

## A THE USE OF LARGE LANGUAGE MODELS (LLMs)

In preparing this paper, Large Language Models (LLMs) were used solely as general-purpose writing assistants. Specifically, we employed LLMs to polish the text by carefully correcting grammar, spelling, and occasionally translating our original sentences into fluent academic English when necessary. Importantly, the LLMs were not used to generate new research ideas, formulate novel claims, or alter the meaning of our original sentences. All outputs generated by LLMs were carefully reviewed and verified by the authors for consistency with our intended meaning and to avoid any issues of plagiarism or scientific misconduct.

## B RELATED WORK

**GNN Explainability and Out-Of-Distribution.** Current GNN explanatory methods (Shin et al., 2024; Wang & Shen, 2022; Yu & Gao, 2024; Zhang et al., 2022; Wu et al., 2023; Li et al., 2023) aim to enhance model decision transparency through model-level or instance-level methods (Yuan et al., 2022). Model-level methods (Yuan et al., 2020) enhance GNN explainability by generating input-independent explanations, while instance-level methods (Yuan et al., 2021) identify input features most relevant to the model prediction. In this work, we focus on post-hoc instance-level methods. However, such methods face growing OOD challenges (Zheng et al., 2023; Fang et al., 2023). While some works (Fang et al., 2024; Faber et al., 2020) have already made attempts, but they require that the model be retrained. Separately, traditional mixup-based data augmentation methods (Zhang et al., 2023; Chen et al., 2024; Zhang et al., 2024) also alleviate the OOD problem, while still struggle to ensure the quality of mixup graphs for approximating the original distribution.

**Graph Contrastive Learning.** Graph Contrastive Learning (GCL) has become a dominant paradigm for graph self-supervised learning, aiming to learn discriminative label-independent representations by maximizing mutual information across different views. For example, GraphCL (You et al., 2020) leverages various data augmentation strategies to enhance the robustness and transferability of unsupervised representations. AutoGCL (Yin et al., 2022) employs learnable view generators, guided by auto-augmentation, to produce semantically consistent yet diverse graph views. SimGRACE (Xia et al., 2022) generates contrastive views via perturbed encoders for representation alignment. Among current post-hoc explanatory methods, RegExplainer (Zhang et al., 2024) is a pioneering work that introduces contrastive learning in the graph regression task. However, a unified task-agnostic explanation framework remains unexplored.

**Graph Pooling.** Graph pooling methods typically capture global information through hierarchical compression to obtain generalizable and expressive representations. Some methods learn a score function from node features and select the top  $r$  nodes to form a new subgraph (Gao & Ji, 2019). For example, SAGPooling (Lee et al., 2019) refines node scoring by integrating structural context via GNNs. GSAPool (Zhang et al., 2021) enhances pooling process by dynamically fusing node features and topology information. DMIPool (Zhao et al., 2023) further captures multi-level dependencies from dual-view representations to improve robustness. MLAP (Itoh et al., 2022) preserves cross-layer structural information through layer-wise pooling, enhancing node distinguishability. Our method introduces pooling strategies to apply hard perturbations to graph structure, extracting more distinguishable explanatory and label-irrelevant subgraphs, ultimately enabling the mixup graphs to more faithfully approximate the original distribution.

## C NOTATIONS

In Table 4, we summarized the important notations we used and their descriptions in this paper.

Table 4: Important symbols and notations.

Symbols	Descriptions
$G$	Graph instance
$\mathbf{X}$	Node feature matrix
$\mathbf{A}$	Adjacency matrix
$Y$	Label for $G$
$\mathcal{G}$	Graph dataset
$\mathcal{V}, \mathcal{E}$	Node set and edge set
$\mathcal{Y}$	Label set
$v_i$	The $i$ -th node
$x_i$	The feature of node $v_i$
$n$	Number of nodes in $G$
$\text{edge}(i, j)$	Edge from node $i$ to node $j$
$i$	The $i$ -th node index
$j$	The $j$ -th node index
$d$	Feature dimension
$G^*$	Optimal explanatory subgraph
$Y^*$	Label for $G^*$
$\mathbf{H}^*, \mathbf{A}^*$	Node embeddings and adjacency matrix of $G^*$
$I(\cdot)$	Mutual information
$\alpha$	Hyper parameter for mutual information term in GIB
$\mathcal{N}$	Number of graphs in $\mathcal{G}$
$f(\cdot)$	To-be-explained GNN model
$\mathbf{H}$	Node embeddings matrix
$f_\psi(\cdot)$	Explainer network
$f_{\text{enc}}(\cdot)$	Encoder component of $f$ that generates graph representations
$\mathbf{h}_G$	Graph representation of $G$
$DV(\cdot)$	Donsker–Varadhan variational representation
$T(\cdot)$	Critic function
$G_Y$	The graph paired with label $Y$
$\text{sim}(\cdot, \cdot)$	Dot-product similarity between $\ell_2$ -normalized vectors
$r_l$	Node preservation ratio at layer $l$
$I^L$	Final preserved node indices after $L$ pooling layers
score	Node scores
$p_l$	Learnable projection vector at layer $l$
$N$	Number of sampled neighbors
$G'$	Sampled neighbor
$\mathbf{A}_{\text{our}}^{(\text{mix})}$	Adjacency matrix obtained via structural mixup strategy
$\mathbf{H}_{\text{our}}^{(\text{mix})}$	Node embeddings matrix obtained via structural mixup strategy
$G_{\text{our}}^{(\text{mix})}$	Graph obtained via structural mixup strategy
$S'$	Binary node preserved matrix from $G'$
$\mathbf{W}$	Edge weights matrix
$\mathbf{M}_{\text{our}}^{(\text{mix})}$	Edge weights mask for mixup graph
$\mathcal{L}_{\text{InfoNCE}}$	InfoNCE loss
$\mathcal{L}_{\text{BCE}}$	Binary cross-entropy loss
$\beta$	Hyper parameter balancing $\mathcal{L}_{\text{InfoNCE}}$ and $\mathcal{L}_{\text{BCE}}$

## D ALGORITHMS

**Algorithm 1** Training Explainer

---

**Input:** A graph dataset  $\mathcal{G}$ , pretrained GNN model  $f$ , explainer network  $f_\psi(\cdot)$ .  
**Output:** Trained explainer network  $f_\psi(\cdot)$ .

- 1: Initialize explainer network  $f_\psi(\cdot)$ .
- 2: **for**  $e \in \text{epochs}$  **do**
- 3:   **for**  $G \in \mathcal{G}$  **do**
- 4:     Randomly sample  $N$  graphs  $\{G_i\}_{i=0}^{N-1}$  from graph dataset  $\mathcal{G}$
- 5:     Compare similarity between  $\{G_i\}_{i=0}^{N-1}$  and  $G$  to obtain  $G^+$  and  $\{G_i^-\}_{i=1}^{N-1}$
- 6:      $G_{\text{our}}^{(\text{mix})+} \leftarrow \text{Structural Mixup}(G, G^+)$
- 7:      $\{G_{\text{our},i}^{(\text{mix})-}\}_{i=1}^{N-1} \leftarrow \text{Structural Mixup}(G, \{G_i^-\}_{i=1}^{N-1})$
- 8:     Generate  $M_{(\text{our})}^{(\text{mix})+}$  and  $\{M_{\text{our},i}^{(\text{mix})-}\}_{i=1}^{N-1}$  with  $G_{\text{our}}^{(\text{mix})+}$  and  $\{G_{\text{our},i}^{(\text{mix})-}\}_{i=1}^{N-1}$
- 9:     Compute  $\mathcal{L}_{\text{InfoNCE}}(G, G^+, \{G_i^-\}_{i=1}^{N-1})$  with Equation 10
- 10:     Compute  $\mathcal{L}_{\text{BCE}}$  with Equation 11
- 11:     Compute overall loss  $\mathcal{L}$  with Equation 12
- 12:   **end for**
- 13:   Update  $f_\psi(\cdot)$  with back propagation
- 14: **end for**
- 15: **Return** Trained explainer network  $f_\psi(\cdot)$

---

**Algorithm 2** Graph Structural Mixup Algorithm

---

**Input:** To-be-explained graph  $G, G'$  sampled from graph dataset  $\mathcal{G}$ .  
**Output:** Graph  $G_{\text{our}}^{(\text{mix})}$ .

- 1: Obtain  $I^L$  via  $L$ -layer graph pooling following Equation 8
- 2: Obtain  $(I')^L$  via  $L$ -layer graph pooling following Equation 8
- 3: Generate explanatory subgraphs  $G^*$  and  $(G')^*$  based on  $I^L$  and  $(I')^L$
- 4: Structural Mixup adjacency  $A_{\text{our}}^{(\text{mix})}$  and node embeddings  $H_{\text{our}}^{(\text{mix})}$  matrix with Equation 9
- 5: **Return**  $G_{\text{our}}^{(\text{mix})} = (H_{\text{our}}^{(\text{mix})}, A_{\text{our}}^{(\text{mix})})$

---

To address the distribution shift issue in the GIB objective, a graph-pooling-based structural mixup strategy is employed to generate in-distribution, naturally connected mixup graphs  $G_{\text{our}}^{(\text{mix})}$  that preserve the information of explanatory subgraphs  $G^*$ . Built upon this, contrastive learning is introduced as an optimization objective within GIB to capture underlying representations, thereby approximating the mutual information between  $G^*$  and the label  $Y$  in a task-agnostic manner. For completeness, we present the pseudocode of our method below, which encompasses both explainer training and the structural mixup strategy.

In Algorithm 1, at each epoch, we first randomly sample  $N$  neighbors  $\{G_i\}_{i=0}^{N-1}$  from the graph dataset  $\mathcal{G}$  for each graph  $G$ . Based on the similarity between their representations and that of  $G$ , we identify the positive neighbor  $G^+$  and the negative neighbors  $\{G_i^-\}_{i=1}^{N-1}$ . Then, the structural mixup strategy combines the explanatory subgraph  $G^*$  with the label-irrelevant parts of  $G^+$  and  $\{G_i^-\}_{i=1}^{N-1}$ , producing the positive and negative mixup graphs  $G_{\text{our}}^{(\text{mix})+}$  and  $\{G_{\text{our},i}^{(\text{mix})-}\}_{i=1}^{N-1}$ . For each mixup graph, an MLP is employed to predict the edge weight of each edge, and the structural information obtained from graph pooling process is incorporated to generate the final masks  $M_{\text{our}}^{(\text{mix})+}$  and  $\{M_{\text{our},i}^{(\text{mix})-}\}_{i=1}^{N-1}$ . The explainer is optimized by minimizing the InfoNCE loss Equation 10 and the BCE loss Equation 11, while the overall loss function Equation 12 sums the two objectives with a trade-off parameter. Finally, the explainer parameters  $f_\psi(\cdot)$  are updated through backpropagation with the overall loss. Specifically, Algorithm 2 shows the structural mixup process. Given a to-be-explained graph  $G$  and a sampled neighbor graph  $G'$ , where  $G'$  can be instantiated as the positive neighbor  $G^+$  or one of the negative neighbors  $\{G_i^-\}_{i=1}^{N-1}$ , we perform  $L$ -layer graph pooling to obtain the preserved node indices  $I^L$  and  $(I')^L$ . Then, the explanatory subgraph  $G^*$  is represented as  $G^* = (H^*, A^*)$ , where  $H^* = H(I^L, :)$  and  $A^* = A(I^L, I^L)$ . Similarly, the explanatory subgraph of  $G'$  is obtained in the same way based on  $(I')^L$ . The adjacency matrix  $A_{\text{our}}^{(\text{mix})}$  and node

864 embeddings  $\mathbf{H}_{\text{our}}^{(\text{mix})}$  of the mixup graph are obtained by replacing  $(\mathbf{A}')^*$  with  $\mathbf{A}^*$ , as defined in  
 865 Equation 9. Finally, the mixup graph  $G_{\text{our}}^{(\text{mix})} = (\mathbf{H}_{\text{our}}^{(\text{mix})}, \mathbf{A}_{\text{our}}^{(\text{mix})})$  is returned.  
 866

867 **Computational Complexity Analysis.** Given graph  $G$ , the time complexity of  $L$ -layer graph pool-  
 868 ing is  $O(L \cdot (|\mathcal{E}|d + nd))$ , where  $d$  is the feature dimension,  $\mathcal{E}$  denotes the edge set, and  $n$  is the  
 869 number of nodes in  $G$ . Then, structural replacement replaces the explanatory subgraph via  $m$  node  
 870 indices, with a complexity of  $O(m^2)$ , where  $m$  is the number of final preserved nodes. Edge weights  
 871 for the mixup graph are generated by an MLP and the time complexity is  $O(|\mathcal{E}^{(\text{mix})}| \cdot d)$ , where  $\mathcal{E}^{(\text{mix})}$   
 872 denotes the edge set of the mixup graph. Since  $m \ll n$ ,  $|\mathcal{E}| \ll n^2$ , and the sizes of  $|\mathcal{E}|$  and  $|\mathcal{E}^{(\text{mix})}|$   
 873 are comparable, the overall time complexity of our method is dominated by  $O(L \cdot |\mathcal{E}|d)$ .  
 874

## 875 E FULL EXPERIMENTAL SETUP

877 Detailed experimental settings are provided in this section, including implementation details,  
 878 datasets, and baseline methods. All experiments are conducted on a Linux workstation running  
 879 Ubuntu 22.10 (kernel 5.19.0-46-generic) equipped with 8 NVIDIA GeForce RTX 4090 GPUs (24  
 880 GB each). The system uses NVIDIA driver version 535.86.05 and CUDA 12.2. All code is im-  
 881 plemented in Python 3.9.21, using PyTorch 2.0.1+cu118, PyTorch Geometric 2.6.1, torch-scatter  
 882 2.1.2+pt20cu118, and torch-sparse 0.6.18+pt20cu118. The complete experimental details can be  
 883 found in our anonymous repository: <https://anonymous.4open.science/r/TAME-main-2FB7>.  
 884

### 885 E.1 IMPLEMENTATION DETAILS

886 Following the configuration in previous work (Ying et al., 2019; Zhang et al., 2023), we divide  
 887 each dataset into training, validation, and test sets with a ratio of 0.8, 0.1, and 0.1, respectively.  
 888 We choose GCN as the backbone model for all experiments, as it demonstrates the most stable  
 889 performance across different datasets. The detailed comparison of the prediction performance of  
 890 GCN and other GNN architectures are provided in Appendix H.3. A three-layer GCN is used as  
 891 the target model, and for graph regression tasks, an additional linear layer is appended. We train  
 892 the GCN model to a reasonable performance, as shown in the GCN column of Tables 5 and 6.  
 893 For explanation, we follow the sample selection strategy used in GNNExplainer (Ying et al., 2019)  
 894 and randomly select 200 graph instances per dataset to reduce computational cost. All explana-  
 895 tion methods are optimized using the Adam optimizer with a weight decay of  $5 \times 10^{-4}$  (Kingma,  
 896 2014). Regarding hyperparameters, we apply grid search to determine the loss weights  $\beta$ , and set  
 897 the top- $r$  ratios according to the ground truth explanations. For baseline methods with overlapping  
 898 hyperparameters, we adopt a unified setting; otherwise, we retain their default configurations. We  
 899 ensure consistent learning rates and training epochs across all explanation methods. We evaluate the  
 900 performance of our proposed method against baseline approaches on ground-truth explanation tasks  
 901 using the AUC-ROC score. To quantitatively assess the effectiveness of our method in addressing  
 902 distributional shifts, we measure the distances between graph representations using cosine similarity  
 903 and Euclidean distance.

### 904 E.2 DATASETS

905 We use nine synthetic datasets and four real-world datasets in our experiments.  
 906

907 **BA-2Motifs (Luo et al., 2020).** The BA-2Motifs dataset comprises 1,000 synthetic graphs, each  
 908 generated from a Barabási–Albert graph and extended with either a house or a five-cycle motif. The  
 909 label of each graph is determined by its attached motif, forming a binary classification task in which  
 910 the motif serves as the ground-truth explanation.  
 911

912 **BA-HouseGrid (Bui et al., 2024).** The BA-HouseGrid dataset employs the house motif and the  
 913  $3 \times 3$  grid motif. These motifs are chosen to minimize structural overlap and ensure that models  
 914 learn the full motif structure instead of relying on local substructures for prediction.  
 915

916 **SPMotif (Wu et al., 2022).** In the SPMotif dataset, each graph consists of a base structure (Tree,  
 917 Ladder, or Wheel) and a motif (Cycle, House, or Crane). A parameter  $b$  controls the degree of the  
 spurious correlation between the base structure and the motif, with  $b = \frac{1}{3}$  indicating no spurious

918 correlation. In our experiments,  $b = 0.7$ . The label and ground-truth explanation for each graph are  
919 determined solely by the motif it contains.  
920

921 **BA-HouseAndGrid (Bui et al., 2024)**. The BA-HouseAndGrid dataset contains Barabási–Albert  
922 graph graphs extended with the house motif, the grid motif, or both. Graphs are labeled 1 if they  
923 contain both motifs, otherwise, they are labeled 0.

924 **Alkane-Carbonyl (Sanchez-Lengeling et al., 2020)**. The Alkane-Carbonyl dataset comprises  
925 4,326 molecular graphs partitioned into two classes based on the functional groups. A molecule is  
926 labeled positive when it contains both alkane and carbonyl groups, which also serve as the ground-  
927 truth explanation.  
928

929 **Fluorid-Carbonyl (Sanchez-Lengeling et al., 2020)**. The Fluorid-Carbonyl dataset comprises  
930 8,671 molecular graphs partitioned into two classes based on functional groups. A molecule is  
931 labeled positive if it contains both fluoride atoms and a carbonyl group, which also serve as the  
932 ground-truth explanation.  
933

934 **Benzene (Sanchez-Lengeling et al., 2020)**. The Benzene dataset comprises 12,000 molecular  
935 graphs from the ZINC15 (Sterling & Irwin, 2015) database, partitioned into two classes based on  
936 the presence of benzene rings. A molecule is labeled positive if it contains at least one benzene ring.  
937 Each benzene ring in the molecule serves as a distinct ground-truth explanation.

938 **BA-Motif-Volume (Zhang et al., 2024)**. In BA-Motif-Volume dataset, each graph is constructed  
939 from a Barabási–Albert graph with an attached five-cycle motif. Node features are assigned random  
940 float values in the range  $[0.00, 100.00]$ , and the regression label for each graph is defined as the sum  
941 of node feature values over the motif.  
942

943 **House-Grid-Volume**. The House-Grid-Volume dataset is derived from BA-HouseGrid (Bui et al.,  
944 2024) by replacing all node features with random floats sampled from  $[0.00, 100.00]$ . The regression  
945 label for each graph is defined as the sum of node features within its motif.

946 **BA-Motif-Counting (Zhang et al., 2024)**. The BA-Motif-Counting dataset consists of graphs cre-  
947 ated by attaching a randomly sampled number of five-cycle motifs (with the number varying from  
948  $\{0, \dots, 10\}$  in our experiment) to a Barabási–Albert random graph. The number of motifs in each  
949 graph serving as its regression label.  
950

951 **Triangles (Chen et al., 2020b)**. The Triangles dataset is constructed following the prior work (Chen  
952 et al., 2020b), consisting of 5,000 Erdős–Rényi random graphs denoted as  $ER(m, p)$ , where  $m = 30$   
953 is the number of nodes in each graph and  $p = 0.1$  is the edge existence probability. The regression  
954 label for each graph is the number of triangles it contains. In our experiments, we use a subset of  
955 1,000 graphs randomly sampled from this dataset.

956 **Crippen (Delaney, 2004)**. The Crippen dataset contains 1,127 molecules with corresponding aque-  
957 ous solubility measurements from the Delaney solubility (Delaney, 2004) dataset and assigns node  
958 weights using the Crippen model (Wildman & Crippen, 1999). Following prior work (Sanchez-  
959 Lengeling et al., 2020), we adopt this dataset and generate edge weights as the average of incident  
960 node weights.  
961

962 **House-OrGrid-Volume**. The House-OrGrid-Volume dataset is derived from BA-HouseOrGrid (Bui  
963 et al., 2024) by replacing all node features with random floats sampled from  $[0.00, 100.00]$ . The  
964 regression label for each graph is defined as the sum of node features within its motif.  
965

### 967 E.3 BASELINES

968  
969  
970 To assess effectiveness, we incorporate various post-hoc methods, including GNNExplainer, PGEx-  
971 plainer, MetaGNN, MatchExplainer, MixupExplainer, ProxyExplainer, and RegExplainer, as well  
as the gradient-based GRAD and the attention-based ATT.

972 **GRAD (Ying et al., 2019)**. GRAD is a gradient-based method that generates weights to edges and  
 973 nodes by computing the gradients of the loss function of GNN with respect to the adjacency matrix  
 974 and node features.

975  
 976 **ATT (Veličković et al., 2017)**. ATT is a graph attention network that learns attention weights for  
 977 edges in the input graph, with these weights are adopted as a proxy measure of edge importance.

978  
 979 **GNNExplainer (Ying et al., 2019)**. GNNExplainer learns soft masks over edges and node features  
 980 for each instance by minimizing the mutual information between the original graph and the predic-  
 981 tion results. The explanatory subgraphs are obtained via element-wise multiplication of the learned  
 982 soft masks with the original graph.

983 **PGExplainer (Luo et al., 2020)**. PGExplainer extends GNNExplainer by parameterizing the ex-  
 984 planation generation process with a trainable explainer, and generates a soft mask to produce the  
 985 explanatory subgraph.

986  
 987 **MetaGNN (Spinelli et al., 2022)**. MetaGNN trains GNNs to be inherently interpretable by incor-  
 988 porating a meta-explainer, which generates post-hoc explanations during training.

989 **MatchExplainer (Wu et al., 2023)**. MatchExplainer explains GNN predictions by matching shared  
 990 subgraph patterns using graph edit distance (GED) as the similarity metric, and this non-parametric  
 991 subgraph matching approach inherently avoids optimization bias.

992  
 993 **MixupExplainer (Zhang et al., 2023)**. MixupExplainer mitigates the distribution shift present in  
 994 previous methods by mixing explanatory subgraphs with the label-irrelevant parts of other randomly  
 995 sampled graphs in a non-parametric manner.

996  
 997 **ProxyExplainer (Chen et al., 2024)**. ProxyExplainer performs autoencoders to reconstruct the  
 998 explanatory and label-irrelevant parts, then combines them to generate proxy graphs in a parametric  
 999 manner, approximating the original distribution to mitigate the OOD issue.

1000 **RegExplainer (Zhang et al., 2024)**. RegExplainer addresses interpretability in graph regression  
 1001 tasks by combining mixup with contrastive learning to mitigate distribution shift and tackle the  
 1002 challenge of continuously ordered labels.

## 1003 F PROOF OF PROPERTY 1

1004  
 1005 *Proof.* The mutual information between the explanatory subgraph  $G^*$  and the label  $Y$  can be for-  
 1006 mulated as the Kullback–Leibler (KL) divergence between the joint distribution  $p(G^*, Y)$  and the  
 1007 product of its marginals:

$$1008 I(G^*; Y) = \text{KL}(p(G^*, Y) \parallel p(G^*)p(Y)).$$

1009  
 1010 The Donsker–Varadhan (DV) variational representation of KL divergence Poole et al. (2019) pro-  
 1011 vides the following inequality:

$$1012 \text{KL}(p \parallel q) = \sup_T \{ \mathbb{E}_p[T] - \log \mathbb{E}_q[e^T] \},$$

1013  
 1014 where the supremum is taken over all measurable functions  $T : \mathcal{G} \times \mathcal{Y} \rightarrow \mathbb{R}$  such that the expecta-  
 1015 tions exist. The bound is tight when  $T(x) = \log \frac{p(x)}{q(x)}$ .

1016  
 1017 Applying this representation to the mutual information yields:

$$1018 I(G^*; Y) \geq \mathbb{E}_{p(G^*, Y)}[T(G^*, Y)] - \log \mathbb{E}_{p(G^*)p(Y)} \left[ e^{T(G^*, Y)} \right].$$

1019  
 1020 This variational characterization establishes a theoretical basis for estimating mutual information by  
 1021 choosing tractable parameterizations of the critic function  $T$ .

1022  
 1023  
 1024  
 1025 □

## 1026 G PROOF OF PROPERTY 2

1027  
1028  
1029 *Proof.* To instantiate the critic function  $T(G^*, Y)$ , we leverage the fact that each label  $Y$  in the  
1030 dataset is paired with a real input graph  $G_Y$ , i.e.,  $(G_Y, Y) \sim p(G, Y)$ . Instead of encoding the label  
1031  $Y$  directly, we obtain its representation via the associated graph  $G_Y$ , which reflects the semantics  
1032 of  $Y$  in the same space as  $G^*$ . This design ensures that both  $G^*$  and  $Y$  are represented in the  
1033 same embedding space, while preserving the theoretical foundation of estimating  $I(G^*; Y)$ . Let  
1034  $f(\cdot)$  denote an encoder that maps a graph to its representation, and define  $\mathbf{h}_{G^*} = f(G^*)$  and  $\mathbf{h}_Y =$   
1035  $f(G_Y)$ . The similarity between two representations is defined as  $\text{sim}(u, v) := u^\top v$ , where  $u$  and  $v$   
1036 are  $\ell_2$ -normalized vectors.

1037 We instantiate the Donsker–Varadhan critic function  $T : \mathcal{G} \times \mathcal{Y} \rightarrow \mathbb{R}$  as:

$$1038 \quad T(G^*, Y) := \text{sim}(\mathbf{h}_{G^*}, \mathbf{h}_Y) + \log(N-1),$$

1039 where  $N$  is the total number of samples (including one positive and  $N-1$  negatives) in each estima-  
1040 tion batch. We first evaluate the expectation under the joint distribution  $p(G^*, Y)$  (positive pairs):

$$1041 \quad \mathbb{E}_{p(G^*, Y)}[T(G^*, Y)] = \mathbb{E}_{p(G^*, Y)}[\text{sim}^+ + \log(N-1)],$$

1042 where  $\text{sim}^+ := \text{sim}(\mathbf{h}_{g^*}, \mathbf{h}_{y^+})$  denotes the similarity score of the positive pair.

1043 To evaluate the negative term, we apply Jensen’s inequality to the expectation under the product of  
1044 marginals:

$$1045 \quad -\log \mathbb{E}_{p(G^*)p(Y)}[e^{T(G^*, Y)}] \geq -\mathbb{E}_{G^* \sim p(G^*)}[\log \mathbb{E}_{Y \sim p(Y)}[e^{T(G^*, Y)}]].$$

1046 We approximate the inner expectation using  $N-1$  i.i.d. samples  $\{y_j^-\}_{j=1}^{N-1}$  drawn from  $p(Y)$ :

$$1047 \quad \mathbb{E}_Y[e^{T(G^*, Y)}] \approx \frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^- + \log(N-1)) = (N-1) \cdot \frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-),$$

1048 where  $\text{sim}_j^- := \text{sim}(\mathbf{h}_{g^*}, \mathbf{h}_{y_j^-})$  denotes the similarity between  $G^*$  and each negative sample  $y_j^-$ .

1049 Taking the logarithm, we obtain:

$$1050 \quad \log \mathbb{E}_Y[e^{T(G^*, Y)}] \approx \log(N-1) + \log\left(\frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-)\right),$$

1051 and thus, the negative term becomes:

$$1052 \quad -\log \mathbb{E}_{p(G^*)p(Y)}[e^{T(G^*, Y)}] \geq -\log(N-1) - \mathbb{E}_{\mathcal{P}}\left[\log\left(\frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-)\right)\right],$$

1053 where  $\mathcal{P} := p(G^*, Y_0) \prod_{j=1}^{N-1} p(Y_j^-)$  denotes the joint sampling distribution comprising one posi-  
1054 tive pair  $(G^*, Y_0)$  and  $N-1$  independent negative samples  $\{Y_j^-\}$ .

Combining the positive and negative components of the DV bound:

$$\begin{aligned}
I(G^*; Y) &\geq \mathbb{E}_{\mathcal{P}} [\text{sim}^+ + \log(N-1)] - \log(N-1) - \mathbb{E}_{\mathcal{P}} \left[ \log \left( \frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-) \right) \right] \\
&= \mathbb{E}_{\mathcal{P}} [\text{sim}^+] - \mathbb{E}_{\mathcal{P}} \left[ \log \left( \frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-) \right) \right] \\
&= \mathbb{E}_{\mathcal{P}} \left[ \log(\exp(\text{sim}^+)) - \log \left( \frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-) \right) \right] \\
&= \mathbb{E}_{\mathcal{P}} \left[ \log \left( \frac{\exp(\text{sim}^+)}{\frac{1}{N-1} \sum_{j=1}^{N-1} \exp(\text{sim}_j^-)} \right) \right] \\
&= \log(N-1) + \mathbb{E}_{\mathcal{P}} \left[ \log \left( \frac{\exp(\text{sim}^+)}{\sum_{j=1}^{N-1} \exp(\text{sim}_j^-)} \right) \right].
\end{aligned}$$

Since the logarithm is monotonically increasing, augmenting the denominator with the positive term yields a smaller ratio and hence a looser (but valid) bound:

$$I(G^*; Y) \geq \log(N-1) + \mathbb{E}_{\mathcal{P}} \left[ \log \left( \frac{\exp(\text{sim}^+)}{\exp(\text{sim}^+) + \sum_{j=1}^{N-1} \exp(\text{sim}_j^-)} \right) \right].$$

We define the InfoNCE estimator as:

$$\ell(G^*, Y) := \log \frac{\exp(\text{sim}(\mathbf{h}_{g^*}, \mathbf{h}_{y^+}))}{\sum_{j=0}^{N-1} \exp(\text{sim}(\mathbf{h}_{g^*}, \mathbf{h}_{y_j}))},$$

where  $\{y_j\}_{j=0}^{N-1}$  includes one positive label  $y_0 = y^+$  and  $N-1$  negative samples  $y_j^-$ .

Thus, we arrive at the final variational lower bound:

$$I(G^*; Y) \geq \log(N-1) + \mathbb{E}_{\mathcal{P}} [\ell(G^*, Y)].$$

□

## H EXTENSIVE EXPERIMENTS

### H.1 ABLATION STUDY

We conduct ablation studies to evaluate the contribution of each component in the TAME framework. Specifically, we design the following variants: 1) *w/o Mix*: Removes the structural mixup step after extracting the explanation, directly feeding the explanation subgraph into the objective function. 2) *w/o CTL*: Retains the structural mixup strategy and BCE loss, but replaces the contrastive learning term with cross-entropy loss for classification and mean squared error (MSE) loss for regression. 3) *w/o BCE*: Removes only the binary mask behavior regularization term from the training loss. All variants are configured with the same hyperparameters as the original TAME, including the learning rate, number of training epochs, and the loss weights  $\beta$ .

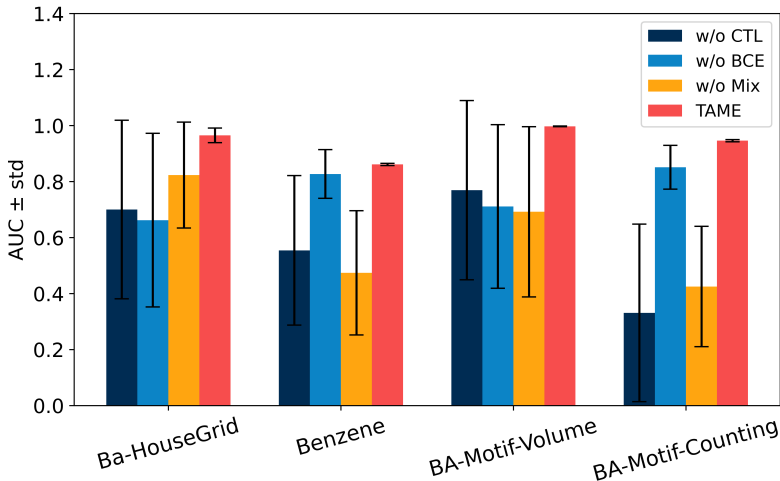


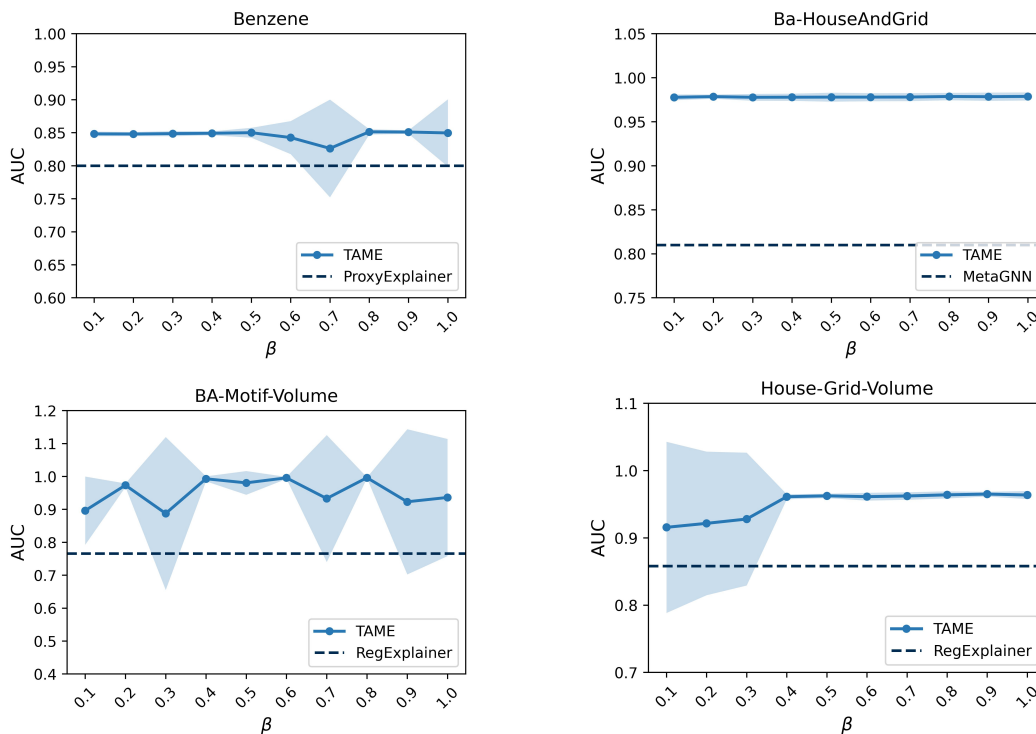
Figure 6: Ablation study of TAME across classification and regression tasks. We evaluate the AUC-ROC performance of the original TAME and its variants. Standard deviations are indicated by black error bars.

We conduct experiments on representative datasets covering both classification and regression tasks, with BA-HouseGrid and Benzene used for classification, and BA-motif-volume and BA-Motif-Counting used for regression. The results are presented in Figure 6. Experimental results show that TAME consistently outperforms all its variants across all datasets, indicating that each component contributes positively to the overall performance. Notably, the performance of the *w/o CTL* variant drops significantly on both classification and regression datasets, indicating that our proposed contrastive learning mechanism plays a critical role in both tasks.

### H.2 HYPER-PARAMETER SENSITIVITY STUDY

We further investigate the sensitivity of our proposed method to the hyperparameter  $\beta$ , which controls the weight of the BCE loss. Specifically, we conduct experiments on four representative datasets, including the Ba-HouseAndGrid and Benzene classification datasets, as well as the BA-Motif-Volume and House-Grid-Volume regression datasets. We vary  $\beta$  from 0.1 to 1.0 in increments of 0.1 while keeping the contrastive learning weight fixed at 1.0. The experiments are conducted under ten random seeds, where the solid lines represent the average AUC of our method and the dashed lines correspond to the second-best performing method. The experimental results are presented in Figure 7. The results indicate that TAME maintains stable performance across both classification and regression tasks. For some datasets, the AUC slightly decreases and exhibits larger variance when the weight of the BCE loss is small. This phenomenon can be attributed to the role of the BCE loss in constraining the edge weights to remain faithful to the explanation structure extracted by graph pooling, while simultaneously reducing the weights of non-explanatory structures. Without

1188 this constraint, the edge weights may not be effectively optimized, leading to a decline in perfor-  
 1189 mance.  
 1190



1215 Figure 7: Sensitivity analysis of hyperparameter  $\beta$  for BCE loss.  
 1216  
 1217  
 1218  
 1219  
 1220  
 1221  
 1222  
 1223  
 1224  
 1225  
 1226  
 1227  
 1228  
 1229  
 1230  
 1231  
 1232  
 1233  
 1234  
 1235  
 1236  
 1237  
 1238  
 1239  
 1240  
 1241

### 1242 H.3 BACKBONE MODELS ANALYSIS

1243  
1244 We conduct a performance analysis of common backbone models employed by explainability meth-  
1245 ods on both classification and regression tasks. For the classification task, model performance  
1246 is evaluated using ACC and Macro\_F1 metrics on four synthetic datasets, including BA-2Motifs,  
1247 BA-HouseGrid, BA-HouseAndGrid, and SPMotif, as well as three real-world datasets including  
1248 Alkane-Carbonyl, Fluoride-Carbonyl, and Benzene. For the regression task, RMSE and MAPE are  
1249 adopted as evaluation metrics on five synthetic datasets including BA-Motif-Volume, BA-Motif-  
1250 Counting, Triangles, House-Grid-Volume, and House-OrGrid-Volume, along with one real-world  
1251 dataset, Crippen.

1252 Experimental results demonstrate that GCN and GIN achieve strong performance in classification  
1253 tasks, while GCN consistently attains the best performance across regression benchmarks. There-  
1254 fore, to ensure that the explanation quality is not confounded by the performance of the backbone  
1255 model, we choose GCN as the backbone model for explainer evaluation.

Dataset	GCN		GIN		GAT	
	ACC	Macro_F1	ACC	Macro_F1	ACC	Macro_F1
BA-2motifs	99.00	98.98	<b>100.00</b>	<b>100.00</b>	44.00	30.55
BA-HouseGrid	99.90	99.89	<b>100.00</b>	<b>100.00</b>	49.10	32.93
BA-HouseAndGrid	98.40	98.39	<b>99.90</b>	<b>99.89</b>	48.90	32.84
SPMotif	<b>96.16</b>	<b>96.15</b>	95.88	95.89	34.05	18.98
Alkane-Carbonyl	<b>95.57</b>	<b>95.13</b>	92.03	91.24	94.69	94.31
Fluoride-Carbonyl	94.35	89.99	<b>97.23</b>	<b>95.03</b>	83.29	45.44
Benzene	90.58	90.53	<b>92.66</b>	<b>92.65</b>	79.08	79.06

1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267 Table 5: Comparison of GCN, GIN, and GAT on classification datasets using ACC and Macro-F1  
1268 metrics.  
1269

Dataset	GCN		GIN		GAT	
	RMSE	MAPE	RMSE	MAPE	RMSE	MAPE
BA-Motif-Volume	<b>234.20</b>	<b>7.47</b>	398.96	14.94	485.31	17.45
BA-Motif-Counting	<b>0.39</b>	—	1.94	—	2.84	—
Triangles	<b>2.00</b>	—	2.46	—	2.55	—
House-Grid-Volume	<b>11.32</b>	<b>2.45</b>	94.54	29.56	49.85	13.88
House-OrGrid-Volume	<b>34.15</b>	<b>6.69</b>	152.14	28.17	53.64	10.27
Crippen	<b>0.92</b>	<b>30.71</b>	1.91	123.17	1.16	74.68

1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281 Table 6: Comparison of GCN, GIN, and GAT on regression datasets using RMSE and MAPE met-  
1282 rics. MAPE values are marked with dashes for BA-Motif-Counting and Triangles due to the presence  
1283 of label-zero samples that make MAPE undefined.  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

## H.4 EXTENSIVE CASE STUDY

### H.4.1 VISUALIZATION OF ADDITIONAL EXPLANATION RESULTS

We further present extensive visualization experiments on both classification and regression datasets. The classification datasets include BA-2Motifs, BA-HouseAndGrid, and Benzene, as shown in Figure 8, Figure 9, and Figure 10, respectively. The regression datasets include BA-Motif-Volume, and House-OrGrid-Volume, as presented in Figure 11, and Figure 12, respectively. For each dataset, we randomly select three target graphs to evaluate different explainers.

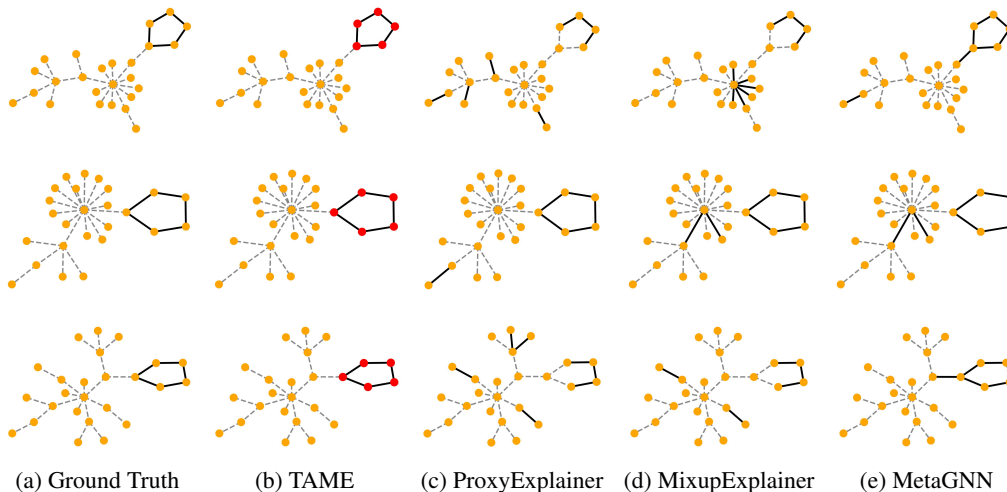


Figure 8: Visualization of explanation on BA-2motifs.

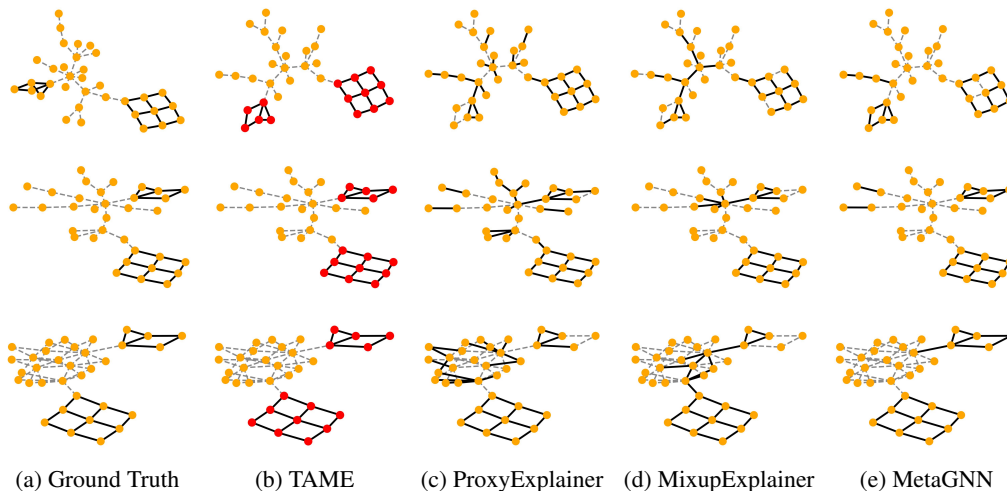


Figure 9: Visualization of explanation on BA-HouseAndGrid.

1350

1351

1352

1353

1354

1355

1356

1357

1358

1359

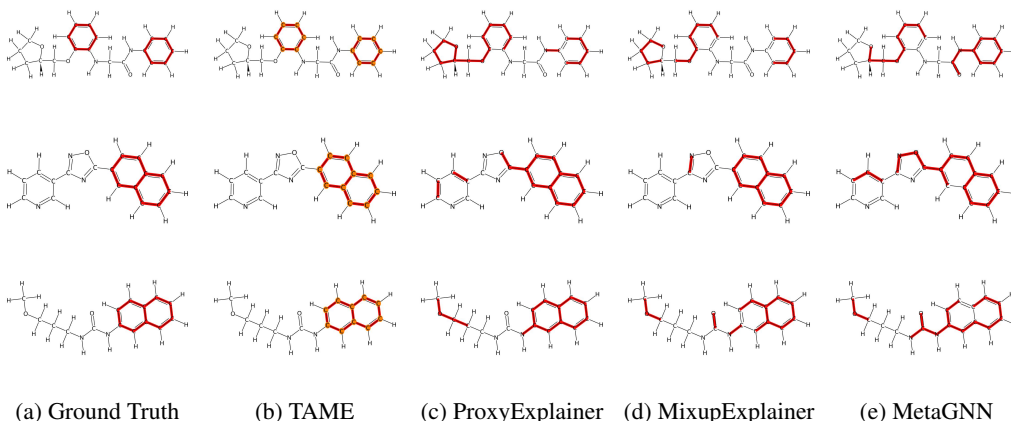
1360

1361

1362

1363

1364



1365

Figure 10: Visualization of explanation on Benzene.

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

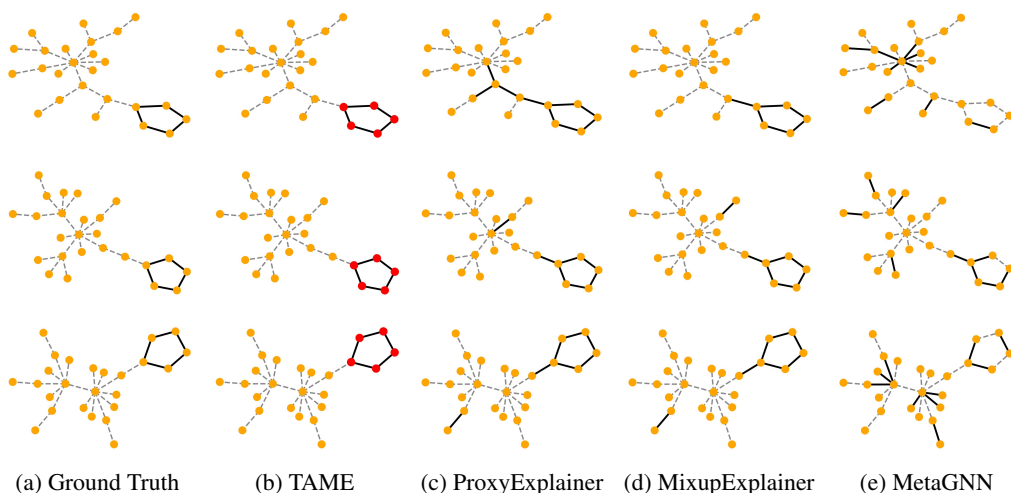
1379

1380

1381

1382

1383



1384

Figure 11: Visualization of explanation on BA-Motif-Volume.

1385

1386

1387

1388

1389

1390

1391

1392

1393

1394

1395

1396

1397

1398

1399

1400

1401

1402

1403

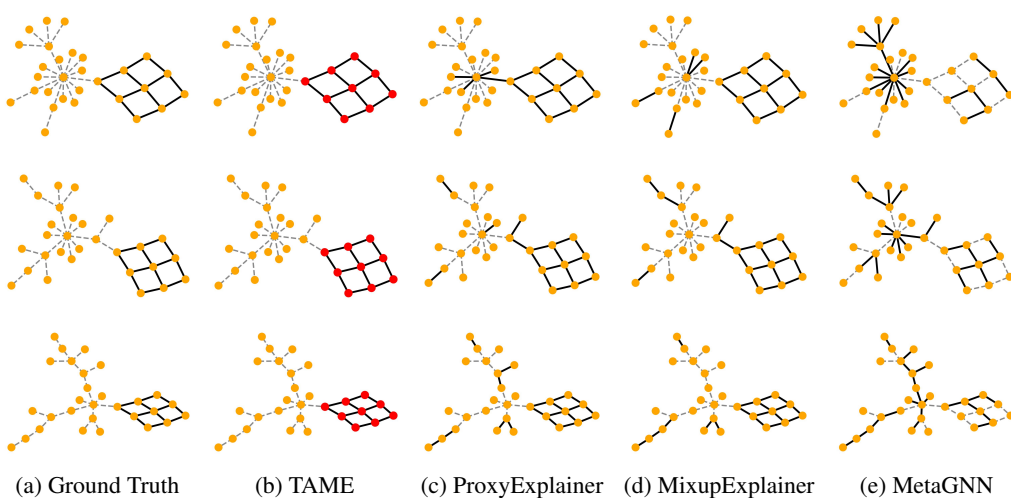
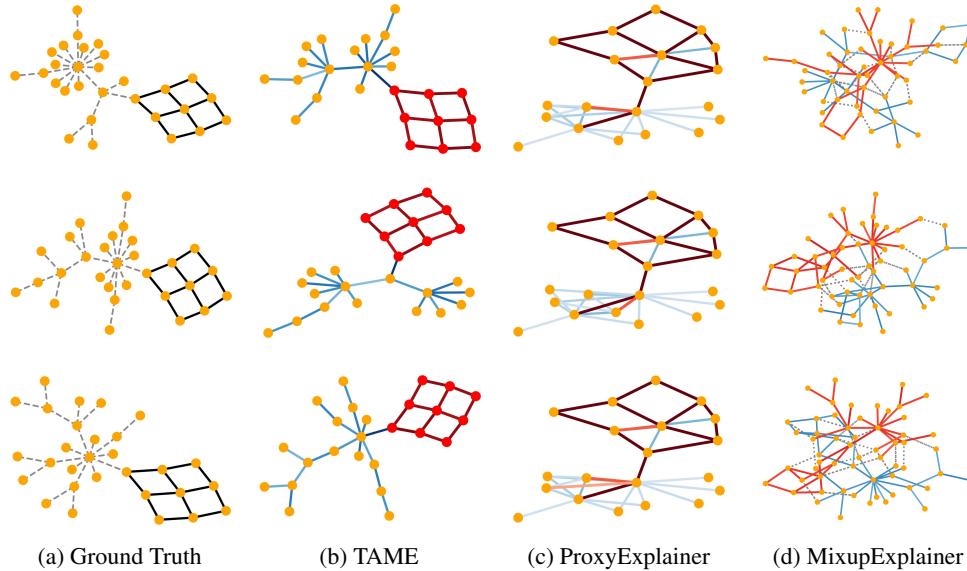


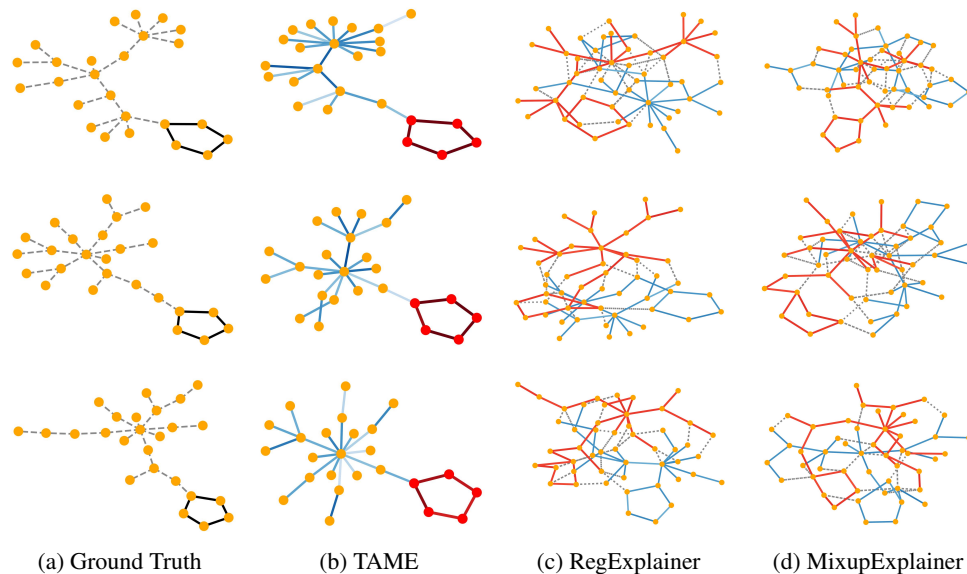
Figure 12: Visualization of explanation on House-Grid-Volume.

#### 1404 H.4.2 VISUALIZATION OF ADDITIONAL MIXUP RESULTS

1405 We present the mixup results on the BA-HouseGrid classification dataset and the BA-Motif-Volume  
 1406 regression dataset in Figure 13 and Figure 14, respectively. For each dataset, three graph samples  
 1407 are randomly selected to generate the mixup results, where the edge color intensity indicates the cor-  
 1408 responding weight magnitude. It should be noted that ProxyExplainer employs a graph generation  
 1409 strategy to produce mixup results, which leads to fully connected edges. To facilitate visualization,  
 1410 we display only the top 64 edges with the highest weights, corresponding to the maximum number of  
 1411 edges considered in this experiment.  
 1412



1432 Figure 13: Visualization of mixup graphs on BA-HouseGrid.



1454 Figure 14: Visualization of mixup graphs on BA-Motif-Volume.

1455

1456

1457