

1 We would like to thank all the reviewers for your insightful and constructive reviews and your commendation for the  
2 theoretical soundness, novelty, and effectiveness of our method. Here are the responses to some of your concerns:

3 **Reviewer #1: (1) Concerns on the quality of generated explanations.** Good point! The cohen’s kappa score of  
4 0.51 is moderate. Kappa scores for human evaluation of natural language generation tasks are generally moderate  
5 (between 0.4-0.6) as it is relatively difficult for human annotators to accurately rate machine-generated text. We also  
6 conducted another quantitative human evaluation based on pairwise comparison and observe similar results. We have  
7 presented some qualitative results (generated explanations) in the Appendix. (2) **Merge the results of supervised  
8 and semi-supervised setting.** Thanks for the suggestion. We will merge this in the revised version. (3) **Data Splits,  
9 SoTA, and Baselines.** We follow the dataset split used in [7] throughout all experiments. We chose the baselines  
10 that use the same backbone model and did not include the ones with additional resources such as knowledge bases.  
11 Our model compares similarly or slightly worse compared with the SoTA results that leverage additional resources  
12 or larger backbone models. We will add the SoTA results for reference in the revised version. (4) **Clarification for  
13 Fig2&3 .** The results in Fig2&3 are conducted on the Restaurant datasets of the ASC task and are compared against the  
14 most competitive baselines. (5) **More clarification on the experiments.** orig + rand exp refers to the model with the  
15 combination of original explanations and the corrupted explanations. The corrupted explanations (w. 80% Rand Word)  
16 is worse than the baseline model, showing the importance of good explanations. The BERT+self-training baseline  
17 denotes the baseline BERT-base model trained with the self-training method.

18 **Reviewer #2: (1) Technical Contributions.** We develop a novel EM algorithm to jointly train a natural language  
19 explanation generation module and explanation-augmented prediction module, which mutually enhance each other.  
20 Specifically, the explanation-based classifier can help improve the explanation generator as it is able to identify some  
21 high quality of explanations for training the generator (see Eq.(4)). (2) **Small samples of explanations.** The number  
22 of annotated explanations is indeed small because we want to focus on the more realistic scenario where only a few  
23 annotated explanations are available. However, the annotated explanations cover all the relation types/sentiments in  
24 RE/ASC tasks so that annotated explanations are still representative. (3) **Details on the explanation generation.** For  
25 generating NL explanation after prediction, we use the variational distribution  $q(e|x, y)$  for generating the explanation  
26 given the input sentence  $x$  and predicted label  $y$ . The same language model backbone and beam search are used for the  
27 generation.

28 **Reviewer #3: (1) Experiments on Data Sets e-SNLI [3] or Cos-E [4].** Thanks very much for pointing out these very  
29 relevant data sets. We will add the results on these data sets in the revised version. (2) **Performance w.r.t. Number of  
30 Retrieved Explanations.** This is a good point. We change the number of retrieved explanations to 5 and 15, and test  
31 them on the Restaurant dataset. The F1 are 77.3 and 75.1 respectively, which are worse than 77.8 with 10 explanations  
32 reported in this paper. We will add the sensitivity analysis w.r.t. the number of retrieved explanations in the revised  
33 version.

34 **Reviewer #4: (1) Explanations as Latent Variables.** Note that our goal lies in both generating explanations and  
35 providing extra supervision through natural language explanations. Treating explanations as latent variables allow  
36 the learning and inference process to be connected in a principled way via explanation generation and explanation-  
37 augmented prediction. The label-based explanation used during training is actually used for approximating the posterior  
38 distribution of explanations for training the explanation generator. The predictor is based on retrieved explanations that  
39 are not based on the label of the input. The final explanation of a prediction is ad-hoc (generated based on both the  
40 input and the prediction), this enables generating more informative explanations. The explanation retrieval process is  
41 important for the predictor in M-step since it conveys additional supervision. Therefore, explanations not only serve  
42 the role of explaining but also serve as additional supervision. For the human evaluation of the quality of generated  
43 explanations, a kappa coefficient of 0.51 is ok for evaluating natural language generation results as it is hard for human  
44 annotators to perfectly rate machine-generated texts. Thanks for the caveat, we admit the necessity of investigating  
45 other aspects of explainability to ensure our model’s alignment with societal desiderata and leave debiasing, risks  
46 mitigation and adversarial defense as our future work. (2) **Overclaiming on explanations:** In our generic explainable  
47 framework, we use existing explanations in [7] for simplicity. Following works like [7,13], we use “logic rules” and  
48 “labeling function” interchangeably since they convey information about labels. Our model could potentially learn to  
49 explain during testing in the example you provided if faithful natural language explanations in more complicated logical  
50 forms are provided, e.g. sent: “anything here is wonderful but the food isn’t”, exp: The word “wonderful” and “but”  
51 occur before the term “food” and “isn’t” occurs after the term “food”. We’ll clarify this issue in the final version to  
52 avoid overclaiming. Thanks for the comments and we have revised the label in Fig 1 accordingly in our next draft.