

# Augment to Interpret: Supplementary Materials

## 1. Reusing a Module with Batch Normalisation

All our competitors use approximately the same GIN architecture for their experimentation (Suresh et al., 2021; Gao et al., 2022; Miao et al., 2022). Specifically, they all use batch normalisation layers, and so do we. However, considering the context it is used in, we argue that the way it should be used to be mathematically correct is not as straightforward as the way it is used in these approaches. This may have unexpected effects on their results and conclusions. Therefore, we propose a batch-norm switch to fix the issue, and we use it for our approach. As this is not within the scope of this paper, only a preliminary experiment is presented to demonstrate the relevance of a switch in practice. Further investigation into alternative methods to address the problem and the resulting effects on performances is left for future work.

### 1.1. Theoretical Analysis

In more detail, all mentioned approaches have one point in common: they use the same neural network to embed both the original graph (*raw*) and an augmented view of it (*aug*). However, the original graph and its augmented view come from different distributions, with significant differences. For example, the augmented view usually has lower node degrees, as only a subset of the original graph edges is kept. Ultimately, this difference has an impact on the input of the batch normalisation layers of this neural network.

As a reminder, the batch normalisation layer uses the batch mean and variance to normalise its input, and stores a running mean of these values to be used at test time as an estimate of said values. That is, the normalisation performed at train time is meant to be roughly equivalent to the normalisation performed at inference time. In that way, given an input, it produces roughly the same output, whether it is in train or in inference mode. This is important as the end of the network is trained on that output.

In our case, the mean and variance of the input of the layer are different whether the input of the network is a raw graph ( $\mu_{raw}, \sigma_{raw}$ ) or an augmented view of it ( $\mu_{aug}, \sigma_{aug}$ ). At train time, the normalisation will therefore be different in both cases. However, at our competitors' inference time, the estimate to perform the normalisation is something close to the average of the real values ( $\frac{1}{2}(\mu_{raw} + \mu_{aug}), \frac{1}{2}(\sigma_{raw} + \sigma_{aug})$ ), for both the original graph and its augmented view. As a result, if the means and variances are very different, the outputs will not be similar in any way to the ones at train time. The final output of the network will therefore be highly unpredictable.

## 1.2. Batch-Norm Switch and Preliminary Experiment

A way to avoid this issue is to have two different running means: one for the original graph distribution, and one for the augmented graph distribution. We call this solution a batch-norm switch, as it switches between running means. With a batch-norm switch, the training of the model remains the same as without the switch, but the inference is mathematically correct. To assess the relevance of the switch, a quick experiment is presented hereunder, based on the code provided by GSAT authors for the *BA2Motifs* dataset with a *GIN* architecture, using seed 0 only.

In GSAT, the issue is mitigated by having augmented views close to the original graphs, i.e. removing on average less than half the edges. It is further mitigated by an early stopping criterion based on a metric on the validation set, which is unlikely to give good results if the network produces unpredictable results. The early stopping will therefore tend to select epochs for which the issue has no negative impact.

For our analysis, we change these two parameters. First, we run the experiments for 500 epochs instead of 100 and show the full curves observed during training rather than the final result of the selected epoch. Second, for half the experiments, we use a final  $r$  value of 0.1 instead of 0.5. This means that starting from epoch 80, roughly 10% of edges are kept rather than 50%. Note the first 40 epochs are not affected by this second change, as GSAT adopt a step decay for  $r$ . The result is shown in Figure 1, measured on the test set.

As can be seen, the downstream ACC is greatly improved by the switch. On the other hand, the interpretability AUC is improved by the switch when interpretations are sparse, but it is not as good when interpretations are not sparse. We think it is because the unexpected effect aligns with some dataset properties. For example, it could favour highly connected interpretations, which is good in the case of the *BA2Motifs* dataset but may not be for other datasets.

Additionally, we can see that the model with a switch and sparse interpretations (the orange curves) has on average almost null attention weights on non-important edges (0.0016 on average at the end), but significant weights on important edges (0.12 on average at the end). This indicates that the model misses important edges, but never includes non-important edges. For example, it could select length-3 cycles from class 1 and nothing from class 0, as shown in Figure 2. This could be considered a perfect interpretation, but it obtains a not-so-good interpretability AUC, given the expected interpretations include the length-5 cycles for each class.

A deeper analysis could be done to verify our hypotheses, challenge the batch-norm switch against alternatives, analyse results obtained with a batch-norm switch on each of our competitors or even challenge the *BA2Motifs* dataset itself, but this is out of the scope of this paper and is therefore left for future work.

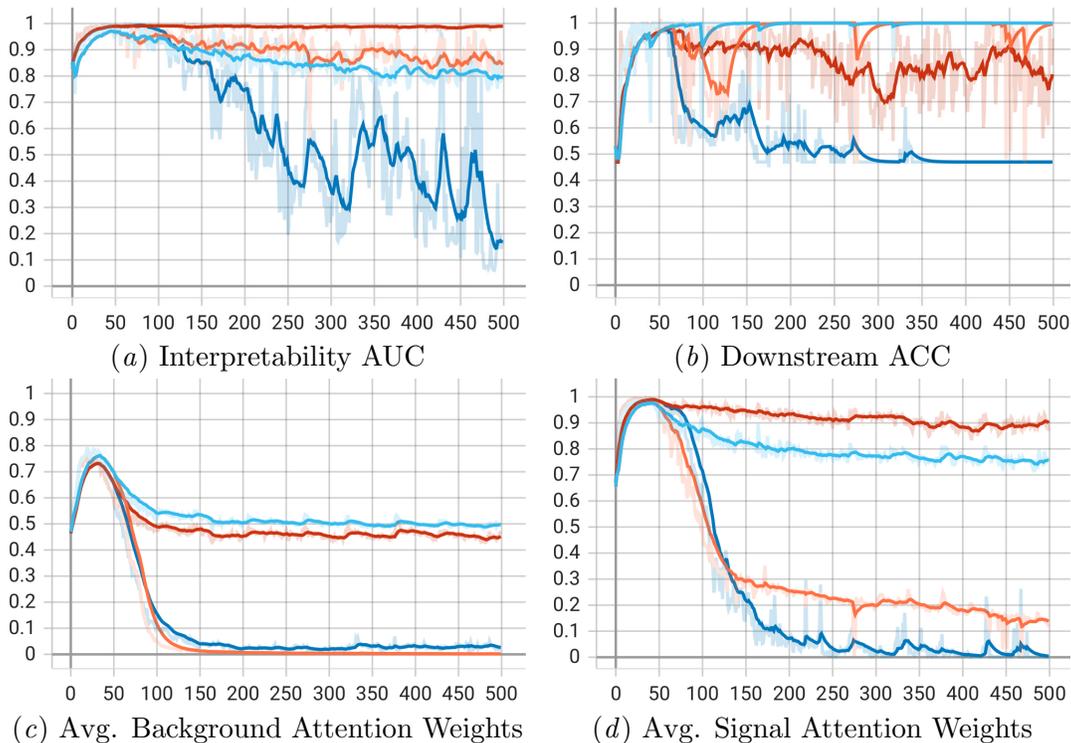


Figure 1: Comparison of GSAT with a batch-norm switch (orange and teal) and without a batch-norm switch (red and dark blue) on *BA2Motifs*. The final value of  $r$  is set to 0.1 (dark blue and orange) or 0.5 (teal and red). A smoothing effect is applied to ease the analysis. Non-smoothed results are shown in transparent.

## 2. Beyond Topology-Based Interpretations

Our results indicate that there is a link between contrastive graph augmentation and interpretability. In such approaches, however, augmentations only focus on edges, and thus interpretations only highlight topology. A natural question is whether these conclusions stand for node features too. This section presents some early work on the matter.

We enhance the framework with an additional 2-layer MLP  $\nu$ , trained to learn importance weights for node features. To make these weights understandable, we apply the feature augmentations in the input space, before being processed by the encoder.

The augmented feature matrix ( $\tilde{X}$ ) is obtained as  $\tilde{X} = (\nu(X) + \epsilon) \cdot X$  and with  $\epsilon$  sampled as indicated in Section 3.1 with  $\nu(X)$  a weight per node per feature, i.e.  $\nu(X) \in \mathbb{R}^{N \times d}$ .

We add an information loss on the importance scores to regularise the optimization, using the same  $r$  as the one for the topology information loss. As no ground truth is available for feature importance, we artificially add 20% of noisy features. We then average the importance score obtained by each feature on the test set and rank them. On *Mutag*, with our complete loss, we obtain the following results:

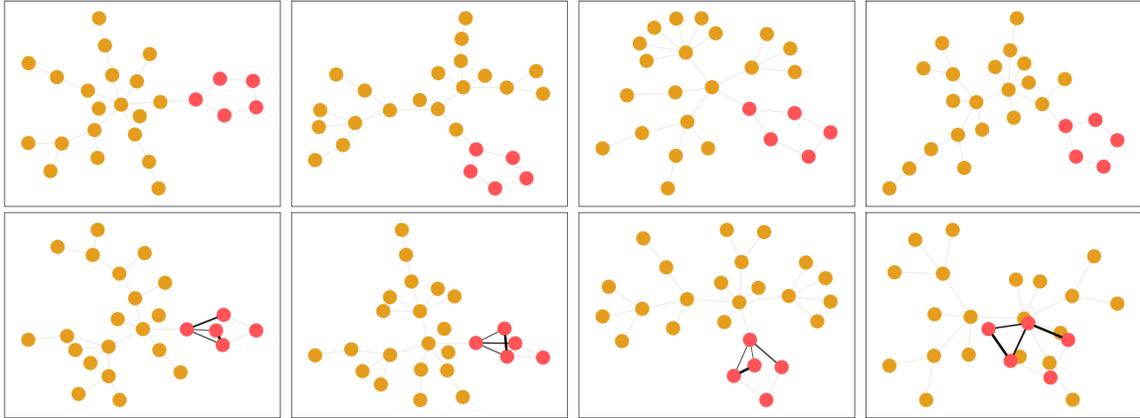


Figure 2: Interpretations of GSAT with a batch-norm switch and final  $r$  at 0.1 at epoch 370 for class 0 (top) and class 1 (bottom). Red nodes are part of the motif, and bold lines are predicted as important.

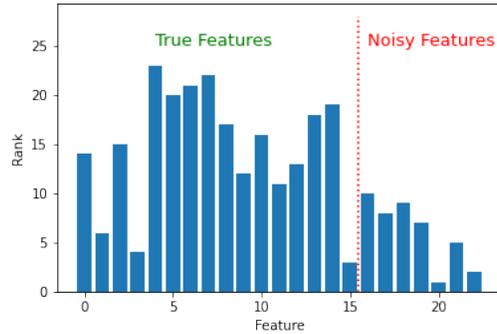


Figure 3: Ranking according to the average importance score given by the feature-augmentation module to each feature. It shows that noisy features are given less importance on average.

Figure 3 shows that the noisy features are given less importance on average, thus demonstrating an intrinsic feature selection ability which can be used as a feature importance scheme.

### 3. Impact of Sparsity on Human Readability

Figure 4(a) show some interpretations where important scores and unimportant ones are close. While truly important edges are highlighted, it is harder to distinguish them than in Figure 4(b)

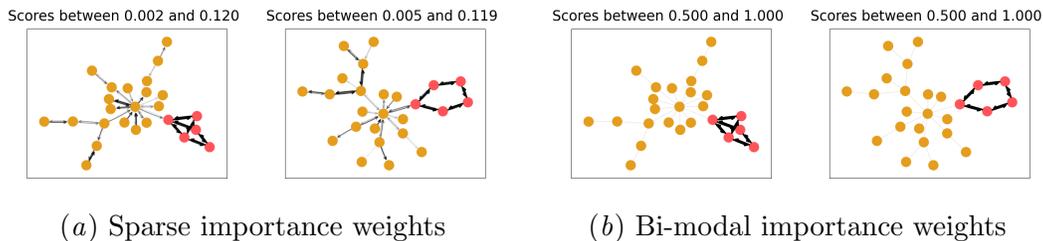


Figure 4: Examples of importance weights. A min-max normalisation is applied for clarity. Original weight limits are shown as the title of the figures.

#### 4. Ablation Study of the Watchman Module

Our framework introduces a watchman module, with the purpose to stabilize the optimization. We observe in Table 3 that the presence of this module prevents quick drops in both interpretability AUC and downstream ACC, except for a few rare cases (as visible in Table 3 and Table 1).

This pattern is either positive or not significantly different across datasets, except for a few: INGENIOUS - Negative on *BA2Motifs* and certain methods on *SPMotifs.5* which significantly decrease due to conflicting optimization objectives, however we advise to always use the watchman for our regularised loss INGENIOUS ( $\mathcal{L}$ ) as it improves the results on all datasets.

Interestingly, the loss which benefits the most from the watchman is the simple unregularised simclr loss ( $\mathcal{L}$ -info-negative), probably because in this case, the problem is under-constrained. It acts as a regularizer by enforcing the recovery of eigenvalues from the embeddings. This result shows the importance of regularisations of any kind for contrastive learning.

We conclude that while performance diminishes with overfitting as we just discussed, both the watchman and early stopping help alleviate this.

Moreover, defining the right number of training epochs is challenging in unsupervised settings. Our training curves (Figure 5) show that these unsupervised frameworks are prone to overfitting. This module can attenuate or delay this overfitting.

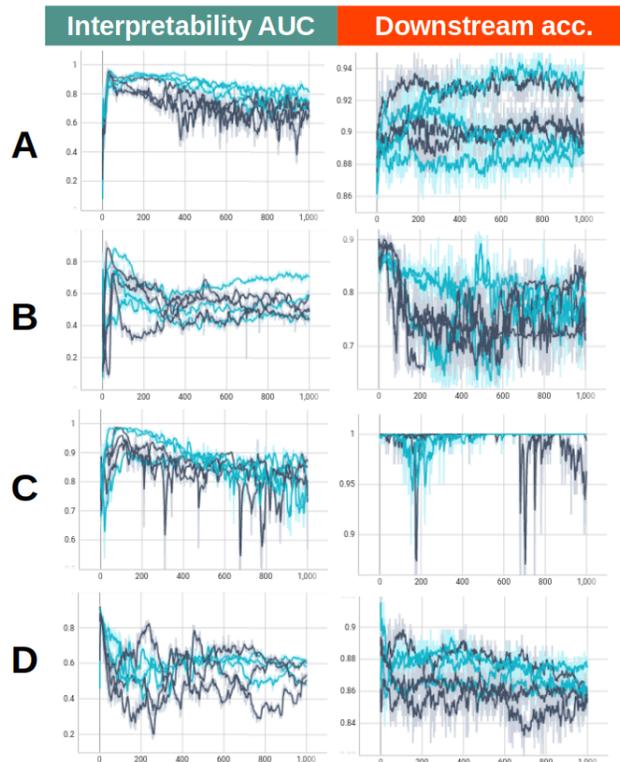


Figure 5: Example of training curves, giving interpretability AUC and downstream accuracy. Light blue is run with Watchman, dark blue without it. Those curves are: (A) for GSAT on *Mutag*, (B) for INGENIOUS on *Mutag*, (C) for INGENIOUS on *BA2Motifs*, (D) for the simple double augmentation *simclr* without regularisation on *Mutag*. The watchman helps prevent overtraining, on both supervised (GSAT) and unsupervised (INGENIOUS). The simple double-aug loss benefits too, showing that any form of regularisation is good.

## 5. Schematics and Other Figures

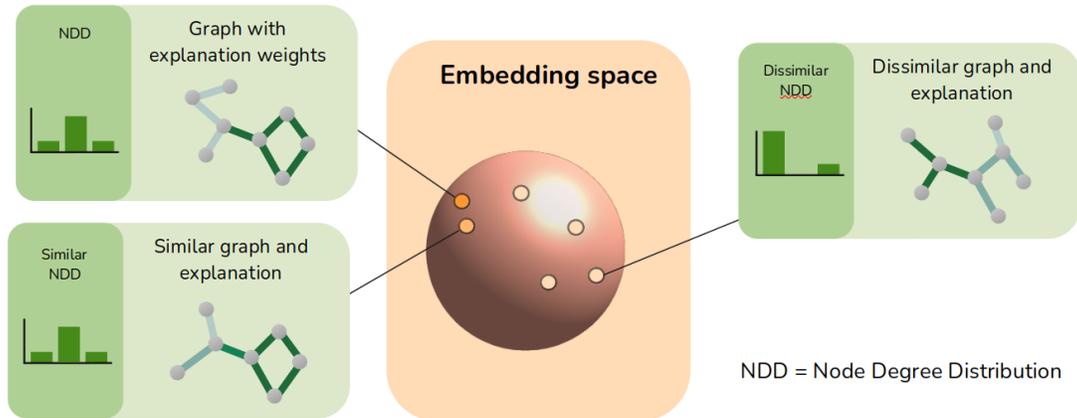


Figure 6: Principle of the Wasserstein distance metric. Our goal is to quantify the continuity of the interpretability functions. Therefore we evaluate whether similar graph embeddings obtain similar interpretations. Graph embedding proximity can be quantified by any vector similarity metric (we used Euclidian distance) and interpretation similarity can be quantified with any graph similarity metric (we use the Wasserstein distance between degree distributions).

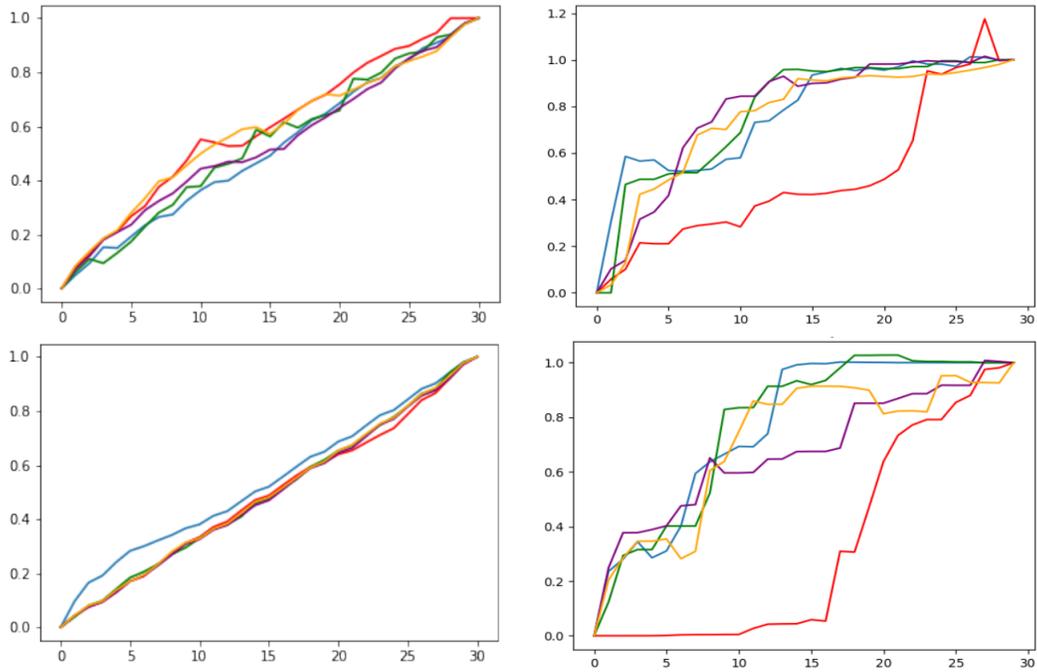


Figure 7: Examples of faithfulness curves on one graph of the *BA2Motifs* dataset. Clockwise starting from the top left: AD-GCL, GSAT, INGENIOUS, MEGA. We show the faithfulness (blue), faithfulness by removing edges according to ground truth (green), shuffled faithfulness (orange), random faithfulness (purple) and opposite faithfulness (red).

## 6. Full Tables

Wm.	Methods	Faithfulness Gap			Wasserstein Gap		
		BA2Motifs	Mutag	SPMotifs.5	BA2Motifs	Mutag	SPMotifs.5
	GSAT	.34 ± .00	.21 ± .08	.06 ± .02	.04 ± .00	.02 ± .00	.04 ± .00
	AD-GCL	.01 ± .13	.03 ± .09	.08 ± .06	.10 ± .01	.01 ± .00	.03 ± .00
	MEGA	.06 ± .02	.39 ± .12	.68 ± .03	.04 ± .01	.01 ± .00	.01 ± .00
	INGENIOUS	.39 ± .05	.24 ± .01	.50 ± 0.00	.06 ± .00	.01 ± .00	.06 ± 0.00
	INGENIOUS - Info	.54 ± .03	.58 ± .04	.46 ± .07	.02 ± .00	.01 ± .00	.05 ± .01
	INGENIOUS - Negative	.13 ± .02	.19 ± .02	.23 ± .04	.02 ± .00	.00 ± .00	.02 ± .00
	INGENIOUS - Negative - Info	.58 ± .13	.17 ± .13	.27 ± .08	.04 ± .03	.00 ± .00	.02 ± .00
✓	GSAT	.25 ± .03	.06 ± .10	-.03 ± .01	.04 ± .00	.02 ± .00	.04 ± .00
✓	INGENIOUS	.60 ± .02	.41 ± .18	.45 ± .03	.06 ± .01	.01 ± .00	.05 ± .01
✓	INGENIOUS - Info	.68 ± .02	.48 ± .03	.48 ± .06	.05 ± .02	.01 ± .00	.05 ± .01
✓	INGENIOUS - Negative	.13 ± .02	.18 ± .00	.22 ± .06	.02 ± .01	.01 ± .00	.02 ± .01
✓	INGENIOUS - Negative - Info	.61 ± .10	.15 ± .09	.32 ± .00	.03 ± .01	.01 ± .00	.03 ± .00

Table 1: Table of full results. A check in the Wm column means the model uses the watchman loss. A higher faithfulness gap means the interpretations are more faithful, a higher wasserstein gap means the embedding space is more continuous in terms of interpretations.

Wm.	Loss	Downstream ACC			Intepretability AUC		
		Cora	Tree-cycle	Tree-grid	Cora	Tree-cycle	Tree-grid
	GSAT	.610 ± .07	.985 ± .01	.984 ± .01	.000 ± .00	.630 ± .17	.844 ± .01
	$\mathcal{L}$	.572 ± .06	.848 ± .06	.954 ± .02	.000 ± .00	.232 ± .11	.532 ± .28
	$\mathcal{L}$ - Info	.593 ± .05	.777 ± .11	.952 ± .03	.000 ± .00	.256 ± .09	.588 ± .34
	$\mathcal{L}$ - Negative	.575 ± .06	.951 ± .01	.868 ± .08	.000 ± .00	.293 ± .01	.549 ± .04
	$\mathcal{L}$ - Negative - Info	.606 ± .04	.856 ± .02	.887 ± .07	.000 ± .00	.319 ± .08	.702 ± .07

Table 2: Downstream ACC and interpretability AUC on node classification. Higher is better

Wm.	Loss	Downstream ACC			Intepretability AUC		
		BA2Motifs	Mutag	SPMotifs.5	BA2Motifs	Mutag	SPMotifs.5
	GSAT	1 ± .00	.921 ± .02	.393 ± .02	.998 ± .00	.843 ± .15	.895 ± .01
	AD-GCL	1 ± .00	.902 ± .02	.337 ± .00	.378 ± .06	.416 ± .19	.472 ± .03
	MEGA	1 ± .00	.886 ± .01	.335 ± .00	.459 ± .28	.566 ± .46	.506 ± .03
	$\mathcal{L}$	1 ± .00	.900 ± .01	.366 ± .01	.910 ± .04	.745 ± .11	.467 ± .01
	$\mathcal{L}$ - Info	.990 ± .00	.889 ± .03	.361 ± .02	.860 ± .07	.606 ± .23	.491 ± .10
	$\mathcal{L}$ - Negative	1 ± .00	.887 ± .01	.340 ± .02	.745 ± .24	.546 ± .04	.494 ± .02
	$\mathcal{L}$ - Negative - Info	1 ± .00	.895 ± .01	.337 ± .01	.462 ± .38	.478 ± .23	.494 ± .04
✓	GSAT	1 ± .00	.912 ± .00	.389 ± .01	.999 ± .00	.904 ± .04	.882 ± .03
✓	$\mathcal{L}$	.993 ± .01	.900 ± .03	.363 ± .00	.959 ± .01	.771 ± .16	.510 ± .06
✓	$\mathcal{L}$ - Info	1 ± .00	.863 ± .04	.360 ± .02	.900 ± .02	.563 ± .20	.499 ± .07
✓	$\mathcal{L}$ - Negative	.995 ± .01	.895 ± .01	.338 ± .00	.148 ± .00	.620 ± .03	.581 ± .03
✓	$\mathcal{L}$ - Negative - Info	.997 ± .01	.903 ± .01	.337 ± .00	.583 ± .30	.650 ± .18	.498 ± .04

Table 3: Downstream ACC and interpretability AUC on graph classification Higher is better.

Wm.	Loss	Random faithfulness			Sparsity		
		BA2Motifs	Mutag	SPMotifs.5	BA2Motifs	Mutag	SPMotifs.5
	GSAT	.840 ± .04	.731 ± .03	.708 ± .07	.883 ± .01	.812 ± .14	.874 ± .00
	AD-GCL	.606 ± .02	.774 ± .03	.679 ± .01	.608 ± .08	.955 ± .05	1 ± .00
	MEGA	.486 ± .01	.638 ± .02	.614 ± .01	.566 ± .13	.549 ± .07	.275 ± .05
	$\mathcal{L}$	.875 ± .02	.804 ± .03	.649 ± .01	.523 ± .03	.897 ± .11	.180 ± .01
	$\mathcal{L}$ - Info	.714 ± .09	.666 ± .08	.651 ± .03	.266 ± .12	.191 ± .20	.166 ± .01
	$\mathcal{L}$ - Negative	1.020 ± .04	.859 ± .01	.647 ± .06	.991 ± .00	.980 ± .00	.714 ± .24
	$\mathcal{L}$ - Negative - Info	.897 ± .13	.879 ± .02	.641 ± .06	.417 ± .13	.958 ± .07	.657 ± .29
✓	GSAT	.817 ± .02	.822 ± .01	.538 ± .05	.855 ± .04	.893 ± .00	.848 ± .00
✓	$\mathcal{L}$	.719 ± .01	.783 ± .02	.636 ± .01	.402 ± .04	.591 ± .31	.205 ± .06
✓	$\mathcal{L}$ - Info	.743 ± .03	.670 ± .09	.606 ± .02	.288 ± .04	.269 ± .23	.168 ± .01
✓	$\mathcal{L}$ - Negative	.877 ± .03	.808 ± .02	.663 ± .07	.840 ± .02	.969 ± .00	.808 ± .31
✓	$\mathcal{L}$ - Negative - Info	.767 ± .10	.789 ± .05	.594 ± .02	.402 ± .09	.971 ± .04	.466 ± .10

Table 4: Random faithfulness and sparsity on graph classification. For random faithfulness, higher is better and for sparsity lower is better.

Wm.	Loss	faithfulness			Opposite faithfulness		
		Cora	Tree-cycle	Tree-grid	Cora	Tree-cycle	Tree-grid
	GSAT	.756 ± .01	.728 ± .21	.960 ± .03	.550 ± .07	.761 ± .05	.691 ± .01
	$\mathcal{L}$	.747 ± .03	.839 ± .07	.965 ± .07	.628 ± .13	.699 ± .08	.707 ± .05
	$\mathcal{L}$ - Info	.752 ± .03	.767 ± .03	.868 ± .06	.611 ± .10	.680 ± .17	.584 ± .11
	$\mathcal{L}$ - Negative	.732 ± .01	.853 ± .08	.902 ± .03	.531 ± .06	.846 ± .05	.610 ± .15
	$\mathcal{L}$ - Negative - Info	.783 ± .04	.780 ± .04	.847 ± .01	.477 ± .03	.608 ± .05	.598 ± .11

Table 5: Faithfulness and opposite faithfulness on node classification. For faithfulness, higher is better and for opposite faithfulness lower is better.

AUGMENT TO INTERPRET

Wm.	Loss	Random faithfulness			Sparsity		
		Cora	Tree-cycle	Tree-grid	Cora	Tree-cycle	Tree-grid
	GSAT	.693 ± .06	.814 ± .12	.911 ± .02	.708 ± .01	.827 ± .04	.801 ± .00
	$\mathcal{L}$	.727 ± .02	.823 ± .13	.877 ± .04	.754 ± .06	.629 ± .08	.639 ± .06
	$\mathcal{L}$ - Info	.727 ± .03	.757 ± .07	.816 ± .10	.711 ± .04	.481 ± .03	.565 ± .11
	$\mathcal{L}$ - Negative	.718 ± .03	.818 ± .08	.849 ± .01	.669 ± .04	.827 ± .01	.609 ± .12
	$\mathcal{L}$ - Negative - Info	.714 ± .04	.735 ± .06	.816 ± .03	.740 ± .01	.564 ± .10	.592 ± .05

Table 6: Random Faithfulness and Sparsity on node classification. For random faithfulness, higher is better and for sparsity lower is better.

Wm.	Method	$W_1$ global			$W_1$ local		
		BA2Motifs	Mutag	SPMotifs.5	BA2Motifs	Mutag	SPMotifs.5
	GSAT	.065 ± .02	.097 ± .00	.080 ± .00	.019 ± .01	.071 ± .01	.036 ± .00
	AD-GCL	.240 ± .04	.099 ± .00	.117 ± .03	.134 ± .02	.085 ± .00	.082 ± .02
	MEGA	.139 ± .04	.128 ± .01	.115 ± .00	.095 ± .03	.113 ± .01	.099 ± .01
	INGENIOUS	.154 ± .01	.099 ± .00	.113 ± .01	.093 ± .02	.087 ± .01	.044 ± .01
	INGENIOUS - Info	.144 ± .01	.096 ± .01	.105 ± .02	.122 ± .01	.084 ± .00	.047 ± .01
	INGENIOUS - Negative	.123 ± .01	.100 ± .00	.132 ± .01	.101 ± .01	.095 ± .00	.103 ± .01
	INGENIOUS - Negative - Info	.144 ± .02	.100 ± .00	.117 ± .01	.098 ± .02	.092 ± .00	.092 ± .02
✓	GSAT	.055 ± .01	.097 ± .00	.080 ± .00	.012 ± .00	.077 ± .00	.035 ± .00
✓	INGENIOUS	.143 ± .01	.099 ± .00	.095 ± .01	.080 ± .01	.084 ± .00	.044 ± .00
✓	INGENIOUS - Info	.151 ± .01	.100 ± .00	.094 ± .02	.098 ± .01	.087 ± .00	.042 ± .00
✓	INGENIOUS - Negative	.136 ± .00	.101 ± .00	.116 ± .01	.107 ± .02	.089 ± .00	.088 ± .02
✓	INGENIOUS - Negative - Info	.140 ± .00	.101 ± .00	.110 ± .01	.105 ± .01	.087 ± .00	.074 ± .01

Table 7: Global and local Wasserstein distances on graph embeddings.

Wm.	Method	Faithfulness			Opposite Faithfulness		
		BA2Motifs	Mutag	SPMotifs.5	BA2Motifs	Mutag	SPMotifs.5
	GSAT	.843 ± .02	.778 ± .05	.702 ± .06	.497 ± .02	.565 ± .04	.643 ± .03
	AD-GCL	.622 ± .07	.779 ± .00	.707 ± .02	.609 ± .06	.747 ± .09	.626 ± .05
	MEGA	.531 ± .01	.827 ± .04	.898 ± .01	.467 ± .02	.436 ± .09	.214 ± .02
	INGENIOUS	.831 ± .03	.830 ± .03	.832 ± 0.01	.432 ± .05	.588 ± .04	.325 ± 0.01
	INGENIOUS - Info	.823 ± .01	.864 ± .03	.817 ± .06	.274 ± .04	.279 ± .03	.356 ± .03
	INGENIOUS - Negative	.978 ± .05	.834 ± .02	.744 ± .03	.844 ± .03	.637 ± .00	.506 ± .07
	INGENIOUS - Negative - Info	.912 ± .04	.838 ± .04	.746 ± .01	.323 ± .12	.664 ± .10	.474 ± .10
✓	GSAT	.782 ± .02	.747 ± .07	.511 ± .04	.523 ± .02	.684 ± .04	.548 ± .05
✓	INGENIOUS	.856 ± .01	.862 ± .07	.797 ± .02	.257 ± .01	.443 ± .12	.345 ± .01
✓	INGENIOUS - Info	.884 ± .02	.819 ± .07	.796 ± .02	.200 ± .01	.336 ± .03	.316 ± .04
✓	INGENIOUS - Negative	.801 ± .00	.803 ± .01	.738 ± .02	.663 ± .02	.617 ± .02	.518 ± .08
✓	INGENIOUS - Negative - Info	.893 ± .05	.782 ± .03	.730 ± .01	.279 ± .08	.629 ± .08	.409 ± .02

Table 8: Faithfulness depending on experiments. For faithfulness, higher is better, for opposite faithfulness, lower is better.

## References

- Hang Gao, Jiangmeng Li, Wenwen Qiang, Lingyu Si, Fuchun Sun, and Changwen Zheng. Bootstrapping informative graph augmentation via a meta learning approach. In Lud De Raedt, editor, *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 3001–3007. International Joint Conferences on Artificial Intelligence Organization, 7 2022. Main Track.
- Siqi Miao, Mia Liu, and Pan Li. Interpretable and generalizable graph learning via stochastic attention mechanism. *International Conference on Machine Learning*, 2022.
- Susheel Suresh, Pan Li, Cong Hao, and Jennifer Neville. Adversarial graph augmentation to improve graph contrastive learning. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.