# PAGE: Parametric Generative Explainer for Graph Neural Network

Yang Qiu<sup>a</sup>, Wei Liu<sup>a</sup>, Jun Wang b,\* and Ruixuan Li<sup>a,\*\*</sup>

<sup>a</sup>School of Computer Science and Technology, Huazhong University of Science and Technology, China <sup>b</sup>iWudao Tech

Abstract. This article introduces PAGE, a parameterized generative interpretive framework. PAGE is capable of providing faithful explanations for any graph neural network without necessitating prior knowledge or internal details. Specifically, we train the autoencoder to generate explanatory substructures by designing appropriate training strategy. Due to the dimensionality reduction of features in the latent space of the autoencoder, it becomes easier to extract causal features leading to the model's output, which can be easily employed to generate explanations. To accomplish this, we introduce an additional discriminator to capture the causality between latent causal features and the model's output. By designing appropriate optimization objectives, the well-trained discriminator can be employed to constrain the encoder in generating enhanced causal features. Finally, these features are mapped to substructures of the input graph through the decoder to serve as explanations. Compared to existing methods, PAGE operates at the sample scale rather than nodes or edges, eliminating the need for perturbation or encoding processes as seen in previous methods. Experimental results on both artificially synthesized and real-world datasets demonstrate that our approach not only exhibits the highest faithfulness and accuracy but also significantly outperforms baseline models in terms of efficiency.

#### 1 Introduction

Graph Neural Networks (GNNs) have consistently demonstrated promising results across a wide variety of tasks [25, 26, 33]. As GNNs find applications in crucial domains where trustworthy AI is imperative, the need for their interpretability has become increasingly paramount [30]. GNNExplainer is the first general modelagnostic approach for interpreting GNNs, and searches for soft masks for edges and node features to explain the predictions via mask optimization [27]. Since GNNExplainer largely focuses on providing the local interpretability by generating a painstakingly customized explanation for a single instance individually and independently, it is not sufficient to provide a global understanding of the trained model [14]. Furthermore, GNNExplainer has to be retrained for every single explanation. So, in real-world scenarios, GNNExplainer would be time-consuming and impractical. To address the above issues, PGExplainer was proposed to learn succinct underlying structures as the explanations from the observed graph data [14]. It also models the underlying structure as edge distributions, where the explanatory subgraph is sampled using binary hard masks for edges. PGExplainer utilizes a deep neural network to parameterize the explanation selection process, providing the ability to collectively explain multiple instances, so it has better generalization ability and can be utilized in an inductive setting easily [14]. Since there is no ground truth for explanatory subgraphs, PGExplainer has to employ reinforcement learning or Gumbel-Softmax sampling to search for informative explanatory subgraphs within a vast space that grows exponentially with the number of edges. Consequently, the computational complexity remains relatively high [12].

Instead of using edge masks, recent methods have attempted to produce the adjacency matrix of explanatory subgraphs directly through a Variational Graph Auto-encoder (VGAE) [9], eliminating the need to consider each node or edge individually and achieving greater efficiency. GEM is the pioneering work in this domain. It introduced a distillation process grounded in the concept of Granger causality [4] to generate ground-truth explanatory subgraphs for the training of VGAE [12]. However, the distillation process intrinsically assumes independence between the edges. This could be problematic as graph-structured data is inherently interdependent [13]. Differing from GEM, which quantifies the causal influence from edges, OrphicX [13] suggests identifying the underlying causal factors from the latent space, allowing it to bypass direct interaction with intricately interdependent edges. It divides the latent representation produced by the encoder into causal features and spurious features and identifies the underlying causal features by harnessing information flow measurements [1], quantifying the causal information emanating from the latent features. It then constructs the adjacency matrix for the explanatory subgraph based on these causal features [13].

The fundamental architecture of Gem and OrphicX is depicted in Figure 1. Gem, while illustrating the viability of training VGAE as a graph interpreter, is subject to certain constraints. Gem's methodology involves a distillation process aimed at generating ground-truth subgraphs for training the autoencoder. This process perturbs the input graph at the edge level and assesses the importance of each edge based on the resulting reduction in model error. However, this approach assumes edge independence, failing to adequately explore the ability to discern causal features within the autoencoder's latent space. OrphicX endeavors to utilize information flow estimation to evaluate the causal relationships between each feature dimension in the latent space and the model predictions. This strategy aims to bypass perturbations from input dimensions. Nonetheless, due to computational constraints and complexity considerations, OrphicX operates within a severely restricted sampling range for each sample. For instance, regardless of the size of a single input graph, OrphicX sam-

<sup>\*</sup> Corresponding Author (Email: jwang@iwudao.tech).

<sup>\*\*</sup> Corresponding Author (Email: rxli@hust.edu.cn). This paper is a collaboration between Intelligent and Distributed Computing Laboratory, Huazhong University of Science and Technology and iWudao Tech.

ples only five nodes of the graph to estimatie the causal effects of the whole graph on the output, which is obviously a limitation that undermines practicality and necessary adaptability. Moreover, the efficacy of information flow measurements in OrphicX appears limited. This assertion is supported by our experimental findings, which indicate that excluding the information flow term from the loss function has minimal impact on the model's performance. Furthermore, both the distillation process in Gem and the information flow measurement in OrphicX necessitate numerous perturbations or samplings on individual samples. These procedures incur significant time and space complexity, posing challenges to the application of these methods on large-scale graphs and extensive datasets.

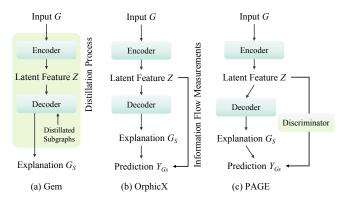
To address these problems in OrphicX, we introduce a novel and more efficient model termed the Parametric Generative Explainer (PAGE<sup>1</sup>). Specifically, our optimization objective remains to maximize the mutual information between explanations and original predictions, but we introduce an extra parameterized discriminator to acquire the global understanding of the causal features in the latent space. As illustrated in Figure 1(c) and Figure 2, our model consists of an autoencoder and an additional discriminator. The features in the latent space compressed by the encoder are divided into two parts: causally relevant features related to predictions are used to generate explanations, while non-causal features are discarded. The training process consists of two stages: In the first stage, autoencoder and discriminator are trained together. The well-trained parameterized discriminator is designed to maximize the mutual information between causal features and prediction results. Parameterized discriminator eliminates the need for information flow estimation through sampling. Furthermore, since the discriminator is well-trained, it possesses a global comprehension on the entire dataset. In the second stage, the parameters of the discriminator are frozen to constrain the encoder to learn better causal features. Gem requires calculating causal contributions for each edge of a instance, while OrphicX involves extensive sampling across different dimensions including samples and features. In comparison, our approach is significantly more efficient than these methods. The trained explainer can collectively generate explanations, eliminating the need for extensive sampling and computations.

Our main contributions can be outlined as follows:

- 1) We introduce PAGE, a novel generative GNN explanation method. This method replaces complicated sampling processes with a streamlined discriminator, enhancing effectiveness.
- 2) We further refine the optimization objective, eliminating the need of complex perturbation or sampling process. Instead, we directly align predictions from the original graph with those from the explanatory subgraph.
- 3) Our research included experiments spanning both node and edge levels, conducted on both synthesized and real-world datasets. Experimental results clearly show that our approach outperforms existing methods across various metrics. This includes, but is not limited to, confidence and accuracy. Moreover, our method boasts significantly higher efficiency when juxtaposed with prior methods.

#### 2 Related Work

Post-hoc GNN Explanation aims to produce an explanation for a GNN prediction on a given graph, usually as a substructure of the graph. Various explaining approaches focus on different aspects of the model and also provide different views. Many sur-



**Figure 1.** Methods using autoencoder as explainer: (a) Gem[12], (b) OrphicX[13], (c) PAGE(Ours)

veys categorized and summarized existing GNN explaining methods [29, 17, 11, 3]. Some methods provide explanations for each individual input instance separately, such as GNNExplainer [27], while some train a parameterized explainer to provide explanations collectively for multiple instances, like PGExplainer [14]. Depending on whether the explainer generates explanations separately or collectively, we categorize existing explanation methods into non-parameterized and parameterized approaches.

Non-parameterized: Non-parameterized methods like GNNExplainer generate explanations for individual instances and predictions through perturbation or optimization processes. Some gradientbased methods compute gradients of the target prediction with respect to input features through back-propagation. Moreover, featurebased methods map hidden features to the input space via interpolation to calculate importance scores. Examples of such methods include SA [3], GuidedBP [3], and Grad-CAM [15]. Additionally, there are some methods based on cooperative game theory, which assign importance scores to input nodes and edges by introducing concepts from cooperative games like Shapley Value. Examples include SubgraphX [29] and GStarX [31]. Moreover, some approaches integrate reinforcement learning by treating the addition or removal of nodes and edges as strategies, while considering the probabilities of the target GNN as rewards, and then utilize a policy network to generate explanations, aiming to maximize the attained rewards, like RC-Explainer [24] and RG-Explainer [19].

Parameterized: Parameterized methods like PGExplainer train an explainer on the entire training set to provide explanations for multiple instances collectively. These methods utilize a parameterized explaining network, typically trained to maximize the mutual information or causal attribution between the explanatory subgraph/structure and the prediction. Subsequently, the network is employed to collectively generate explanations, like Gem [12], OrphicX [13] and our method PAGE. Some methods give explanations by developing interpretable surrogate models, which are subsequently employed to approximate the behavior of the target model to be explained. The explanations derived from these surrogate models are then used to interpret the behavior of the target GNN, such as GraphLime [8], PGM-Explainer [21], and GraphSVX [6].

Beyond these two types of methods, there are also model-level explanation approaches that seek to explore the patterns of the input that can induce specific behaviors in the GNN. These approaches offer more general and high-level insights into the GNN model. Examples include XGNN [28], GLGExplainer [2] and GNNInterpreter [23]. Additionally, some methods specifically focus on offering counterfactual explanations, such as CF-GNNExplainer, CF2,

<sup>&</sup>lt;sup>1</sup> Code is available at https://github.com/anders1123/PAGE/

and RCExplainer, and some self-explanatory models aim to develop the capability to generate predictions alongside corresponding explanations, like ProtGNN [32].

# 3 Methodology

#### 3.1 Preliminary

#### 3.1.1 Graph neural network

The work of graph neural network mainly depends on two mechanisms: message passing and representation aggregation [20] [10]. A graph G can be described as G=(V,E) with node attribute X and adjacency matrix A. Taking the Graph Convolutional Network(GCN) [10] as an example, the computation process is as follows:

$$h_v^{k+1} = \sigma \left( \sum_{u \in \mathcal{N}(v)} \left( W^k h_u^k \tilde{A}_{uv} \right) \right) \tag{1}$$

where  $h_u^k$  is the representation of node u at the k th layer in GCN and  $\tilde{A} = \hat{D}^{-\frac{1}{2}}\hat{A}\hat{D}^{-\frac{1}{2}}$  is the normalized adjacency matrix.  $\hat{A} = A + I$  is the adjacency matrix of the graph G with self loops added and  $\hat{D}$  is a diagonal matrix with  $\hat{D}_{ii} = \sum_j \hat{A}_{ij}$ .  $\mathcal{N}(v)$  denote the neighbors of node v and  $W^k$  is the weight matrix to be trained of the k-th layer. Eq. (1) represents that node v will collect and aggregate the information of all neighboring nodes, which will be used as the representation of node v in the next layer through the activation function  $\sigma(*)$ . Nodes' representation will be used to perform node-level tasks or complete graph-level tasks by using readout functions.

## 3.1.2 Graph Auto-encoder

Graph auto-encoder [9] is a self-supervised learning method. The model consists of a inference network and a generative network, and we call then the encoder and decoder respectively. It mainly retains the following two procedures: It maps the input sample to the hidden space through the GCN encoder to obtain the low-dimensional representation of the sample features, and then latent features are used to reconstruct the original graph structure through a inner product decoder. This process can be formulated as follows:

$$\hat{A} = \sigma\left(ZZ^{\top}\right), \text{ with } Z = GCN(X, A)$$
 (2)

For graph auto-encoder(GAE), the adjacency matrix reconstructed from Z should be as comparable as possible to the original, so GAE uses cross-entropy as the loss function:

$$\mathcal{L}_{GAE} = E_{q(Z|X,A)}[\log p(A \mid Z)]$$
(3)

q and p are the encoding and decoding functions, respectively. For variational graph auto-encoder(VGAE), Z is obtained by sampling from a Gaussian distribution instead of by a definite function. In VGAE (Variational Graph Autoencoder), the reparameterization technique is employed to facilitate the backpropagation of gradients. For more detailed information, you can refer to the original literature [9]. It uses two GCNs that share the same weight in the first layer to generates the means and variances. An additional Kullback-Leibler divergence more than GAE is employed in the loss function:

$$\mathcal{L}_{VGAE} = E_{q(Z|X,A)}[\log p(A \mid Z)] - KL[q(Z \mid X,A)||p(Z)]$$
(4)

#### 3.1.3 Auto-encoder as an explainer

One application of autoencoder is to learn low-dimensional representations of the input data utilizing the nonlinear expressive power of neural networks. These low-dimensional representations(denoted as Z) in the latent space should be more distinguishable. Based on this principle, we aim to separate the causal feature component (denoted as  $Z_c$ ) related to the model's predictions from the entire features in the latent space and discard the non-causal part(denoted as  $Z_s$ ). This approach avoids separating the causal and non-causal parts in the input spaces as in previous methods. By imposing constraints on the latent space feature along the feature dimensions and setting an appropriate learning objective, we intend to isolate the causal subfeature  $Z_c$  from Z in terms of the feature dimension. Then the inner product decoder can map the sub-matrix  $Z_c$  into an adjacency mask, serving as a substructure for explanation.

#### 3.2 Framework

As shown in Figure 1 , our model mainly consists of three components: the encoder (denoted as  $f(\cdot):\mathcal{G}\mapsto\mathcal{Z}$ ), decoder (denoted as  $g(\cdot):\mathcal{Z}\mapsto\mathcal{G}$ ), and discriminator (denoted as  $h(\cdot):\mathcal{Z}\mapsto\mathcal{Y}$ ). We mark the GNN to be explained as  $\mathcal{F}(\cdot):\mathcal{G}\mapsto\mathcal{Y}$ , which gives a predicted label  $Y\in\mathcal{Y}$  for each input graph  $G\in\mathcal{G}$ . The model learns the low-dimensional representation of the input graph through a two-layer GCN encoder. In order to maintain the structural consistency of the reconstructed graph with the original input, latent features  $Z=f(A,X)\in\mathcal{Z}$  is employed to calculate the autoencoder reconstruction loss. Meanwhile, we partition Z into causal and non-causal parts based on the feature dimensions(denoted as  $Z=cat(Z_c,Z_s)$ ).  $Z_c$  will be concatenated with a zero matrix to restore the original dimensions of Z, and then utilized to generate an adjacency mask, serving as the explanation subgraph(denoted as  $G_S,G_S=(A\odot g(Z_c),X)$ ). Then we have the following assumptions:

**Proposition 1.** For any  $Z_p$  that  $Z_p \subseteq Z$ ,  $I(Y; Z_p) \leq I(Y; Z)$ 

Proof.  $I(Y;Z_p) \leq I(Y;Z)$  is equivalent to  $H(Y|Z_p) \geq H(Y|Z)$  where H denotes entropy, since H(Y) is a constant.  $Z_p \subseteq Z$  means  $Z_p$  is a subset of Z, i.e.,  $Z = Z_p \cup Z_+$ . Therefore,  $H(Y|Z_p) \geq H(Y|Z_p,Z_+) = H(Y|Z)$ , i.e.,  $I(Y;Z_p) \leq I(Y;Z)$ 

**Proposition 2.** For causal features in the latent space, they should be strongly correlated with the model predictions (i.e., mutual information  $I(Y; Z_c) = I(Y; Z)$ ). Then the training criterion of the encoder should be:

$$\theta_E^* = \underset{\theta_E}{\operatorname{argmax}} I(Y; Z_c), \text{ since } I(Y; Z_c) \le I(Y; Z)$$
 (5)

**Proposition 3.** For non-causal features in the latent space, they should be independent of the model predictions( $Z_s \perp Y$ , i.e., mutual information  $I(Y; Z_s) = 0$ ).

The above hypotheses can only hold true when the parameter of encoder  $\theta_E$  are optimal. Similar to previous work, we refer to the features in the latent space that are correlated with the predictions made by the GNN being explained as causal features (denoted as  $Z_c$ ), and the parts that are not correlated as non-causal (denoted as  $Z_s$ ). In other words, these features are the causal attribution behind the predictions made by the GNN model on the original input graph G.

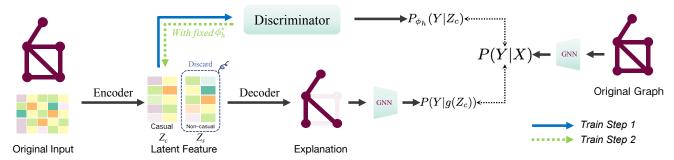


Figure 2. Basic framework of PAGE. The causal part of the latent features related to GNN prediction is used to generate explanations, while the non-causal part is discarded. An additional discriminator is employed to maximize the mutual information between causal features and prediction results. The training process consists of two steps: the first step (indicated by the blue arrows) involves training the autoencoder and the discriminator, and the parameters of the discriminator will be fixed to constrain the encoder in generating better causal features in the second step.

However, the mutual information  $I(Y;Z_s)$  and  $I(Y;Z_c)$  cannot be directly calculated, since  $I(Y;Z_c)=\mathbb{E}\left[\log\frac{P(Y|Z_c)}{P(Y)}\right]$  and  $P(Y|Z_c)$  is unknown (We can only obtain the distribution  $P(Y|G_s)$ ) given by the target GNN). Following Proposition 1 in GAN [7], we introduce an optimal discriminator  $h:\mathcal{Z}\mapsto\mathcal{Y}$  with parameter  $\phi_h$  to approximate  $P(Y|Z_c)$  as  $P_{\phi_h}(Y|Z_c)$ . We then have the following proposition:

**Definition 1.** For  $\theta$  fixed, the optimal discriminator  $\phi_h^*$  is:

$$\phi_h^* = \underset{\phi_h}{\operatorname{argmax}} \mathbb{E}\left[\log P(Y|Z_c)\right] \tag{6}$$

**Proposition 4.** Denoting Kullback-Leibler divergence as  $KL[\cdot||\cdot]$ , for  $\theta$  fixed, the optimal discriminator  $\phi_h$  is  $\phi_h^*$ , s.t.:

$$KL[P(Y|g(Z_c))||P_{\phi_h}(Y|Z_c)] = 0$$
 (7)

The proof of Proposition 4 is as follows:

Proof. According to the definition of mutual information, we have:

$$I(Y; Z_c) = \mathbb{E} \left[ \log P_{\phi_h}(Y|Z_c) \right] + H(Y) + KL \left[ P(Y|Z_c) || P_{\phi_h}(Y|Z_c) \right]$$
(8)

For  $\theta$  fixed,  $I(Y; Z_c)$  and H(Y) are constant, so we have:

$$\phi_h^* = \underset{\phi_h}{\operatorname{argmin}} KL\left[P(Y|Z_c) \mid\mid P_{\phi_h}(Y|Z_c)\right] \tag{9}$$

Thus, by disregarding the constant H(Y), the training criterion becomes:

$$\theta_E^* = \underset{\theta_E}{\operatorname{argmax}} \left\{ \underset{\phi_h}{\operatorname{max}} \mathbb{E} \left[ \log P_{\phi_h}(Y|Z_c) \right] \right\}$$
 (10)

It is difficult to directly find the optimal parameters that satisfy the above expression, so we propose a two-step training strategy: In the first step, the discriminator is trained to make predictions consistent with the target GNN based on the causal sub-features  $Z_c$  in the latent space. The discriminator is a two-layer Multi-Layer Perceptron (MLP). More details of our model can be found in the Appendix [16]. Since the latent features are dimensionality-reduced representations,

the discriminator can achieve high accuracy. In the second step, we fix the parameters of the discriminator to enforce the causal features to be distributed as much as possible in the dimensions we require. Then, the features are used by the decoder to generate explanations based on these causal features.

#### 3.3 Learning Objectives

Post-hoc explanation refers to the process of providing an explanation or justification for a well-trained model. The outputs of GNN  $\mathcal{F}$  are regarded as the training label (formalized as  $Y = \mathcal{F}(A, X)$ ). The learning objective is designed to maximizing the mutual information between the model predictions and the underlying structure  $G_S$ :

$$\max_{G_s} I(Y; G_s) = \max_{G_s} \left\{ H(Y) - H(Y \mid G = G_s) \right\}$$
 (11)

Our training process consists of two stages. In the first stage, the discriminator and auto-encoder are trained together. The motivation is to ensure that the auto-encoder can completely rejuvenate the original graph structure and train the discriminator to learn the causal attribution from causal features to the predicted outcome. The learning objective can be formulized as:

$$\mathcal{L} = \mathcal{L}_{AE} + \lambda_1 * KL [P_{\phi_h}(Y|Z_c) || P(Y|g(Z_c))] + \lambda_2 * KL [P(Y|g(Z_c)) || P(Y|X)]$$
(12)

 $\mathcal{L}_{AE}$  is the loss of the autocoder (GAE or VGAE).

In the second stage, the parameters of the trained discriminator are fixed. The learning objective of the second stage is:

$$\mathcal{L} = \mathcal{L}_{AE} + \mathcal{L}_{size} + \lambda_3 * KL \left[ P(Y|g(Z_c)) || P_{\phi_{h_{fixed}}}(Y|Z_c) \right]$$
(13)

$$\mathcal{L}_{\text{size}} = \left| \frac{\|A * g(Z_c)\|_1}{\|A\|_1} - \gamma \right|$$
 (14)

 $\mathcal{L}_{\mathrm{size}}$  is the size loss to ensure that the mask generated by the decoder are within a reasonable range. In our experiment, the value of  $\gamma$  is set to 0.5 in our experiments. A is the adjacency matrix and  $\gamma$  is a hyper parameter. Causal features  $Z_c$  are transformed into masks of adjacency matrix format by the inner product decoder, which should be multiplied with the original adjacency matrix to serve as explanations.

## 4 Evaluation

#### 4.1 Datasets

To verify the effectiveness of our proposed model, we conducted many experiments on various datasets. Like other GNN interpretation methods, we employed widely used synthetic and real-world datasets. For the node classification task, we use **BA-Shapes** and **Tree-Cycles** presented in [27] with gound-truth explanations, and for the graph classification task, we use real-world datasets **MUTAG** [5] and **NCI1** [22]. More details can be found in Appendix [16].

#### 4.2 Metrics

#### 4.2.1 Fidelity

Fidelity is a commonly used metric to evaluate the faithfulness of the explanations to the model, which is defined as the difference of predicted probability/accuracy between the original predictions and the new predictions of masked input features given by the explainer [18] [30]. Intuitively, the local important input features identified by the interpreter are discriminative to the GNN model. In that case, the model's prediction should change significantly when these local features are removed. We measure this changement with the fidelity score. Analogously, keeping only discriminative features should lead to similar predictions as the original, even if we remove the other unimportant features. We measure that variation by the infidelity metric. In our experiment, we used fidelity based on the predicted probability to verify how much the model can fit the behavior of the original GNN, and it is computed as:

Fidelity prob = 
$$\frac{1}{N} \sum_{i=1}^{N} \left( \mathcal{F} \left( \mathcal{G}_{i} \right)_{y_{i}} - \mathcal{F} \left( \mathcal{G}_{i}^{1-m_{i}} \right)_{y_{i}} \right)$$
(15)

$$\text{Infidelity }^{\text{prob}} = \frac{1}{N} \sum_{i=1}^{N} \left( \mathcal{F} \left( \mathcal{G}_{i} \right)_{y_{i}} - \mathcal{F} \left( \mathcal{G}_{i}^{m_{i}} \right)_{y_{i}} \right) \qquad (16)$$

Here  $\mathcal{G}_i$  is the original graph and  $\mathcal{F}$  is the GNN model to be explained.  $\mathcal{G}_i^{1-m_i}$  represents the new graph obtained by keeping features of  $\mathcal{G}_i$  based on the complementary mask  $1-m_i$ .  $\mathcal{G}_i^{m_i}$  is the new graph by keeping important features of  $\mathcal{G}_i$  based on hard mask  $m_i$  of the explanation.

#### 4.2.2 Accuracy

For synthetic datasets with ground truths, we can leverage the underlying rules of building these datasets to identify important edges or nodes, such as motifs of the graph. Using these important features as references, we can compare the explanations generated by the explainer with the ground truth. Accuracy, ROC, and F1 scores are commonly employed metrics for such evaluations.

The model accuracy measures the predictive accuracy of the generated explanations in relation to the original inputs. Specifically, we input both the original graph and the explanation into the target GNN model and compare the resulting predictions. A higher model accuracy indicates a closer alignment between the explanatory subgraph and the predicted outcomes of the original inputs. This demonstrates that the explainer is more proficient at identifying the most relevant subgraph for the pre-trained GNN.

#### 4.2.3 Sparsity

Sparsity measures the conciseness of explanations, as different explanation methods yield various forms of explanations. Some methods may select a subset of edges or nodes as explanations, while others may generate global importance scores for edges or nodes. Sparsity calculation refers to the proportion of explanations (e.g., edges or nodes) compared to the original input graph. The accuracy and fidelity of explanations are often closely related to sparsity. When sparsity is low, meaning the explanations are more comprehensive and closer to the original input, the accuracy and faithfulness of those explanations tend to be higher. To ensure a fair comparison, it is necessary to set an appropriate threshold that guarantees the sparsity of explanations remains within the same order of magnitude.

# 4.3 Baselines

We consider several baseline models, including perturbation-based method (GNNExplainer), parameterized model-based method (PG-Explainer), and generative models based on autoencoder (Gem and OrphicX). PGExplainer, OrphicX, and Gem are all methods that train an interpreter to explain a target GNN model. On the other hand, GN-NExplainer requires multiple perturbation processes on each input to generate explanations. Gem and OrphicX are the closest baselines to our proposed method in this paper. We set the hyperparameters of these baseline models according to the reported settings in their respective papers.

## 4.4 Quantitive Analysis

K		BA	-SHAI	PES		TREE-CYCLES				
# of edges	5	6	7	8	9	6	7	8	9	10
GNNExp.	67.6	82.4	82.4	88.2	85.3	64.3	66.5	74.3	88.6	97.1
PGExp.	53.4	59.5	60.8	65.5	68.5	76.2	81.5	91.3	95.4	97.1
Gem	64.7	76.4	89.5	91.1	91.1	74.2	85.7	100	100	100
OrphicX	61.7	61.7	73.5	76.4	76.4	74.2	82.8	97.1	100	100
OrphicX-0	61.7	73.5	73.5	76.4	76.4	74.2	85.7	97.1	97.1	100
PAGE-GAE	76.4	76.4	76.4	89.5	89.5	74.2	82.8	97.1	97.1	100
PAGE	91.1	91.1	91.1	91.1	91.1	74.3	88.6	100	100	100

Table 1. Explanation Accuracy on Synthetic Datasets (%)

1-Sparsity	Mutagenicity   0.5   0.6   0.7   0.8   0.9						NCI1				
GNNExp.	65.0	66.6	66.4	71.0	78.3	64.2	65.7	68.6	75.2	81.8	
PGExp.								65.2			
Gem	66.4	67.7	71.4	76.5	81.8	61.8	65.6	70.6	74.9	83.9	
OrphicX	61.9	66.6	72.8	76.0	80.1	59.1	60.8	67.1	76.3	79.3	
OrphicX-0	65.6	66.1	70.8	73.5	78.5	58.1	61.3	66.4	72.7	79.1	
PAGE-GAE	61.7	66.9	70.5	76.4	80.4	56.2	60.8	65.1	74.9	79.1	
PAGE	63.5	68.4	72.3	76.5	83.4	58.1	66.4	70.8	77.4	85.4	

Table 2. Explanation Accuracy on Real-World Datasets (%)

Datasets	GNNExp.	PGExp.	Gem	OrphicX	PAGE
BA-SHAPES	90.8	76.3	65.4	94.2	94.2
TREE-CYCLES	91.2	77.2	76.3	94.2	96.7

Table 3. Edge Accuracy of Explanation (%)

Similar to previous works [30], our evaluation of the explainer's performance primarily focuses on two performance metrics: model

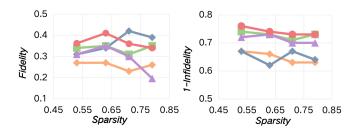


Figure 3. Fidelity and 1-Infidelity vs. Sparsity on BA-Shapes.

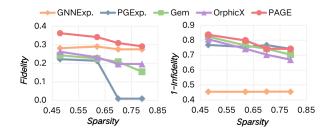


Figure 4. Fidelity and 1-Infidelity vs. Sparsity on Mutagenicity.

accuracy and fidelity. Accuracy reflects how closely the GNN's explanations align with the original inputs. Table 1 and 2 respectively illustrate the performance of different explanation models on artificial synthetic datasets and real-world datasets under different sparsity constraints.

It should be noted that for the artificial synthetic dataset, we used a constraint of the top K edges for different models, instead of sparsity. This practice aligns with conventions in previous works. A unified K allows us to compare the concordance between explanations and ground truths. We didn't choose a value smaller than 5 for K because explanations formed by selecting fewer than 5 edges wouldn't comprehensively cover the ground-truth motifs, and such explanations would lack meaningful context. For instance, for the Tree-cycle dataset, choosing any k edges (k<5) from a cycle motif for explanation would not make any sense.

According to Table 1 and 2, we observe that the accuracy of our explanations surpasses that of baseline models both on the synthetic datasets and real-world datasets. For the artificial synthetic dataset, our explanations achieve optimal accuracy under different constraints of K. For real-world datasets, our method outperforms other approaches under most sparsity constraints and maintains a higher average accuracy. We conducted experiments using both GAE and VGAE as interpreters. The experimental results indicated that the performance of GAE-based model was surpassed by the VGAEbased. Consequently, for the subsequent experiments, we consistently chose VGAE as the foundational model for our interpreter. In addition, OrphicX-0 in the table represents the results obtained by removing the calculation of information flow from the loss function of OrphicX. The table shows that the information flow within OrphicX does not effectively capture the causal effects between hidden features and model predictions. Even after it was removed, the performance of the explainer did not significantly decline.

Indeed, different explanations could potentially lead to similar classification results as the original samples. Therefore, to further compare the performances among different explaining methods, we reported the evaluation of edge accuracy on the synthetic dataset. By transforming the inclusion of each edge in the explanation into a binary classification problem, we can assess the concordance between the generated edges and the ground truth motifs. The results are presented in Table 3. A higher accuracy indicates that the explainer tends to assign higher importance scores to edges of the ground truth motifs/subgraphs. Experimental results demonstrate that explanations generated by our approach are superior to the others.

Furthermore, we use the fidelity and infidelity metrics to compare the quality of explanations generated by different methods under four different sparsity levels. Fidelity and Infidelity metrics can provide an alternative perspective on the quality of explanations. We compare our model and baseline methods on the BA-Shapes and Mutagenicity datasets, as shown in Figure 4 and 3 (using 1-Infidelity as the y-axis for easy comparison), indicating that our method generally exhibits better fidelity in most cases.

## 4.5 Qualitative Analysis

To further demonstrate the interpretability of the explanations, Figure 5 illustrates some visualized explanation instances generated by different methods on the Mutagenicity dataset. The figure presents three different mutagenic molecules. The first column represents the initial molecule, while the remaining five columns showcase explanations generated by different methods, all adhering to a common sparsity constraint. Black edges indicate the edges the interpreter selects, while gray edges signify those not selected. The value of p represents the probability that the initial molecule (or the explained) possesses mutagenic potential given by the target GNN. In the first two cases, our method is the only one capable of capturing the complete carbon ring and the connected -NO2 functional group attaching to it, which is always regarded as a feature of being mutagenic. In the second case, OrphicX failed to recognize the -NH2 functional group, and the explanation is predicted as non-mutagenic. In the third example, we report a mutagenic molecule with no explicit motifs. PAGE produced an explanation that most closely matched the original prediction. Such explanations could aid in uncovering novel chemical patterns. In summary, compared to baseline models, our approach generates explanations that best reflect the behavior of the target GNN.

#### 4.6 Ablation Study

	BA-SHAPES(top k) 6   7   8	Mutagenicity(1-Sparsity) 0.6   0.7   0.8					
None	59.5   59.5   59.5	61.9   62.4   65.2					
ERM	76.4   88.5   88.5	64.6   69.7   75.4					
No Pre-training	77.3   89.5   89.5	65.5   69.7   76.7					
FULL	91.1   91.1   91.1	66.4   70.8   77.4					

**Table 4.** Ablation Study Result. (1, None: only autoencoder loss and size loss.2, ERM: the empirical term  $KL\left[P(Y|g(Z_c))||P(Y|X)\right]$  added.3, No pre-training: do not pre-train the discriminator.4, FULL: fully trained.)

To further investigate the effectiveness of designed components in PAGE, we conducted ablation experiments on both the artificially synthesized dataset BAShapes and the real-world dataset Mutagenicity. Specifically, we validated the performance of the interpreter under four conditions. The result is shown in Table 4, showing that

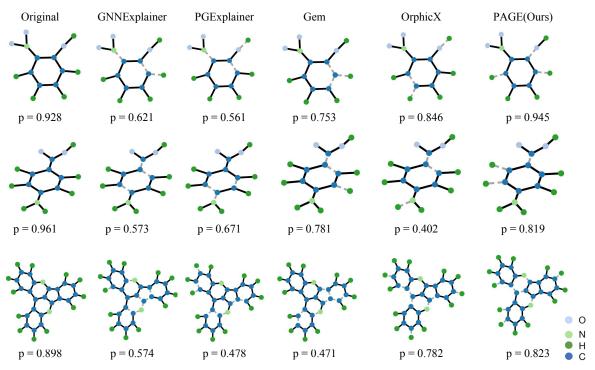


Figure 5. Visualized explanations on Mutagenicity. The explanation results of different methods are highlighted with black edges, where the gray edges are regarded as non-casual parts for the prediction. "P" under each graph/subgraph denotes the probability of being classified into the mutagenic class, which is obtained by feeding the associated graph/subgraph into the target GNN.

applying only partial component cannot obtain optimal performance. Additionally, we conducted experiments on hyperparameter sensitivity. Please refer to the supplementary materials for more details.

## 4.7 Efficiency Study

time(ms)		GNNExp.		PGExp.	Gem	OrphicX	PAGE
training		-		340,148	24,041	368,718	78,233
inference		533,681		354	79	77	91
total	T	533,681	1	340,502	24,120	368,795	78,324

Table 5. Training and inference time(ms) on Mutagenicity.

GNNExplainer requires multiple perturbations on a single sample. PGExplainer necessitates generating soft masks individually for each edge. GEM involves considering each edge to obtain a subgraph for the "guidance" of the training. OrphicX computes information flow through sampling at different scales. Compared to these methods, our approach eliminates any perturbation or sampling processes. Each inference and backpropagation involve only a single computation, resulting in a time complexity of O(1). As a result, our method remains significantly more efficient than these baseline models. Table 4 presents the training and inference times for the mentioned models on the Mutagenicity dataset. All models are configured with hyperparameters, learning rates, and epoch numbers as described in their original papers. It needs to be specified that, although experiments demonstrate that Gem has shorter training and inference times compared to ours, Gem requires an additional distillation process to generate guidances for the training of the autoencoder. This distillation incurs significant costs (more than 150,000 ms on the Mutagenicity dataset). Consequently, considering the overall expenses, our method's efficiency still surpasses all other baseline models.

# 5 Conclusion

In this article, we introduce PAGE, a parametric generative Graph Neural Network (GNN) explaining method designed to generate concise and reliable causal explanations for any graph neural network. PAGE optimize a generative autoencoder with a learning objective of maximizing mutual information between latent features and outputs. Compared to existing methods, PAGE offers several advantages: its computation and inference processes do not require any perturbation or sampling processes, ensuring high explanation accuracy while maintaining greater efficiency than previous approaches. Moreover, as a model-agnostic post-hoc explanation approach, it can offer causal explanations for different types of GNNs without relying on any prior assumptions or internal model details. We demonstrate the superiority of our approach over baseline models through experiments and data analysis. A limitation of PAGE is that we only explored autoencoders as interpreters. An avenue for potential improvement could involve using more powerful generative models as interpreters. We leave this for future investigation.

#### Acknowledgements

This work is supported by National Natural Science Foundation of China under grants 62376103, 62302184, 62206102, Science and Technology Support Program of Hubei Province under grant 2022BAA046, and CCF-AFSG Research Fund.

## References

- N. Ay and D. Polani. Information Flows in Causal Networks. Advances in Complex Systems, 11(01):17–41, 2008.
- [2] S. Azzolin, A. Longa, P. Barbiero, P. Liò, and A. Passerini. Global explainability of gnns via logic combination of learned concepts. arXiv preprint arXiv:2210.07147, 2023.
- [3] F. Baldassarre and H. Azizpour. Explainability Techniques for Graph Convolutional Networks. *arXiv preprint arXiv:1905.13686*, 2019.
- [4] S. L. Bressler and A. K. Seth. Wiener–Granger Causality: A well established methodology. *NeuroImage*, 58(2):323–329, 2011.
- [5] A. K. Debnath, R. L. Lopez De Compadre, G. Debnath, A. J. Shuster-man, and C. Hansch. Structure-activity relationship of mutagenic aromatic and heteroaromatic nitro compounds. Correlation with molecular orbital energies and hydrophobicity. *Journal of Medicinal Chemistry*, 34(2):786–797, 1991.
- [6] A. Duval and F. D. Malliaros. GraphSVX: Shapley Value Explanations for Graph Neural Networks. In N. Oliver, F. Pérez-Cruz, S. Kramer, J. Read, and J. A. Lozano, editors, *Machine Learning and Knowledge Discovery in Databases. Research Track*, pages 302–318, Cham, 2021. Springer International Publishing. ISBN 978-3-030-86520-7.
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative Adversarial Nets. In Advances in Neural Information Processing Systems, volume 27. Curran Associates, Inc., 2014.
- [8] Q. Huang, M. Yamada, Y. Tian, D. Singh, and Y. Chang. GraphLIME: Local Interpretable Model Explanations for Graph Neural Networks. *IEEE Transactions on Knowledge and Data Engineering*, 35(7):6968–6972, 2023.
- [9] T. N. Kipf and M. Welling. Variational Graph Auto-Encoders. arXiv preprint arXiv:1611.07308, 2016.
- [10] T. N. Kipf and M. Welling. Semi-Supervised Classification with Graph Convolutional Networks. *arXiv preprint arXiv:1609.02907*, 2017.
- [11] P. Li, Y. Yang, M. Pagnucco, and Y. Song. Explainability in Graph Neural Networks: An Experimental Survey. arXiv preprint arXiv:2203.09258, 2022.
- [12] W. Lin, H. Lan, and B. Li. Generative Causal Explanations for Graph Neural Networks. In *Proceedings of the 38th International Conference on Machine Learning*, pages 6666–6679. PMLR, 2021.
  [13] W. Lin, H. Lan, H. Wang, and B. Li. OrphicX: A Causality-Inspired Learn Virginia M. Lin, Physics Phys
- [13] W. Lin, H. Lan, H. Wang, and B. Li. OrphicX: A Causality-Inspired Latent Variable Model for Interpreting Graph Neural Networks. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 13719–13728, New Orleans, LA, USA, 2022. IEEE. ISBN 978-1-66546-946-3.
- [14] D. Luo, W. Cheng, D. Xu, W. Yu, B. Zong, H. Chen, and X. Zhang. Parameterized Explainer for Graph Neural Network. In *Advances in Neural Information Processing Systems*, volume 33, pages 19620–19631. Curran Associates, Inc., 2020.
- [15] P. E. Pope, S. Kolouri, M. Rostami, C. E. Martin, and H. Hoffmann. Explainability Methods for Graph Convolutional Neural Networks. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 10764–10773, Long Beach, CA, USA, 2019. IEEE. ISBN 978-1-72813-293-8.
- [16] Y. Qiu. PAGE: Parametric Generative Explainer for Graph Neural Network, Aug. 2024. URL https://doi.org/10.5281/zenodo.13375562.
- [17] M. Rathee, T. Funke, A. Anand, and M. Khosla. BAGEL: A Benchmark for Assessing Graph Neural Network Explanations. arXiv preprint arXiv:2206.13983, 2022.
- [18] M. T. Ribeiro, S. Singh, and C. Guestrin. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16, pages 1135–1144, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 978-1-4503-4232-2.
- [19] C. Shan, Y. Shen, Y. Zhang, X. Li, and D. Li. Reinforcement Learning Enhanced Explainer for Graph Neural Networks. Advances in Neural Information Processing Systems, pages 22523–22533, 2021.
- [20] N. Shervashidze, P. Schweitzer, É. J. van Leeuwen, K. Mehlhorn, and K. M. Borgwardt. Weisfeiler-lehman graph kernels. *Journal of Machine Learning Research*, 12(77):2539–2561, 2011.
- [21] M. Vu and M. T. Thai. PGM-Explainer: Probabilistic Graphical Model Explanations for Graph Neural Networks. Advances in Neural Information Processing Systems, 33:12225–12235, 2020.
- [22] N. Wale, I. A. Watson, and G. Karypis. Comparison of descriptor spaces for chemical compound retrieval and classification. *Knowledge and In*formation Systems, 14(3):347–375, 2008.
- [23] X. Wang and H.-W. Shen. GNNInterpreter: A Probabilistic Generative Model-Level Explanation for Graph Neural Networks. arXiv preprint arXiv:2209.07924, 2023.

- [24] X. Wang, Y. Wu, A. Zhang, F. Feng, X. He, and T.-S. Chua. Reinforced Causal Explainer for Graph Neural Networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(2):2297–2309, 2023.
- [25] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu. A Comprehensive Survey on Graph Neural Networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2021.
- [26] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How Powerful are Graph Neural Networks? arXiv preprint arXiv:1810.00826, 2019.
- [27] Z. Ying, D. Bourgeois, J. You, M. Zitnik, and J. Leskovec. GNNExplainer: Generating Explanations for Graph Neural Networks. In Advances in Neural Information Processing Systems, volume 32. Curran Associates, Inc., 2019.
- [28] H. Yuan, J. Tang, X. Hu, and S. Ji. XGNN: Towards Model-Level Explanations of Graph Neural Networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 430–438, Virtual Event CA USA, 2020. ACM. ISBN 978-1-4503-7998-4.
- [29] H. Yuan, H. Yu, J. Wang, K. Li, and S. Ji. On Explainability of Graph Neural Networks via Subgraph Explorations. In *Proceedings of the 38th International Conference on Machine Learning*, pages 12241–12252. PMLR, 2021.
- [30] H. Yuan, H. Yu, S. Gui, and S. Ji. Explainability in Graph Neural Networks: A Taxonomic Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):5782–5799, 2023.
- [31] S. Zhang, Y. Liu, N. Shah, and Y. Sun. GStarX: Explaining Graph Neural Networks with Structure-Aware Cooperative Games. Advances in Neural Information Processing Systems, 35:19810–19823, 2022.
- [32] Z. Zhang, Q. Liu, H. Wang, C. Lu, and C. Lee. ProtGNN: Towards Self-Explaining Graph Neural Networks. Proceedings of the AAAI Conference on Artificial Intelligence, 36(8):9127–9135, 2022.
- [33] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, and M. Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.

## A Experiments details in Section 3.

#### A.1 Datasets

BA-Shapes and Tree-cycles are synthetic node classification datasets. BA-Shapes is composed of 80 house motifs and a base BarabasiAlbert (BA) graph containing 300 nodes. Node labels are divided into four categories, representing vertices, middle, and bottom parts of the houses, or not belonging to any motifs. Tree-cycles consist of a base eight-level binary balanced trees with 80 six-node cyclic motifs. Node labels are binary and denote whether the node belongs to the cycle motif. Mutagenicity and NCI1 are real-world graph classification datasets. Mutagenicity comprises 4337 unique chemical molecules, where nodes represent individual atoms and edges denote different chemical bonds. The labels indicate whether the chemical molecule is mutagenic. NCI1 includes 4110 chemical compounds with labels indicating whether the compound inhibits cancer cell growth. Table 5 provides additional details about these datasets.

# A.2 Experimental Hardware

All experiments were conducted on a PC equipped with an NVIDIA RTX 4070 Ti GPU and an Intel Core i5-13600KF processor, supported by 32GB of RAM and 12GB of graphics memory.

## A.3 Details about the target GNN

For the target GNN to be explained, we followed the same experimental setup as previous works. Specifically, for node classification, we employed three layers of Graph Convolutional Networks (GCNs) with output dimensions set to 20. We concatenated the outputs from these three layers and subsequently subjected them to a linear transformation to derive the node labels. In the case of graph classification, we utilized three layers of GCNs with dimensions of 20 and conducted global max-pooling to generate the graph representations. A subsequent linear transformation layer was employed to derive the graph labels. The target GNN achieved accuracies of 94.1, 97.1, 88.5, and 78.6 on the BA-Shapes, Tree-cycles, Mutagenicity, and NCI1 datasets, respectively.

## A.4 Details about the explainer

For the encoder of explainer, we applied a three-layer GCN with output dimensions 32, 32, and 16. The decoder is equipped with a two-layer MLP and an inner product decoder. The discriminator was implemented using a two-layer MLP with a hidden layer dimension of 32. We trained the explainers using a learning rate of 0.003 for 300 epochs in the first stage and 50 epochs in the second stage. All of our experiments and models, including the target GNN, our interpreter, and the baseline models, were implemented using PyTorch and trained with Adam optimizer.

#### A.5 Downstream Tasks

**Node-Level Task** For node-level tasks, such as node classification, the interpretation of the target node is designed as the subgraph with the strongest correlation with the prediction result on the computation graph. For a graph with more than thousands of nodes, it is impractical to reconstruct the complete graph for interpreting a single node; Moreover, the auto-encoder can not complete the training based on only one single graph. Therefore, for node-level tasks, we will obtain the subgraph composed of *k*-hop neighbors of each node

Datasets	BA-Shapes	Tree-cycles	Mutagenicity	NCI1
#Graphs	1	1	4337	4110
#Nodes	700	871	30.32(Avg.)	32.30(Avg.)
#Edges	4110	1950	30.77(Avg.)	29.87(Avg.)
#Labels	4	2	2	2
#Training	300	270	3468	3031
#Validation	50	45	434	411
#Testing	50	45	434	410

**Table 6.** Details about the datasets.

for training, leading to a compact study on the whole graph (The value of k is set as the number of layers in the target GNN.). When interpreting a specific node, we will use the subgraph consisting of no more than k-hop reconstructed by the encoder for prediction, and these subgraphs are considered as the most relevant and influential part of the input graph to the prediction.

**Graph-Level Task** For graph-level tasks, such as graph classification, we will let the auto-encoder learn to reconstruct the complete graph structure. When making predictions, a readout function (can be average pooling or maximal pooling) is employed to obtain the graph-level representation from the latent feature Z. Then the discriminator employs the graph representation to make predictions. Experimental results show that this austere readout function is efficient and effective under the constraints of our designed optimizing objective.

#### A.6 Setting of sparsity hyperparameter $\gamma$

Regarding the setting of sparsity, although we recommend setting the sparsity hyperparameter  $\gamma$  to 0.5 to generate an importance score matrix that is relatively balanced, we still suggest adjusting the sparsity differently based on different interpretability requirements to achieve better interpretability (for example, when a concise explanation is needed, lowering the sparsity as much as possible).

## A.7 Hyperparameter sensitivity.

To analyze the impact of different hyperparameter settings on the performance of PAGE, we conducted hyperparameter sensitivity experiments on Mutagenicity. We separately varied the values of the hyperparameters  $\lambda_1,\,\lambda_2,$  and  $\lambda_3$  while keeping the default parameter settings ( $\lambda_1=1.0,\,\lambda_2=0.1,\,\lambda_3=0.01$ ), and observed the explanation accuracy under different sparsity levels. The experimental results, as shown in Figure 6, indicate that all three hyperparameters have varying influence on the model's performances. This not only demonstrates the effectiveness of the discriminator module we designed but also underscores the necessity of selecting appropriate values for hyperparameters.

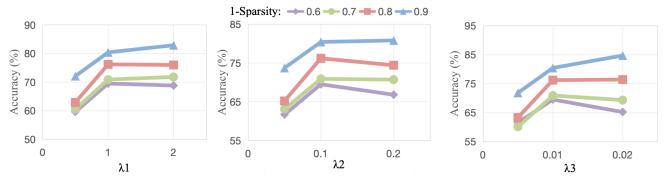


Figure 6. Results about hyperparameter sensitivity.