h_in):
15         if self.training:
16             crossover_mask = self.mm.sample(h.shape).squeeze(-1)
17             h = h_in * crossover_mask + h * (1 - crossover_mask)
18         else:
19             h = h_in * self.f_prob + h * (1 - self.f_prob)
20
21         return h

```

Figure 6: Implementation of the crossover operation in Pytorch.

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

```

25 class Mutation(nn.BatchNorm1d):
26     def __init__(self, num_features, mutate_prob, eps=1e-5, momentum=0.1,
27                 affine=True, track_running_stats=True):
28         super(Mutation, self).__init__(
29             num_features, eps, momentum, affine, track_running_stats)
30         self.mutate_prob = mutate_prob
31
32     def forward(self, input):
33         self._check_input_dim(input)
34         exponential_average_factor = 0.0
35
36         if self.training and self.track_running_stats:
37             if self.num_batches_tracked is not None:
38                 self.num_batches_tracked += 1
39                 if self.momentum is None: # use cumulative moving average
40                     exponential_average_factor = 1.0 / float(self.num_batches_tracked)
41             else: # use exponential moving average
42                 exponential_average_factor = self.momentum
43
44         # calculate running estimates
45         if self.training:
46             mean = input.mean([0])
47             # use biased var in train
48             var = input.var([0], unbiased=False)
49             n = input.numel() / input.size(1)
50             with torch.no_grad():
51                 self.running_mean = exponential_average_factor * mean\
52                     + (1 - exponential_average_factor) * self.running_mean
53                 # update running_var with unbiased var
54                 self.running_var = exponential_average_factor * var * n / (n - 1)\
55                     + (1 - exponential_average_factor) * self.running_var
56         else:
57             mean = self.running_mean
58             var = self.running_var
59         mean = self.running_mean
60         var = self.running_var
61
62         gaussian_noise = torch.randn(input.shape).type_as(input)
63
64         prob_mutate = self.mutate_prob
65         if self.training:
66             mm = torch.distributions.bernoulli.Bernoulli(torch.tensor([prob_mutate],
67 device='cuda'))
68             mutate_mask = mm.sample(input.shape).squeeze(-1)
69             input = (gaussian_noise * var + mean)*mutate_mask + input * (1 - mutate_
70 mask)
71         else:
72             input = mean * prob_mutate + input * (1 - prob_mutate)
73         return input

```

Figure 7: Implementation of the mutation operation in Pytorch.

```

864 Algorithm 3 Pseudo for the cross-generation crossover operation.
865 Input: Crossover probability  $p$ ,  $\mathbf{h}$ ,  $\mathbf{h\_in}$  // embeddings generated by the current layer and the
866 previous layer
867 Output:  $\mathbf{h\_crossover}$  // Crossover of  $\mathbf{h}$  and  $\mathbf{h\_in}$ 
868 1: if model.training == True then
869 2:    $\mathbf{crossover\_mask} = \text{Bernoulli.sample}(\text{prob}=p)$  // each element in  $\mathbf{crossover\_mask}$ 
870   is sampled from the Bernoulli distribution with probability  $p$ 
871 3:    $\mathbf{h\_crossover} = \mathbf{h\_in} * \mathbf{crossover\_mask} + \mathbf{h} * (1 - \mathbf{crossover\_mask})$ 
872 4: else
873 5:    $\mathbf{h\_crossover} = \mathbf{h\_in} * p + \mathbf{h} * (1 - p)$ 
874 6: end if

```

```

876 Algorithm 4 Pseudo for the mutation operation.
877 Input: Node embedding  $\mathbf{h}$ , mutation probability  $r$ 
878 Output:  $\mathbf{h\_mutation}$  // Mutation output of  $\mathbf{h}$ 
879 1:  $\mathbf{running\_mean}, \mathbf{running\_var} = \text{Update}(\mathbf{h})$  // update running mean and var
880 2:  $\mathbf{gaussian\_noise} = \text{Gaussian.Sample}()$  //the reparameterization trick
881 3: if model.training == True then
882 4:    $\mathbf{mutation\_mask} = \text{Bernoulli.sample}(\text{prob}=r)$  // each element in  $\mathbf{mutation\_mask}$ 
883   is sampled from the Bernoulli distribution with probability  $p$ 
884 5:    $\mathbf{h\_mutation} = (\mathbf{gaussian\_noise} * \mathbf{running\_var} + \mathbf{running\_mean}) * \mathbf{mutation\_mask} + \mathbf{h} * (1 - \mathbf{mutation\_mask})$ 
885 6: else
886 7:    $\mathbf{h\_mutation} = \mathbf{running\_mean} * r + \mathbf{h} * (1 - r)$ 
887 8: end if

```

889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917