

1 We appreciate the valuable feedback from all the reviewers and will include the following discussion into our work.

2 **Q1. Discussion on GNNExplainer for multi-instance case (Reviewers #1 and #4).** As discussed in lines 47-52, to  
3 explain a set of instances, GNNExplainer first interprets a representative instance and then adopts ad-hoc post-analysis  
4 to generalize to multiple instances. We believe that this is not an elegant way to have a global view of the GNN model.  
5 “Since the explanatory motifs are not learned end-to-end, the model however may suffer from sub-optimal generalization  
6 performance.” Instead, the explanation in PGExplainer is generated with a parameterized network, where parameters  
7 are shared among populations. PGExplainer is natively designed for collectively explaining multiple instances.

8 **Q2. Clarification for technical detail (Reviewer #1).** The inner workings of the method (how to sample the subgraph)  
9 are mainly described by equations. Note that index notations, such as  $i$ , in different contexts may have different  
10 meanings. Using graph classification as an example, in Algorithm 2,  $i$  is the index of an instance (graph) to be explained.  
11 In line 8 of the algorithm, we iterate over all edges in the graph and each corresponding latent variable is calculated with  
12 Eq. (11). Since the index  $(i, j)$  is used to indicate an edge in Eq. (11), we don’t expand Eq. (11) to avoid confusion.  
13 The source code of PGExplainer can be found in GitHub with the name “PGExplainer”.

14 **Q3. The usage of AUC (Reviewer #1).** We follow the experimental setting in GNNExplainer. In their paper,  
15 “explanation accuracy” is not formally defined. From their source code, we find that they actually used AUC (line  
16 327, explainer/explain.py in their GitHub repository). Besides, comparing to classification accuracy, AUC is a more  
17 comprehensive measurement to evaluate binary classification.

18 **Q4. Std in Table 1 (Reviewer #1).** We didn’t report std for baselines because they don’t have sampling processes in  
19 subgraph generation. GNNExplainer doesn’t report std in their paper either. Since PGExplainer includes sampling, std  
20 is reported to show its stableness. Baselines’ stds are shown in the table below. (GRAD is deterministic with 0 std.)

21 **Q5. Discussion on**

	BA-Shapes	BA-Community	Tree-Cycles	Tree-Grid	BA-2motifs	MUTAG
ATT	0.815 ± 0.005	0.739 ± 0.015	0.824 ± 0.012	0.667 ± 0.005	0.674 ± 0.040	0.765 ± 0.013
GNNExplainer	0.925 ± 0.002	0.836 ± 0.001	0.948 ± 0.005	0.875 ± 0.012	0.742 ± 0.001	0.727 ± 0.014

22 **related work (Re-**

23 **viewer #1).** Methods in the CVPR paper mentioned only explain a specific type of GNN: GCN, on a specific task:  
24 graph classification. PGExplainer is a general model compatible with different GNNs and diverse learning tasks.  
25 Besides, instead of edge-level important scores, they just calculate node-level important scores. Thus, this paper was  
26 not cited and compared by GNNExplainer. We select a method “Gradient” in the CVPR paper which doesn’t require  
27 the global average pooling layer to calculate the importance of each node, then calculate the importance of an edge by  
28 average the connected nodes’ importance scores. The AUC scores on BA-2motifs and MUTAG are 0.773 and 0.653,  
29 respectively, much lower than PGExplainer. The KDD paper mentioned just showed up (June 3). It only applies to  
30 graph classification. Second, it only provides model-level explanations without preserving the local fidelity. Thus,  
31 the generated explanation may not be a substructure of the real input graph. Instead, PGExplainer can provide an  
32 explanation for each instance with a global view of the GNN model, which can preserve the local fidelity.

33 **Q6. Usage of “global” (Reviewers #2 and #3).** We provide an explanation for each instance, which is generated from  
34 a parameterized network. Parameters are shared among instances in the population, which equips PGExplainer with a  
35 global view of the explained GNN model. That’s why we call it a global method. For quantitative evaluation, we follow  
36 the setting in GNNExplainer to evaluate the explanation for each instance and then report the average AUC. In Fig. 1,  
37 we use a case study to show a high frequent motif in mutagens. For MUTAG dataset, we only evaluate with mutagens  
38 containing  $NH_2$  or  $NO_2$ . In Table 1, “house” motif is used as GT in BA-Community (line 287).

39 **Q7. Discussion on “local fidelity” and “general understanding” (Reviewer #2).** Both “local fidelity” and “general  
40 understanding” are important to interpret a GNN model. As discussed in [38], “local fidelity” requires an explanation  
41 corresponds to how the model behaves in the vicinity of the predicted instance. For example, in MUTAG dataset,  
42 ubiquitous motifs for mutagens include  $NH_2$  and  $NO_2$  [50,9]. Some mutagens only have  $NH_2$  and some only have  
43  $NO_2$ . Thus, it is important to give an explanation for each instance. Meanwhile, “global perspective”, or “general  
44 understanding” is important to ascertain trust in the mode [38]. However, by generating a painstakingly customized  
45 explanation for each instance **individually and independently**, it is challenging for GNNExplainer to have a “general  
46 understanding” of the GNN model, because the GNN model is trained with multiple instances. Thus, we proposed  
47 PGExplainer to preserve the “local fidelity”, at the same time, with a global view of the GNN model. Lacking of “global”  
48 view of GNN model explains why some important edges are missing in the explanations given by GNNExplainer.

49 **Q8. Our contributions (Reviewer #4).** GNNExplainer is a pioneer to provide explanations for GNN’s predictions.  
50 We include a parameterized network to enable explainer a global view of the GNN model. By collectively explaining a  
51 set of instances, PGExplainer can achieve more consistent and accurate explanations. Besides, since the parameterized  
52 network is shared among the population, PGExplainer is able to explain new instances without retraining. Empirically,  
53 PGExplainer is much more effective and efficient than the state-of-the-art method. It is non-trivial to significantly  
54 improve the accuracy performance with much less running time.

55 **Q9. Discussion on total training time of PGExplainer (Reviewer #4).** As discussed in Appendix D.1, “PGExplainer  
56 can achieve relatively good performance with a small number of trained instances, which makes PGExplainer more  
57 practical in large datasets. The results also explain why we dismiss the training time of PGExplainer and only count the  
58 inference time in Section 5” (line 541-543).