

# MEGAN: Multi-Explanation Graph Attention Network

Jonas Teufel<sup>1</sup>[0000-0002-9228-9395], Luca Torresi<sup>1</sup>[0000-0003-2205-6753], Patrick Reiser<sup>1</sup>[0000-0002-7052-696X], and Pascal Friederich<sup>\*1</sup>[0000-0003-4465-1465]

Institute of Theoretical Informatics (ITI), Karlsruhe Institute of Technology (KIT),  
Karlsruhe, Germany

{jonas.teufel}@student.kit.edu  
{luca.torresi,patrick.reiser,pascal.friederich}@kit.edu

**Abstract.** We propose a multi-explanation graph attention network (MEGAN). Unlike existing graph explainability methods, our network can produce node and edge attributional explanations along multiple channels, the number of which is independent of task specifications. This proves crucial to improve the interpretability of graph regression predictions, as explanations can be split into positive and negative evidence w.r.t to a reference value. Additionally, our attention-based network is fully differentiable and explanations can actively be trained in an explanation-supervised manner. We first validate our model on a synthetic graph regression dataset with known ground-truth explanations. Our network outperforms existing baseline explainability methods for the single- as well as the multi-explanation case, achieving near-perfect explanation accuracy during explanation supervision. Finally, we demonstrate our model’s capabilities on multiple real-world datasets. We find that our model produces sparse high-fidelity explanations consistent with human intuition about those tasks.

**Keywords:** Graph Neural Network · Self-Explaining Model · Explanation Supervision

## 1 Introduction

Explainable AI (XAI) methods aim to provide explanations complementing a model’s predictions to make it’s complex inner workings more transparent to humans with the intention to improve trust and reliability, provide tools for model analysis, and comply with anti-discrimination laws [8]. The majority of existing work on graph explainability focuses on post-hoc methods, which can be used to generate explanations for already trained models, which have been proven to perform well. While post-hoc methods are an important area of development to add explainability to time-tested models, we want to emphasize the potential of *self-explaining* methods. In their literature review, Jiminez-Luna *et al.* [17] describe these methods as being explainable by design. One example

---

\* corresponding author

of this class are the simpler, traditional machine learning approaches that are naturally interpretable, such as decision tree methods [11]. However, we want to focus on self-explaining graph neural networks, which produce the attributional explanations for the nodes and edges of the input graph directly alongside each prediction. We emphasize this class of methods specifically due to their capability for explanation-supervised training. During explanation-supervised training, a model is additionally trained to produce explanations that are similar to a given set of reference explanations. Recently, there has been promising progress on the topic of explanation supervision in the domains of image processing [19,29,2] and natural language processing [10,28,37]. Previous work is able to improve model interpretability by training models to generate more human-like explanations and even improve main prediction performance by training models on human-generated image saliency maps. In the graph domain, however, there has been little work on explanation supervision [13,21] yet. Inspired by the successes recently demonstrated in other domains, we propose the self-explaining *multi-explanation graph attention network* (MEGAN) architecture. In this work, we demonstrate that our model shows significantly improved capability to learn explanations during explanation-supervised training, outperforming the baseline method [13] from the literature.

In addition to its properties w.r.t. explanation supervision, we design our network to output explanations along *multiple channels*, the number of which is independent of the main prediction task. Like the majority of existing GNN explainability methods, we focus on attributional explanations, which attribute a value of importance to each element of the input graph. For existing methods, the number of these attribution values is dictated by the details of the main prediction task. For single-value graph regression tasks for example a single value would be assigned to each node and edge. For our multi-explanation method, however, this number of attributions is a property of the network rather than restricted by task specifications.

We want to emphasize the importance of this property especially in regard to graph regression problems. For the prediction of a single regression value, existing methods only produce a single attribution for each node and edge. We argue that such explanations are insufficient for the interpretation of regression predictions. In reality, one often encounters structure-property explanations of opposing *polarity*. One practical example of this is the prediction of water solubility, where large non-polar carbon structures generally cause low solubility values and polar functional groups cause higher values. A single attributional explanation may highlight all the important motifs, but is not able to capture this crucial detail about their polarity. For this reason, we decouple the number of explanations from the task specification to be able to produce two explanations (negative and positive influence) for graph regression problems. We introduce an explanation co-training method which uses only the generated explanation masks to solve an approximation of the prediction problem to promote each explanation channel to behave according to their intended interpretation. In our

experiments, we find that this explanation co-training is an effective method to guide the generation of the explanation channels to contribute faithfully to the prediction outcome according to pre-determined interpretations. We validate this finding on several real-world datasets, where our model produces explanations consistent with human intuition about those tasks. Beyond that, we apply our model to one real-world task of molecular property prediction without common human intuition and are able to support previously published hypotheses about structure-property relationships and propose several new potential explanatory motifs.

## 2 Related Work

**GNN Explanation Methods** Yuan *et al.* [41] provide a taxonomic overview of XAI methods for graph neural networks. Some methods have been adapted from similar approaches in other domains, such as GradCAM [26], GraphLIME [16] and LRP [33]. Other methods were developed specifically for graph neural networks. Notable ones include GNNExplainer [40], PGExplainer [20], and Zorro [12]. Jiminez-Luna *et al.* [17] present another literature review about the applications of XAI in drug discovery. Henderson *et al.* [15] for example introduce regularization terms to improve GradCAM-generated explanations for chemical property prediction. Sanchez-Lenglin *et al.* [32] introduce new benchmark datasets for attributional graph explanations based on molecular graphs and compare several existing explanation methods.

Generally, most explanation methods aim to produce attributional explanations, which explain a prediction by assigning importance values to the nodes and edges of the input graph. However, there exists some criticism about this class of explanations [1,18], which is partially why recently different modalities of explanations have been explored for the graph domain as well. Magister *et al.* [22] for example propose GCExplainer, which can be used to generate *concept-based* explanations for graph neural networks in a post-hoc fashion. Shin *et al.* [34] for example propose PAGE, a method to generate *prototype-based* explanations. *Counterfactuals* are yet another popular explanation modality, for which Tan *et al.* [38] and Prado-Romero and Stilo [27] have recently proposed methods for graph neural networks.

**Self-Explaining Graph Neural Networks** In their literature review, Jiminez-Luna *et al.* [17] define *self-explaining* methods as those that are explainable by design. One large fraction of this category is represented by simpler traditional machine learning methods. Friederich *et al.* [11] for example use an interpretable decision tree approach to structure-property relationships for several real-world graph datasets. However, there is also recent progress for more complex self-explaining models such as graph neural networks. Dang and Wang [4] and Zhang *et al.* [43] independently introduce self-explaining graph neural networks for prototype-based explanations. Magister *et al.* [21] introduce a self-explaining network for concept-based explanations. Furthermore, Müller *et al.*

[24] propose DT+GNN, an interesting method that combines the capabilities of GNNs with the inherent interpretability of decision trees.

**Explanation Supervision** During explanation supervision, models are not only trained to perform a main prediction task through ground truth target labels but also to produce explanations that are similar to a given set of reference explanations. Most interestingly explanation supervision provides the possibility to train models to produce more human-like explanations. Beyond that, several works are able to show that the inclusion of human saliency maps has the potential to increase the task performance of the models [19,2]. In that context, Linseley *et al.* [19] for example show that human saliency maps improve the performance of an image classifier. Boyd *et al.* [2] demonstrate that human saliency annotations improve the performance of a deep fake detection model. In the domain of natural language processing, Pruthi *et al.* [28] use explanation-supervised models to substitute human participants in artificial simulatability studies to assess the quality of explanations. Fernandes *et al.* [10] even take this concept one step further and train an explainer to optimize this property of simulatability.

### 3 Multi-Explanation Graph Attention Network

#### 3.1 Task Description

We assume a directed graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  is represented by a set of node indices  $\mathcal{V} \subset \mathbb{N}^V$  and a set of edges  $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V} \subset \mathbb{R}^E$ , where a tuple  $(i, j) \in \mathcal{E}$  denotes an edge from node  $i$  to node  $j$ . Every node  $i$  is associated with a vector of initial node features  $\mathbf{h}_i^{(0)} \in \mathbb{R}^{N_0}$ , combining into the initial node feature tensor  $\mathbf{H}^{(0)} \in \mathbb{R}^{V \times N_0}$ . Each edge is associated with a feature vector  $\mathbf{u}_i \in \mathbb{R}^M$ , combining into the edge feature tensor  $\mathbf{U} \in \mathbb{R}^{E \times M}$ .

We consider graph classification and regression problems, which means graphs are associated with a target vector  $\mathbf{y} \in \mathbb{R}^C$  which is either a one-hot class encoding or continuous regression values. In addition, node and edge attributional explanations for graphs are considered. We define explanations as masks that assign  $[0, 1]$  values to each node and each edge, representing the importance of the corresponding graph element toward the outcome of the prediction. We generally assume that any prediction may be explained by  $K$  individual importance channels, where  $K$  is an independent hyperparameter. The node explanations are given as the *node importance* tensor  $\mathbf{V}^{\text{im}} \in [0, 1]^{V \times K}$  and the edge explanations are given as the *edge importance* tensor  $\mathbf{E}^{\text{im}} \in [0, 1]^{E \times K}$ .

#### 3.2 Architecture Overview

To solve the previously defined task we propose the following *multi-explanation graph attention network* (MEGAN) architecture, for which Figure 1 provides a visual overview. The network consists of  $L$  attention layers, where the number of layers  $L$  and the hidden units of each layer are hyperparameters. Each of these

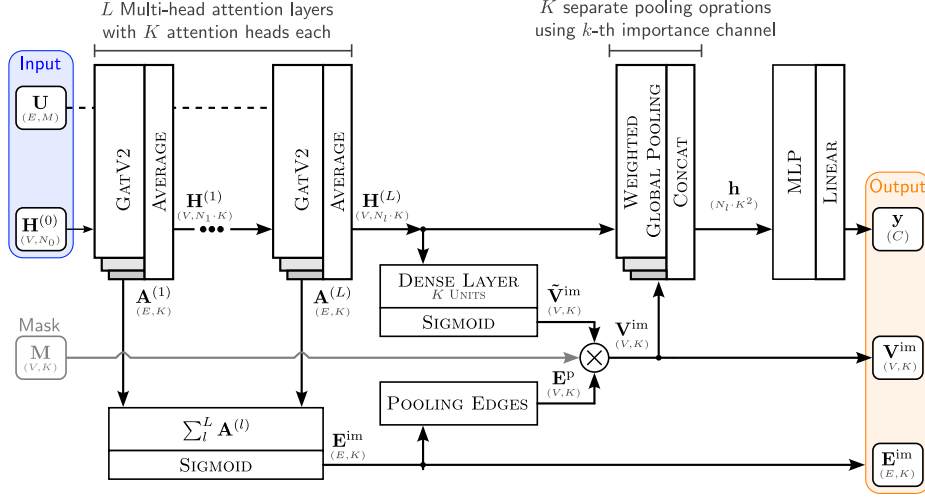


Fig. 1: Multi-explanation graph attention network (MEGAN) architecture overview. Rectangle boxes represent layers; arrows indicate layer interconnections. Rounded boxes represent tensors. Intermediate tensors are also named annotated arrows. Tuples beneath variable names indicate the tensor shape, with batch dimension omitted, but implicitly assumed as the first dimension for all.

layers consists of  $K$  individual, yet structurally identical GATv2 [3] attention heads, one for each of the  $K$  expected explanation channels. Assuming the attention heads in the  $l$ -th layer have  $N_l$  hidden units, then each attention head produces its own node embeddings  $\mathbf{H}^{(l,k)}$ , where  $k \in \{1, \dots, K\}$  is the head index. The final node embeddings  $\mathbf{H}^{(l)} \in \mathbb{R}^{V \times N_l \cdot K}$  of layer  $l$  are then produced by averaging all these individual matrices along the feature dimension:

$$\mathbf{H}^{(l)} = \frac{1}{K} \sum_k \mathbf{H}^{(l,k)} \quad (1)$$

This node embedding tensor is then used as the input to *each* of the  $K$  attention heads of layer  $l + 1$ . Aside from the node embeddings, each attention head also produces a vector  $\mathbf{A}^{(l,k)} \in \mathbb{R}^E$  of attention logits which are used to calculate the attention weights

$$\boldsymbol{\alpha}^{(l,k)} = \text{softmax}(\mathbf{A}^{(l,k)}) \quad (2)$$

of the  $k$ -th attention head in the  $l$ -th layer. The edge importance tensor  $\mathbf{E}^{\text{im}} \in [0, 1]^{E \times K}$  is calculated from the concatenation of these attention logit tensors in the feature dimension and summed up over the number of layers:

$$\mathbf{E}^{\text{im}} = \sigma \left( \sum_{l=1}^L \left( \mathbf{A}^{(l,1)} \parallel \mathbf{A}^{(l,2)} \parallel \dots \parallel \mathbf{A}^{(l,K)} \right) \right) \quad (3)$$

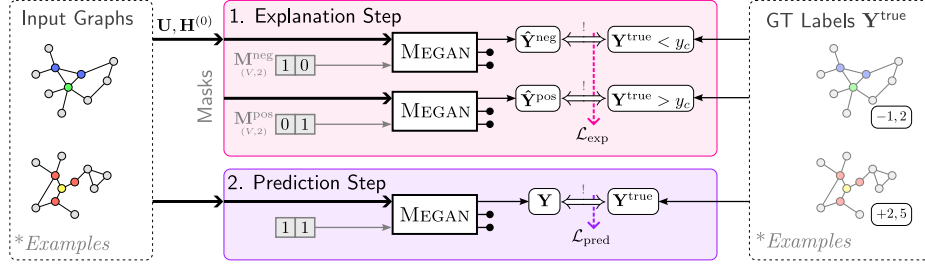


Fig. 2: Illustration of the split training procedure for the regression case. The explanation-only train step attempts to find an approximate solution to the main prediction task, by using only a globally pooled node importance tensor. After the weight update for the explanation step was applied to the model, the prediction step performs another weight update based on the actual output of the model and the ground truth labels.

Based on this, a local pooling operation is used to derive the pooled edge importance tensor  $\mathbf{E}^p \in [0, 1]^{V \times K}$  for the *nodes* of the graph. This local pooling operation can be seen as the aggregation step in a message-passing framework, where the edge importance values are treated as the corresponding messages. The final node embeddings  $\mathbf{H}^{(L)}$  are then used as the input to a dense network, whose final layer is set to have  $K$  hidden units, producing the node importance embeddings  $\tilde{\mathbf{V}}^{\text{im}} \in [0, 1]^{V \times K}$ . The node importance tensor is then calculated as the product of those node importance embeddings  $\tilde{\mathbf{V}}^{\text{im}} \in [0, 1]^{V \times K}$  and the pooled edge importance tensor  $\mathbf{E}^p \in [0, 1]^{V \times K}$ :

$$\mathbf{V}^{\text{im}} = \tilde{\mathbf{V}}^{\text{im}} \cdot \mathbf{E}^p \cdot \mathbf{M}. \quad (4)$$

The mask  $\mathbf{M}$  introduced in Fig. 1 is only optionally used to compute the fidelity metric, which is introduced in Section 3.4. At this point, the edge and node importance matrices, which represent the explanations generated by the network, are already accounted for, which leaves only the primary prediction to be explained. The first remaining step is a global sum pooling operation which turns the node embedding tensor  $\mathbf{H}^{(L)}$  into a vector of global graph embeddings. For this,  $K$  separate weighted global sum pooling operations are performed, one for each explanation channel. Each of these pooling operations uses the same node embeddings  $\mathbf{H}^{(L)}$  as input, but a different slice  $\mathbf{V}_{:,k}^{\text{im}}$  of the node importance tensor as weights. In that way,  $K$  separate graph embedding vectors

$$\mathbf{h}^{(k)} = \sum_{i=0}^V \left( \mathbf{H}^{(L)} \cdot \mathbf{V}_{:,k}^{\text{im}} \right)_{i,:} \quad (5)$$

are created, which are then concatenated into a single graph embedding vector

$$\mathbf{h} = \mathbf{h}^{(1)} \parallel \mathbf{h}^{(2)} \parallel \dots \parallel \mathbf{h}^{(K)} \quad (6)$$

where  $\mathbf{h} \in \mathbb{R}^{N_L \cdot K}$ . This graph embedding vector is then passed through a generic MLP whose final layer either has linear activation for graph regression or softmax activation for graph classification to create an appropriate output

$$\mathbf{y} = \text{MLP}(\mathbf{h}) \quad (7)$$

### 3.3 Explanation Co-Training

With the architecture as explained up to this point, there is no mechanism yet to ensure that individual explanation channels learn the appropriate explanations according to their intended interpretation (for example positive vs negative evidence). We use a special explanation co-training procedure to guide the individual explanation channels to develop according to pre-determined interpretations. This is illustrated in Figure 2. For this purpose, the loss function consists of two parts: The prediction loss and the explanation loss. The explanation loss is based only on the node importance tensor produced by the network. A global sum pooling operation is used to turn the importance values of each separate channel into a single *alternate output tensor*  $\hat{\mathbf{Y}} \in \mathbb{R}^{B \times K}$ , where  $B$  is the training batch size. This alternate output tensor is then used to solve an approximation of the original prediction problem: This can be seen as a reduction of the problem into a set of  $K$  separate and independent subgraph counting problems, where each of those only uses the subset of training batch samples that aligns with the respective channel’s intended interpretation.

**Regression** For regression, we assume  $K = 2$ , where the first channel represents the negative and the second channel the positive influences relative to the reference value  $y_c$ , which is a hyperparameter of the model and usually set as the arithmetic mean of the target value distribution in the train set. We select all samples of the current training batch lesser and greater than the reference value and use these to calculate a mean squared error (MSE) loss:

$$\mathcal{L}_{\text{exp}} = \frac{1}{2 \cdot B} \sum_{b=1}^B \begin{cases} (\hat{\mathbf{Y}}_{b,0} - y_c - \mathbf{Y}_b^{\text{true}})^2 & \text{if } \mathbf{Y}_b^{\text{true}} < y_c \\ (\hat{\mathbf{Y}}_{b,1} - y_c - \mathbf{Y}_b^{\text{true}})^2 & \text{if } \mathbf{Y}_b^{\text{true}} > y_c \end{cases} \quad (8)$$

**Classification** We assume the number of channels  $K = C$  is equal to the number of possible output classes  $C$ . We use the alternate output channel to compute an individual binary cross entropy (BCE) loss for each channel:

$$\mathcal{L}_{\text{exp}} = \frac{1}{C \cdot B} \sum_{b=1}^B \sum_{c=1}^C \mathcal{L}_{\text{BCE}}(\mathbf{Y}_{b,c}^{\text{true}}, \hat{\mathbf{Y}}_{b,c}) \quad (9)$$

For regression as well as classification, the total loss during model training consists of these task-specific terms and an additional term for explanation sparsity:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{pred}} + \gamma \mathcal{L}_{\text{exp}} + \beta \mathcal{L}_{\text{sparsity}} \quad (10)$$

where  $\gamma$  and  $\beta$  are hyperparameters of the training process. Explanation sparsity  $\mathcal{L}_{\text{sparsity}}$  is calculated as L1 regularization over the node importance tensor. Based on this loss the gradients are calculated and the model weights are updated.

We will henceforth use the notation  $\text{MEGAN}_{\gamma}^K$  to refer to specific model configurations with  $K$  explanation channels,  $\gamma$  explanation co-training weight and use the superscript  $\text{MEGAN}^{(S)}$  to indicate when models were trained in an explanation-supervised fashion.

### 3.4 Multi-Channel Fidelity

A particular challenge in the field of explainable AI is the question of how to properly assess the quality of explanations [8]. One commonly used metric is the *fidelity* of explanations w.r.t. the model predictions. It quantifies the extent to which the explanation is responsible for the corresponding prediction. Yuan *et al.* [41] define the Fidelity<sup>+</sup> metric as the deviation of the predicted model output if all the nodes and edges that are part of the explanation are removed from the input. The reasoning is that the higher this resulting output deviation, the more important the explanation must have been for the original prediction. This metric is usually computed by setting all the features of the corresponding nodes and edges of the input graph to zero. However, one issue with this approach is that zero might be an in-distribution value for the input features. Therefore, the masked input elements may have an effect on the model that is different than their intended removal.

To address this issue we introduce the multi-channel Fidelity\* metric to assess the faithfulness of MEGAN’s predictions. Since our network directly incorporates the explanations into the prediction process as weights of the final global pooling operation, we can directly manipulate these explanations to quantify their impact on the prediction. This can be done by providing an additional importance mask  $\mathbf{M} \in [0, 1]^{V \times K}$  during the prediction of the network (see Figure 1). For each explanation channel  $k$ , we construct a mask  $\mathbf{M}^k$  which only suppresses that channel from the final pooling operation. The model is then queried with that mask to produce the modified output  $\hat{\mathbf{y}}^k$ , which we use to calculate the deviation  $\Delta^k = |\mathbf{y} - \hat{\mathbf{y}}^k|$  w.r.t. the original output. The fidelity is then calculated as:

$$\text{Fidelity}^* = \frac{1}{K} \sum_k \begin{cases} +\Delta^k & \text{if deviation as expected for channel } k \\ -\Delta^k & \text{if deviation not as expected for channel } k \end{cases} \quad (11)$$

What kind of deviation counts as *expected* for a given channel  $k$  is defined by the interpretation that is assigned to that channel. In the case of regression, for example, we assign the interpretation of the first explanation channel to be the negatively influencing evidence and the second channel to be the positively influencing evidence. In that case, if all the negative evidence is omitted from the result, it would be expected that the output becomes more positive than



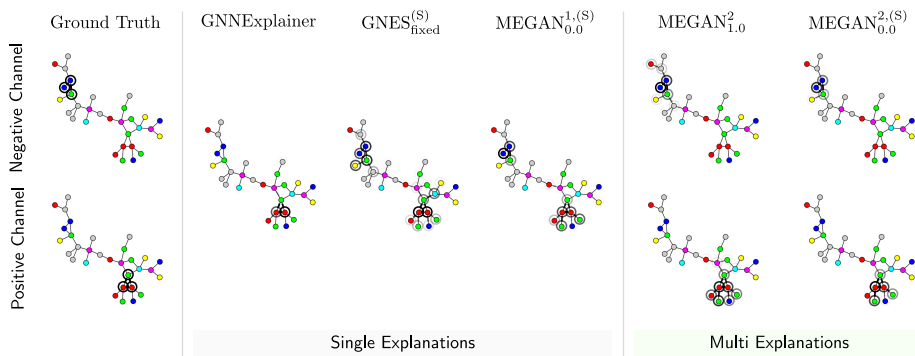


Fig. 3: Examples for explanations generated for one element of the RbMotifs dataset using selected methods. Explanations are represented as bold highlights of the corresponding graph elements. Left: The ground truth explanations split by the polarity of their influence on the graph target value. Middle: Explanations generated by some selected single-explanation methods. Right: Explanations generated by the multi-explanation MEGAN models.

the original prediction and vice versa. For classification on the other hand, if all evidence for one specific class is suppressed it would be expected that the confidence of that respective class decreases.

Consequently, a positive Fidelity\* indicates that the channels of the model generally have an effect on the prediction outcome that matches with their pre-defined interpretation.

## 4 Computational Experiments

We conduct computational experiments to demonstrate the capabilities of our network. Primarily, we emphasize two key strengths of our proposed model: (1) The inherent advantage of multi channel-explanations especially in regard to the interpretability of regression problems. On a specifically designed synthetic dataset we show that, unlike other post-hoc methods, by using explanation co-training our model is able to correctly capture the *polarity* of existing sub-graph evidence. (2) Our model’s significantly increased capability for explanation-supervised training, where our model correctly learns to replicate the ground truth explanations that it was trained on. Additionally, we conduct experiments with real-world graph classification and regression datasets that provide anecdotal evidence for the correctness of the model’s explanations for more complex tasks as well.

#### 4.1 Synthetic Graph Regression

We create a synthetic graph regression dataset called *RbMotifs* consisting of 5000 randomly generated graphs, where each node is associated with 3 node features representing an RGB color code. Graphs are additionally randomly seeded with specific simple sub-graph motifs, which either consist dominantly of red nodes or blue nodes. If a red-based motif exists within a graph, it contributes a constant positive value to the overall target value of a graph. Likewise, a blue-based motif contributes a negative value. Thus, the overall target value associated with each graph is the sum of all the sub-graph contributions and a small random value. The dataset represents a simple motif-based graph regression problem, where the individual sub-graph motifs are considered the perfect ground truth explanations. Most importantly, the explanations have a clear *opposing polarity* which is crucial to the understanding of the dataset’s underlying structure-property relationship.

**Single Explanations** Although many regression tasks may exhibit such explanations of different polarity, existing post-hoc attributional XAI methods are only able to provide a single explanation. These single explanations are only able to point out which parts of the graph are generally important for the prediction but do not capture in what manner they contribute to the outcome. Therefore, to compare our proposed MEGAN model to some established existing post-hoc explanation methods, we conduct a first experiment that only considers such single explanations. For this case, we concatenate all of the relevant sub-graphs into a single channel which will be considered the ground truth explanation for each element of the dataset.

We conduct the experiment for explanations obtained from Gradients [26], GN-NEExplainer [40], GNES [13] and MEGAN. For all the post-hoc methods we train a 3-layer GATv2 network as the basis for the explanations. The results of this experiment can be found in Table 1. We report on the overall prediction performance of the network, the explanation accuracy, the sparsity, and the fidelity of the explanations. The explanation accuracy is given as the node and edge AUROC score resulting from a comparison with the ground truth explanations, as it is proposed by McCloskey *et al.* [23]. The fidelity is given as the relative value  $\text{Fidelity}_{\text{rel}}^+$ , which is the difference between the predicted explanation’s fidelity and the fidelity of random explanations of the same sparsity (see Appendix A). In addition, we perform experiments with explanation supervision. To our knowledge, MEGAN and GNES are currently the only methods capable of explanation supervision for node and edge attributional explanations. For both of these cases, the models are trained with ground truth explanations in addition to the target values.

The results show that the explanations generated by all the methods achieve reasonable results for predictive performance, the node accuracy w.r.t. the explanation ground truth, as well as sparsity and fidelity. The explanation su-

Table 1: Results for 25 independent repetitions of the computational experiments on the RbMotifs dataset. We report the mean in black and the standard deviation in gray. The upper section contains results for the single-explanation experiments and the lower section for multi-explanation experiments. We highlight the best results in each section in bold and underline the second-best.

Explanations	$r^2 \uparrow$	Node AUC $\uparrow$	Edge AUC $\uparrow$	Sparsity $\downarrow$	Fidelity $^+_{\text{rel}} \uparrow$
Gradients	0.89 $\pm$ 0.05	0.73 $\pm$ 0.05	0.60 $\pm$ 0.03	0.12 $\pm$ 0.01	0.57 $\pm$ 0.14
GnnExplainer	0.89 $\pm$ 0.05	0.70 $\pm$ 0.04	0.52 $\pm$ 0.03	0.22 $\pm$ 0.06	<u>0.78</u> $\pm$ 0.20
GNES $^{(S)}$ <sub>original</sub>	0.88 $\pm$ 0.02	0.63 $\pm$ 0.04	0.58 $\pm$ 0.03	0.10 $\pm$ 0.01	0.50 $\pm$ 0.22
GNES $^{(S)}$ <sub>fixed</sub>	0.88 $\pm$ 0.02	<u>0.85</u> $\pm$ 0.04	0.66 $\pm$ 0.02	0.12 $\pm$ 0.01	0.74 $\pm$ 0.13
MEGAN $^1_{0.0}$	<u>0.92</u> $\pm$ 0.05	0.82 $\pm$ 0.12	<u>0.79</u> $\pm$ 0.08	0.14 $\pm$ 0.08	<b>1.10</b> $\pm$ 3.03
MEGAN $^{1,(S)}_{0.0}$	<b>0.95</b> $\pm$ 0.02	<b>0.98</b> $\pm$ 0.00	<b>0.99</b> $\pm$ 0.00	0.18 $\pm$ 0.00	0.53 $\pm$ 0.17
MEGAN $^2_{1.0}$	0.95 $\pm$ 0.01	<u>0.94</u> $\pm$ 0.02	<u>0.85</u> $\pm$ 0.06	0.10 $\pm$ 0.06	<u>2.06</u> $^{(*)}$ $\pm$ 0.85
MEGAN $^{2,(S)}_{0.0}$	<b>0.95</b> $\pm$ 0.03	<b>0.99</b> $\pm$ 0.00	<b>0.99</b> $\pm$ 0.00	0.09 $\pm$ 0.06	<b>2.11</b> $^{(*)}$ $\pm$ 0.36

<sup>(S)</sup> Explanation-supervised models. These models were trained on the ground truth explanation annotations in addition to the main target values.

<sup>(\*)</sup> Values of the multi-channel Fidelity\* metric. Note that these are *not* comparable to the other fidelity values obtained in a single channel setting.

pervised methods show the best results for explanation accuracy. The supervised MEGAN $^{1,(S)}_{0.0}$  model achieves a near-perfect accuracy, with the explanation-supervised GNES method being second-best.

The differences in prediction performance between the baseline methods and MEGAN models can be explained by the slightly different model architectures. However, one particularly interesting result is the small but significant performance difference between MEGAN $^1_{0.0}$  and the supervised MEGAN $^{1,(S)}_{0.0}$  version. In both cases, the same model architecture and hyperparameters are used, the only difference being that the latter additionally receives the explanatory information during training. This indicates that the explanations provide the model with some additional level of information about the task, which is useful for the main prediction task as well.

Aside from the numerical results, Figure 3 illustrates one example for these explanations. It shows that the single-explanation methods are able to capture the ground truth explanations to various degrees of success. However, in the presence of motifs with opposing influence, we often observe the issue that single-explanation methods focus on only one of these motifs and fail to high-

light the other. An example of this can be seen with the explanation generated by GNNExplainer in Figure 3, where it only highlights the positive explanation as being important. Although this is not always the case, we believe this effect contributes to the lower explanation accuracy results of these methods. Explanation-supervised training can be used to effectively counter this property, as is evident from the examples and the numerical results. However, even if all the explanatory motifs are correctly highlighted, we argue that single-explanations still don't provide the crucial information about *how* each motif contributes to the prediction outcome, as the polarity information cannot be retrieved from a single channel.

**Multi-Explanations** To demonstrate the advantages of multi-channel explanations, we conduct an experiment with the RbMotifs dataset, where the ground truth explanatory motifs of each graph are separated into two channels according to their influence on the target values. All blue-based motifs with a negative influence are sorted into one channel and all red-based motifs with positive influence are sorted into another.

We train two models to solve the prediction task: A two-channel  $\text{MEGAN}_{1,0}^2$  model, which uses explanation co-training to promote the generation of explanations according to the previously introduced explanations and a  $\text{MEGAN}_{0,0}^{2,(S)}$  which is explanation-supervised with the ground truth explanations instead. The results can be found in the lower section of Table 1.

Both models achieve nearly equal predictive performance, explanation sparsity, and Fidelity\*. The explanation-supervised model achieves near-perfect explanation accuracy for nodes and edges. However, the explanation co-training model also achieves a very good explanation accuracy. The right-hand side of Figure 3 shows an example of these results. As can be seen, both versions of the model are able to correctly capture the ground truth explanatory motifs according to their respective influence on the target value. The highly positive Fidelity\* results in both cases prove that both of the model's channels actually contribute to the prediction outcome according to their assigned interpretations of negative and positive influence. The results of this experiment present solid evidence that our proposed explanation co-training is an effective method to accurately capture the polarity of ground truth explanations even in the absence of ground truth explanations during training.

## 4.2 Real World Datasets

**MovieReviews - Sentiment Classification** The *MovieReviews* dataset is originally a natural language processing dataset from the ERASER benchmark [7] consisting of 2000 movie reviews from the IMDB database. The general sentiment of each review is labeled as either "positive" or "negative", where both classes are represented equally. Since this is a text classification dataset in its original form, we first process it in a manner similar to Rathee *et al.* [30]. First,

the raw strings are converted into token lists, where tokens are either words or other sentence elements such as punctuation. Each token is converted into a 50-dimensional feature vector through a pre-trained GLOVE model [25]. We finally convert the token list into a graph by applying a sliding window method, where each token is considered to be a node and connected to its four closest neighbors through an undirected edge.

We train a three-layer MEGAN<sub>1.0</sub><sup>2</sup> model to solve the binary sentiment classification task for each graph using the classification version of the explanation co-training procedure. The explanation co-training procedure promotes the first explanation channel of the network to contain evidence for the "negative" class label and the second channel for the "positive" class label.

In terms of classification performance our model achieves similar results (F1  $\approx$  0.85) as previously reported by Rathee *et al.* [30], who also use GNN and GLOVE embeddings. However, these results are significantly worse than results obtained with state-of-the-art NLP models, as they are for example reported by DeYoung *et al.* [7] (F1  $\approx$  0.92). We believe the main reason for this difference to be the use of the token embeddings derived from the 2014 GLOVE model. In the future, it would be interesting to see if GNNs could achieve competitive performance by using a state-of-the-art encoder such as BERT [6].

In regard to the generated explanations, Table 2 shows one example of a movie review. As can be seen, the model correctly learns negative adjectives such as "bad" as evidence for the "negative" class and positive adjectives such as "breathtaking" and "best" as evidence for the "positive" class. Despite this encouraging result, we still find there to be some errors in regard to the model's explanations about sentiment classification. On the one hand, the model also highlights unrelated words as explanations as well, such as "criminal" showing up as an explanation for negative reviews and "director" as positive evidence. On the other hand, the model is also not capable of accurately identifying negations and sarcasm to cause an inversion of sentiment.

**AqSolDB - Molecular Regression** The *AqSolDB* [35] dataset consists of roughly 10000 molecular graphs which are annotated with experimentally determined values of their water solubility. In chemistry, there exists some general intuition about what kinds of molecular structures are responsible for higher solubility values and which are responsible for lower ones. In a simplified manner, one can say that non-polar substructures such as carbon rings and long carbon chains generally result in lower solubility values, while polar structures such as certain nitrogen and oxygen functional groups are associated with higher values. In this experiment, we train a dual-channel three-layer MEGAN<sub>1.0</sub><sup>2</sup> model to predict the continuous solubility values for the molecular graphs. We make use of the previously described regression version of the co-training procedure, which promotes the first channel to highlight negatively influencing motifs and the second channel to highlight positively influencing motifs. Additionally, we train a comparable GATv2 model on the solubility dataset as well and use GNNEx-

Table 2: Example explanations generated for both sentiment classes for a review about the movie "Avengers Endgame". Larger importance values are represented by stronger color highlights.

Negative	Positive
overall avengers endgame was a remarkable	overall avengers endgame was a remarkable
movie and a worthy culmination of the mcu up	movie and a worthy culmination of the mcu up
to this point there were some genuinely	to this point there were some genuinely
heartbreaking moments and breathtaking action	heartbreaking moments and breathtaking action
sequences but to be honest some of the movies	sequences but to be honest some of the movies
i had to sit through to get here were not worth	i had to sit through to get here were not worth
it some of the early mcu movies and series	it some of the early mcu movies and series
leading up to this finale i found rather bland	leading up to this finale i found rather bland
unfunny and sometimes just downright bad but	unfunny and sometimes just downright bad but
this movie was one of the best movies i have	this movie was one of the best movies i have
seen in a while	seen in a while

Table 3: Results for 5 independent repetitions of the experiments with the Aq-SolDB dataset for water solubility. We report the mean in black and the standard deviation in gray.

Model	$R^2 \uparrow$	Sparsity $\downarrow$	Fidelity <sup>(*)</sup> $\uparrow$
GNNX+GATv2	0.93 $\pm$ 0.01	0.34 $\pm$ 0.27	1.26 $\pm$ 0.90
MEGAN <sub>1.0</sub> <sup>2</sup>	0.93 $\pm$ 0.01	0.22 $\pm$ 0.14	2.50 <sup>(*)</sup> $\pm$ 2.29
Consensus Model <sup>†</sup>	0.93	-	-

<sup>†</sup> Previously published results by Sorkun *et al.* [35].

<sup>(\*)</sup> Multi-explanation case measures Fidelity\* metric

plainer to produce single explanations as a comparison.

Both the MEGAN model and the GATv2 model are able to match the predictive performance which was previously reported in the literature by Sorkunen *et al.* [36]. Both approaches also generate explanations with low sparsity and high fidelity values, as it can be seen in Table 3. Figure 4 illustrates some example explanations generated by MEGAN and GNNExplainer. The examples show that the explanations generated by MEGAN match the general human intuition about the structure-property relationships of water solubility. Large carbon structures are consistently highlighted in the negative explanation channel. The positive explanation channel on the other hand mostly contains polar nitrogen and oxygen functional groups. The explanations generated by GNNExplainer on the basis

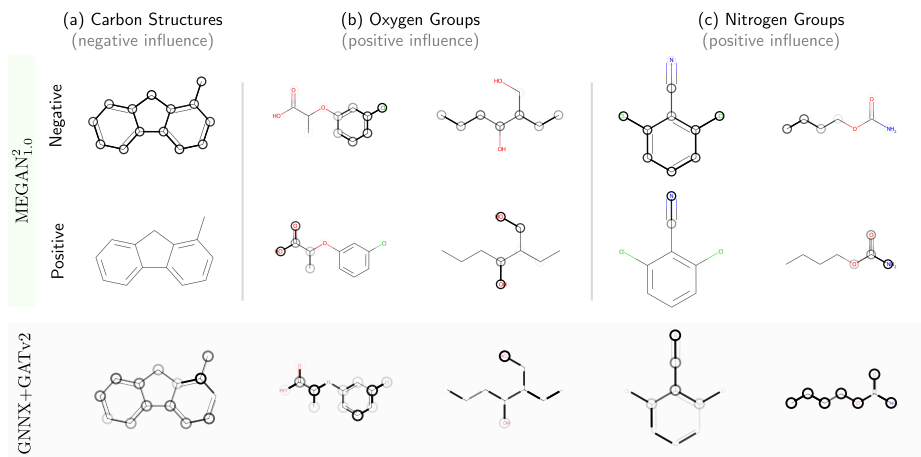


Fig. 4: Example explanations generated by MEGAN and GNNExplainer for the prediction of water solubility. Explanations are represented as bold highlights of the corresponding graph elements. Explanations are represented as bold highlights of the corresponding graph elements. (a) Examples of molecules dominated by large carbon structures which are known as negative influences on water solubility. (b) Examples of molecules containing oxygen functional groups which are known to be a positive influence on water solubility. (c) Examples of molecules containing nitrogen groups which are also known as positive influences.

of the GATv2 model, however, do not show any such discernible pattern.

Despite an equally high predictive performance and high explanation fidelity, we argue that the single-explanation case contributes significantly less useful information for a human understanding of the predictions. We think this example reinforces the importance of the multi-explanation approach, especially for graph regression problems. By considering the polarity of structure-property explanations in graph regression problems, the MEGAN model is able to provide explanations that are more consistent with human intuition and are thus more interpretable.

**TADF - Molecular Regression** Previous experiments were able to provide exemplary evidence for the correctness of MEGAN’s explanations through real-world datasets for which human intuition exists. In this final experiment, we choose a dataset where almost no human intuition exists to investigate potential applications to reveal novel insights about structure-property relationships. The *TADF* dataset consists of roughly half a million molecular graphs. Target value annotations were during a high-throughput virtual screening experiment conducted by Gómez-Bombarelli *et al.* [14] with the objective to discover novel materials for an application in OLED technology. Specifically, the authors aimed

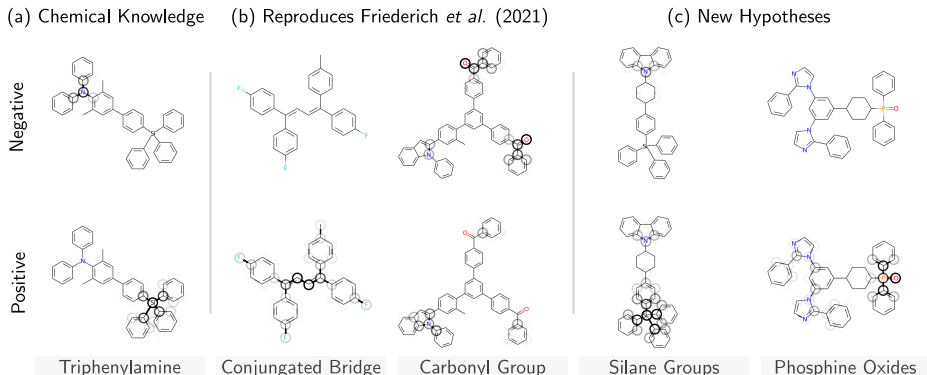


Fig. 5: Example explanations obtained from the MEGAN model for the prediction of the singlet-triplet energy gap of the TADF dataset. (a) Explanations that reproduce known chemical intuition about the task. (b) Explanations that reproduce hypotheses previously published by Friederich *et al.* [11]. (c) New explanatory sub-graph motifs proposed through an observation of the explanations generated by MEGAN.

to discover materials that show a specific characteristic of thermally delayed fluorescence (TADF). This class of materials is a promising approach to avoid the high cost of typically used phosphorescent OLED materials [9,42]. Along the delayed fluorescent rate constant  $k_{\text{TADF}}$ , the elements of the dataset are annotated with the singlet-triplet energy gap  $\Delta E_{\text{st}}$  and the oscillator strength  $f$ .

In this experiment, we train a three-layer MEGAN<sub>1.0</sub><sup>2</sup> model to estimate the singlet-triplet gap  $\Delta E_{\text{st}}$  for each element. As before, the explanation co-training promotes the first channel to contain the negative influences and the second channel to contain the positive influences.

Our model achieves overall good predictivity ( $R^2 \approx 0.90$ ) for the main prediction task and a positive Fidelity\* value validating that the individual channels indeed affect the model prediction according to their pre-determined interpretations. Figure 5 illustrates some example explanations obtained from the model. Most importantly, we show that our model is able to replicate one of the few known structure-property relationships about the singlet-triplet energy. Triphenylamine bridges are known to be associated with low energy gaps, as they cause the necessary twist angles between the fragments, decoupling electron-donating and electron-accepting parts of a molecule to reduce the exchange interaction between the frontier orbitals which would otherwise lower the triplet state compared to the singlet state, thus preventing undesired singlet-triplet splittings. This fact is reflected in Figure 5(a), where a triphenylamine bridge is highlighted as a negative influence on the prediction outcome. Furthermore, our model is able to support hypotheses published in previous work by Friederich *et al.* [11], who use an interpretable decision tree method to generate explanation hypotheses



for the same task. As shown in Figure 5(b) our model replicates their findings of conjugated bridges as a positive influence on the energy gap and carbonyl groups as a negative influence. Beyond that, our model finds several novel hypotheses about structure-property relationships, two of which are shown in Figure 5(c): We can propose silane groups and phosphine oxides as positive influences to the singlet-triplet energy gap.

## 5 Limitations

Despite the encouraging experimental results, there are limitations to the proposed MEGAN architecture: Firstly, there is no hard guarantee that each channel’s explanations align correctly according to their pre-determined interpretations. This alignment is mainly promoted through the explanation co-training, whose influence on the network is dependent on a hyperparameter. We occasionally observed ”explanation leakage” and ”explanation flipping” during training. In those rare cases, explanations factually belonging to one channel may either faintly appear in the opposite channel or a particularly disadvantageous initialization of the network causes explanations to develop in the exact opposite channel relative to their assigned interpretation. Ultimately, the alignment of a particular channel with its intended interpretation has to be tested through a Fidelity\* analysis after the model training.

The second limitation is in the design of the explanation co-training itself, which essentially reduces the problem to a subgraph counting task. While there are many important real-world applications that can be approximated as such, it still presents an important limit to the expressiveness of the models produced by our model.

## 6 Conclusion

In this work, we introduce the self-explaining multi-explanation graph attention network (MEGAN) architecture, which produces node and edge attributional explanations for graph regression and classification tasks. Our model implements the number  $K$  of generated explanations as a hyperparameter of the network itself, instead of being dependent on the task specification. Based on several exemplary synthetic and real-world datasets, we show that this property is especially crucial for graph regression problems. By being able to generate attributional explanations for a single regression target along multiple explanation channels, our model is able to account for the *polarity* of explanations. In many graph regression applications certain sub-graph motifs influence the predicted outcome in opposing directions: Some motifs present a negative influence on the overall prediction, while others are a positive influence. We achieve the alignment of the model’s multiple explanation channels according to these pre-determined interpretations by introducing an explanation co-training procedure. Beside the main prediction loss, an additional explanation loss is generated from an approximate solution of the prediction problem based only on each channels

explanation masks. We can validate the channel’s alignment to their respective intended interpretations through the Fidelity\* metric, which extends the concept of explanation fidelity to our multi-channel case.

Additionally, we demonstrate the capabilities of our model for explanation-supervised training, where a model is trained to produce explanations based on a set of given ground truth explanations. For a synthetic graph regression dataset, we show that our model is able to learn the given ground truth explanations almost perfectly, significantly outperforming an existing baseline method from literature.

One particularly interesting result is the improvement of the prediction performance for the explanation-supervised training during the first synthetic experiment but not during the second one. Similar effects have already been shown in the domain of image processing, where various authors are able to demonstrate a performance increase when models are additionally trained to emulate human saliency maps [19,2]. One promising direction for future work will be to investigate the conditions under which (human) explanations have the potential to improve predictive performance for graph-related tasks as well.

## 7 Reproducibility Statement

We make our experimental code publically available at [https://github.com/aimat-lab/graph\\_attention\\_student](https://github.com/aimat-lab/graph_attention_student). The code is implemented in the Python 3.9 programming language. Our neural networks are built with the KGCNN library by Reiser *et al.* [31], which provides a framework for graph neural network implementations with TensorFlow and Keras. We make all data used in our experiments publically available on a file share provider <https://bwsyncandshare.kit.edu/s/E3MynrfQsLAHzJC>. The datasets can be loaded, processed, and visualized with the visual graph datasets package [https://github.com/aimat-lab/visual\\_graph\\_datasets](https://github.com/aimat-lab/visual_graph_datasets). All experiments were performed on a system with the following specifications: Ubuntu 22.04 operating system, Ryzen 9 5900 processor, RTX 2060 graphics card and 80GB of memory. We have aimed to package the various experiments as independent modules and our code repository contains a brief explanation of how these can be executed.

## References

1. Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., Kim, B.: Sanity Checks for Saliency Maps. In: *Advances in Neural Information Processing Systems*. vol. 31. Curran Associates, Inc. (2018), [https://proceedings.neurips.cc/paper\\_files/paper/2018/hash/294a8ed24b1ad22ec2e7efea049b8737-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2018/hash/294a8ed24b1ad22ec2e7efea049b8737-Abstract.html)
2. Boyd, A., Tinsley, P., Bowyer, K., Czajka, A.: CYBORG: Blending Human Saliency Into the Loss Improves Deep Learning (Aug 2022). <https://doi.org/10.48550/arXiv.2112.00686>, <http://arxiv.org/abs/2112.00686>, arXiv:2112.00686 [cs]
3. Brody, S., Alon, U., Yahav, E.: How Attentive are Graph Attention Networks? (Feb 2022), <https://openreview.net/forum?id=F72ximsx7C1>
4. Dai, E., Wang, S.: Towards Self-Explainable Graph Neural Network. In: *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. pp. 302–311. CIKM '21, Association for Computing Machinery, New York, NY, USA (Oct 2021). <https://doi.org/10.1145/3459637.3482306>, <https://doi.org/10.1145/3459637.3482306>
5. Delaney, J.S.: ESOL: Estimating Aqueous Solubility Directly from Molecular Structure. *Journal of Chemical Information and Computer Sciences* **44**(3), 1000–1005 (May 2004). <https://doi.org/10.1021/ci034243x>, <https://doi.org/10.1021/ci034243x>, publisher: American Chemical Society
6. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>
7. DeYoung, J., Jain, S., Rajani, N.F., Lehman, E., Xiong, C., Socher, R., Wallace, B.C.: ERASER: A Benchmark to Evaluate Rationalized NLP Models. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 4443–4458. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.408>, <https://aclanthology.org/2020.acl-main.408>
8. Doshi-Velez, F., Kim, B.: Towards A Rigorous Science of Interpretable Machine Learning. arXiv:1702.08608 [cs, stat] (Mar 2017), <http://arxiv.org/abs/1702.08608>, arXiv: 1702.08608
9. Endo, A., Sato, K., Yoshimura, K., Kai, T., Kawada, A., Miyazaki, H., Adachi, C.: Efficient up-conversion of triplet excitons into a singlet state and its application for organic light emitting diodes. *Applied Physics Letters* **98**(8), 083302 (Feb 2011). <https://doi.org/10.1063/1.3558906>, <https://aip.scitation.org/doi/full/10.1063/1.3558906>, publisher: American Institute of Physics
10. Fernandes, P., Treviso, M., Pruthi, D., Martins, A.F.T., Neubig, G.: Learning to Scaffold: Optimizing Model Explanations for Teaching (Apr 2022). <https://doi.org/10.48550/arXiv.2204.10810>, <http://arxiv.org/abs/2204.10810>, arXiv:2204.10810 [cs]
11. Friederich, P., Krenn, M., Tamblyn, I., Aspuru-Guzik, A.: Scientific intuition inspired by machine learning-generated hypotheses. *Machine Learning: Science and Technology* **2**(2), 025027 (Apr 2021). <https://doi.org/10.1088/2632-2153/abda08>, <https://doi.org/10.1088/2632-2153/abda08>, publisher: IOP Publishing

12. Funke, T., Khosla, M., Rathee, M., Anand, A.: ZORRO: Valid, Sparse, and Stable Explanations in Graph Neural Networks. *IEEE Transactions on Knowledge and Data Engineering* pp. 1–12 (2023). <https://doi.org/10.1109/TKDE.2022.3201170>, <https://ieeexplore.ieee.org/document/9866587/>
13. Gao, Y., Sun, T., Bhatt, R., Yu, D., Hong, S., Zhao, L.: GNES: Learning to Explain Graph Neural Networks. In: 2021 IEEE International Conference on Data Mining (ICDM). pp. 131–140 (Dec 2021). <https://doi.org/10.1109/ICDM51629.2021.00023>, ISSN: 2374-8486
14. Gómez-Bombarelli, R., Aguilera-Iparraguirre, J., Hirzel, T.D., Duvenaud, D., Maclaurin, D., Blood-Forsythe, M.A., Chae, H.S., Einzinger, M., Ha, D.G., Wu, T., Markopoulos, G., Jeon, S., Kang, H., Miyazaki, H., Numata, M., Kim, S., Huang, W., Hong, S.I., Baldo, M., Adams, R.P., Aspuru-Guzik, A.: Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach. *Nature Materials* **15**(10), 1120–1127 (Oct 2016). <https://doi.org/10.1038/nmat4717>, <https://www.nature.com/articles/nmat4717>, number: 10 Publisher: Nature Publishing Group
15. Henderson, R., Clevert, D.A., Montanari, F.: Improving Molecular Graph Neural Network Explainability with Orthonormalization and Induced Sparsity. In: Proceedings of the 38th International Conference on Machine Learning. pp. 4203–4213. PMLR (Jul 2021), <https://proceedings.mlr.press/v139/henderson21a.html>, ISSN: 2640-3498
16. Huang, Q., Yamada, M., Tian, Y., Singh, D., Chang, Y.: GraphLIME: Local Interpretable Model Explanations for Graph Neural Networks. *IEEE Transactions on Knowledge and Data Engineering* pp. 1–6 (2022). <https://doi.org/10.1109/TKDE.2022.3187455>, conference Name: IEEE Transactions on Knowledge and Data Engineering
17. Jiménez-Luna, J., Grisoni, F., Schneider, G.: Drug discovery with explainable artificial intelligence. *Nature Machine Intelligence* **2**(10), 573–584 (Oct 2020). <https://doi.org/10.1038/s42256-020-00236-4>, <https://www.nature.com/articles/s42256-020-00236-4>
18. Kindermans, P.J., Hooker, S., Adebayo, J., Alber, M., Schütt, K.T., Dähne, S., Erhan, D., Kim, B.: The (Un)reliability of Saliency Methods. In: Samek, W., Montavon, G., Vedaldi, A., Hansen, L.K., Müller, K.R. (eds.) *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pp. 267–280. *Lecture Notes in Computer Science*, Springer International Publishing, Cham (2019). [https://doi.org/10.1007/978-3-030-28954-6\\_14](https://doi.org/10.1007/978-3-030-28954-6_14), [https://doi.org/10.1007/978-3-030-28954-6\\_14](https://doi.org/10.1007/978-3-030-28954-6_14)
19. Linsley, D., Shiebler, D., Eberhardt, S., Serre, T.: Learning what and where to attend (2019), <https://openreview.net/forum?id=BJgLg3R9KQ>
20. Luo, D., Cheng, W., Xu, D., Yu, W., Zong, B., Chen, H., Zhang, X.: Parameterized Explainer for Graph Neural Network. *ArXiv* (Nov 2020), <https://www.semanticscholar.org/paper/Parameterized-Explainer-for-Graph-Neural-Network-Luo-Cheng/d9f5ec342df97e060b527a8bc18ae4e97401f246>
21. Magister, L.C., Barbiero, P., Kazhdan, D., Siciliano, F., Ciravegna, G., Silvestri, F., Jamnik, M., Lio, P.: Encoding Concepts in Graph Neural Networks (Aug 2022). <https://doi.org/10.48550/arXiv.2207.13586>, <http://arxiv.org/abs/2207.13586>, arXiv:2207.13586 [cs]
22. Magister, L.C., Kazhdan, D., Singh, V., Liò, P.: GCEXplainer: Human-in-the-Loop Concept-based Explanations for Graph Neural Networks (Jul 2021).

- <https://doi.org/10.48550/arXiv.2107.11889>, <http://arxiv.org/abs/2107.11889>, arXiv:2107.11889 [cs]
23. McCloskey, K., Taly, A., Monti, F., Brenner, M.P., Colwell, L.J.: Using attribution to decode binding mechanism in neural network models for chemistry. *Proceedings of the National Academy of Sciences* **116**(24), 11624–11629 (Jun 2019). <https://doi.org/10.1073/pnas.1820657116>, <https://www.pnas.org/doi/10.1073/pnas.1820657116>, publisher: Proceedings of the National Academy of Sciences
  24. Müller, P., Faber, L., Martinkus, K., Wattenhofer, R.: DT+GNN: A Fully Explainable Graph Neural Network using Decision Trees (May 2022). <https://doi.org/10.48550/arXiv.2205.13234>, <http://arxiv.org/abs/2205.13234>, arXiv:2205.13234 [cs]
  25. Pennington, J., Socher, R., Manning, C.: GloVe: Global Vectors for Word Representation. In: *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 1532–1543. Association for Computational Linguistics, Doha, Qatar (Oct 2014). <https://doi.org/10.3115/v1/D14-1162>, <https://aclanthology.org/D14-1162>
  26. Pope, P.E., Kolouri, S., Rostami, M., Martin, C.E., Hoffmann, H.: Explainability Methods for Graph Convolutional Neural Networks. pp. 10772–10781 (2019), [https://openaccess.thecvf.com/content\\_CVPR\\_2019/html/Pope\\_Explainability\\_Methods\\_for\\_Graph\\_Convolutional\\_Neural\\_Networks\\_CVPR\\_2019\\_paper.html](https://openaccess.thecvf.com/content_CVPR_2019/html/Pope_Explainability_Methods_for_Graph_Convolutional_Neural_Networks_CVPR_2019_paper.html)
  27. Prado-Romero, M.A., Stilo, G.: GRETEL: Graph Counterfactual Explanation Evaluation Framework. In: *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*. pp. 4389–4393. CIKM '22, Association for Computing Machinery, New York, NY, USA (Oct 2022). <https://doi.org/10.1145/3511808.3557608>, <https://dl.acm.org/doi/10.1145/3511808.3557608>
  28. Pruthi, D., Gupta, M., Dhingra, B., Neubig, G., Lipton, Z.C.: Learning to Deceive with Attention-Based Explanations. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. pp. 4782–4793. Association for Computational Linguistics, Online (Jul 2020). <https://doi.org/10.18653/v1/2020.acl-main.432>, <https://aclanthology.org/2020.acl-main.432>
  29. Qiao, T., Dong, J., Xu, D.: Exploring Human-Like Attention Supervision in Visual Question Answering. In: *Thirty-Second AAAI Conference on Artificial Intelligence* (Apr 2018), <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16485>
  30. Rathee, M., Funke, T., Anand, A., Khosla, M.: BAGEL: A Benchmark for Assessing Graph Neural Network Explanations (Jun 2022). <https://doi.org/10.48550/arXiv.2206.13983>, <http://arxiv.org/abs/2206.13983>, arXiv:2206.13983 [cs]
  31. Reiser, P., Eberhard, A., Friederich, P.: Graph neural networks in TensorFlow-Keras with RaggedTensor representation (kgcnn). *Software Impacts* **9**, 100095 (Aug 2021). <https://doi.org/10.1016/j.simpa.2021.100095>, <https://www.sciencedirect.com/science/article/pii/S266596382100035X>
  32. Sanchez-Lengeling, B., Wei, J., Lee, B., Reif, E., Wang, P., Qian, W., McCloskey, K., Colwell, L., Wiltchko, A.: Evaluating Attribution for Graph Neural Networks. In: *Advances in Neural Information Processing Systems*. vol. 33, pp. 5898–5910. Curran Associates, Inc. (2020), <https://proceedings.neurips.cc/paper/2020/hash/417fbf2e9d5a28a855a11894b2e795a-Abstract.html>

33. Schwarzenberg, R., Hübner, M., Harbecke, D., Alt, C., Hennig, L.: Layerwise Relevance Visualization in Convolutional Text Graph Classifiers. In: Proceedings of the Thirteenth Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-13). pp. 58–62. Association for Computational Linguistics, Hong Kong (Nov 2019). <https://doi.org/10.18653/v1/D19-5308>, <https://aclanthology.org/D19-5308>
34. Shin, Y.M., Kim, S.W., Shin, W.Y.: PAGE: Prototype-Based Model-Level Explanations for Graph Neural Networks (Oct 2022). <https://doi.org/10.48550/arXiv.2210.17159>, <http://arxiv.org/abs/2210.17159>, arXiv:2210.17159 [cs, math]
35. Sorkun, M.C., Khetan, A., Er, S.: AqSolDB, a curated reference set of aqueous solubility and 2D descriptors for a diverse set of compounds. *Scientific Data* **6**(1), 143 (Aug 2019). <https://doi.org/10.1038/s41597-019-0151-1>, <https://www.nature.com/articles/s41597-019-0151-1>, number: 1 Publisher: Nature Publishing Group
36. Sorkun, M.C., Koelman, J.M.V.A., Er, S.: Pushing the limits of solubility prediction via quality-oriented data selection. *iScience* **24**(1), 101961 (Jan 2021). <https://doi.org/10.1016/j.isci.2020.101961>, <https://www.sciencedirect.com/science/article/pii/S2589004220311585>
37. Stacey, J., Belinkov, Y., Rei, M.: Supervising Model Attention with Human Explanations for Robust Natural Language Inference (May 2022). <https://doi.org/10.48550/arXiv.2104.08142>, <http://arxiv.org/abs/2104.08142>, arXiv:2104.08142 [cs]
38. Tan, J., Geng, S., Fu, Z., Ge, Y., Xu, S., Li, Y., Zhang, Y.: Learning and Evaluating Graph Neural Network Explanations based on Counterfactual and Factual Reasoning. In: Proceedings of the ACM Web Conference 2022. pp. 1018–1027. WWW '22, Association for Computing Machinery, New York, NY, USA (Apr 2022). <https://doi.org/10.1145/3485447.3511948>, <https://doi.org/10.1145/3485447.3511948>
39. Wu, Z., Ramsundar, B., Feinberg, E.N., Gomes, J., Geniesse, C., Pappu, A.S., Leswing, K., Pande, V.: MoleculeNet: a benchmark for molecular machine learning. *Chemical Science* **9**(2), 513–530 (Jan 2018). <https://doi.org/10.1039/C7SC02664A>, <https://pubs.rsc.org/en/content/articlelanding/2018/sc/c7sc02664a>, publisher: The Royal Society of Chemistry
40. Ying, Z., Bourgeois, D., You, J., Zitnik, M., Leskovec, J.: GNNExplainer: Generating Explanations for Graph Neural Networks. In: Advances in Neural Information Processing Systems. vol. 32. Curran Associates, Inc. (2019), <https://papers.nips.cc/paper/2019/hash/d80b7040b773199015de6d3b4293c8ff-Abstract.html>
41. Yuan, H., Yu, H., Gui, S., Ji, S.: Explainability in Graph Neural Networks: A Taxonomic Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* pp. 1–19 (2022). <https://doi.org/10.1109/TPAMI.2022.3204236>, conference Name: IEEE Transactions on Pattern Analysis and Machine Intelligence
42. Zhang, Q., Li, J., Shizu, K., Huang, S., Hirata, S., Miyazaki, H., Adachi, C.: Design of Efficient Thermally Activated Delayed Fluorescence Materials for Pure Blue Organic Light Emitting Diodes. *Journal of the American Chemical Society* **134**(36), 14706–14709 (Sep 2012). <https://doi.org/10.1021/ja306538w>, <https://doi.org/10.1021/ja306538w>, publisher: American Chemical Society
43. Zhang, Z., Liu, Q., Wang, H., Lu, C., Lee, C.: ProtGNN: Towards Self-Explaining Graph Neural Networks. Proceedings of the AAAI Conference on Artificial Intelligence **36**(8), 9127–9135 (Jun 2022). <https://doi.org/10.1609/aaai.v36i8.20898>, <https://ojs.aaai.org/index.php/AAAI/article/view/20898>, number: 8

## A Evaluation Metrics

**Fidelity** Fidelity metrics are used to quantify the degree to which explanations are actually responsible for a model’s prediction. In our experiments, we use the definition of the Fidelity<sup>+</sup> metric as defined by Yuan *et al.* [41]. It is calculated as the difference between the original predicted value and the predicted value if the elements of the explanation are removed from the input graph. It is generally assumed the higher this value, the more important those elements are for the prediction. This metric generally works well by itself for classification problems, where confidence values are limited to the range between 0 and 1. In such a case, a fidelity value of 0.8 would be considered quite high because there exists a frame of reference that defines 1 as the maximum possible value. However, for this reason, we find that the metric is not immediately applicable to the regression problems since there exists no frame of reference as to what would be considered a particularly high or low value.

Instead, for our regression experiments, we use a relative fidelity value which is defined relative to a point of reference.

$$\text{Fidelity}_{\text{rel}}^+ = \text{Fidelity}^+ - \text{Fidelity}_{\text{random}}^+ \quad (12)$$

As the frame of reference, we use the fidelity value which results from a purely random input graph mask, which has the *same sparsity* as the given explanation. The random fidelity value is calculated as the arithmetic mean resulting from 10 such randomly sampled input masks per explanation.

## B GNES Implementation

In our experiments, we use the GNES method by Gao *et al.* [13] as a baseline approach from the literature that supports explanation supervision. In their framework, the authors propose using existing differentiable post-hoc explanation methods for explanation supervision. For that, they introduce a generic framework to describe node and edge attributional explanations. For example, they define node the attributional explanation for node  $n$  at layer  $l$  as

$$M_n^{(l)} = \left\| \text{ReLU}\left(g\left(\frac{\partial y_c}{\partial F_n^{(l)}}\right) \cdot h(F_n^{(l)})\right) \right\| \quad (13)$$

where  $F_n^{(l)}$  is the activation of node  $n$  at layer  $l$ .  $g(\cdot)$  and  $h(\cdot)$  are generic functions that can be defined for specific implementations of explanation methods. Edge explanations are defined in a similarly generic way. Explanation supervision is then achieved through additional loss MAE loss terms between these generated explanations and the given reference explanations.

For our experiments, we were not able to use the original code at <https://github.com/YuyangGao/GNES> as that implementation only supports binary classification problems and is limited to a batch size of 1. We re-implement their method

in the KGCNN framework. We follow the original paper as closely as possible for the version we call  $\text{GNES}_{\text{original}}$ . However, we find that the used  $\text{ReLU}(\cdot)$  operation does not work well with regression operations as it cuts off negative values and thus actively discards explanatory motifs with *opposing influence*. Consequently, we modify the method to use an absolute value operation  $\|\cdot\|$  instead of the  $\text{ReLU}(\cdot)$  for the version we call  $\text{GNES}_{\text{fixed}}$ . We find that this version works much better with regression tasks as it is able to properly account for positive and negative influences.

## C GNN Benchmarks

Aside from its capability for explanation supervision, we also find that our model generally shows a good prediction performance as well, when compared to other state-of-the-art GNNs. Figure 6 shows the benchmarking results of the MEGAN model compared to several other GNNs from the literature for two datasets of molecular property prediction. The benchmarking results were obtained from the KGCNN library [https://github.com/aimat-lab/gcmn\\_keras/tree/master/training/results](https://github.com/aimat-lab/gcmn_keras/tree/master/training/results). To produce the results, all models were subjected to a cursory hyperparameter optimization on the respective datasets. The MEGAN models trained for this comparison use neither explanation supervision nor the co-training method.

The results show that MEGAN achieves the second-best results for both tasks.

(a) Results for the ESOL dataset [5] which consists of 1128 molecular graphs and their respective values for water solubility.

Model	MAE ↓	RMSE ↓
GAT	0.49±0.02	0.70±0.04
GIN	0.50±0.02	0.70±0.03
CMPNN	0.48±0.03	0.68±0.02
INorp	0.49±0.01	0.68±0.03
GATv2	0.47±0.03	0.67±0.03
Schnet	0.46±0.03	0.65±0.04
DMPNN	0.45±0.02	0.63±0.02
AttentiveFP	0.46±0.01	0.63±0.03
MEGAN	0.44±0.03	0.60±0.05
PAiNN	0.43±0.02	0.60±0.02

(b) Results for the LIPOP dataset [39] which consists of 4200 molecular graphs and their respective octanol/water distribution coefficient.

Model	MAE ↓	RMSE ↓
GAT	0.50±0.02	0.70±0.04
INorp	0.46±0.01	0.65±0.01
Schnet	0.48±0.00	0.65±0.00
GIN	0.45±0.01	0.64±0.03
AttentiveFP	0.45±0.01	0.62±0.01
GATv2	0.41±0.01	0.59±0.01
PAiNN	0.40±0.01	0.58±0.03
CMPNN	0.41±0.01	0.58±0.01
MEGAN	0.40±0.01	0.56±0.01
DMPNN	0.38±0.01	0.55±0.03

Fig. 6: Benchmarking results obtained from the KGCNN library from a random 5-fold cross-validation.