

1 We thank all the reviewers. Based on the valuable comments, we provide the following rebuttal to address the concerns
2 of each reviewer, respectively.

3 **Response to Reviewer 1** Thanks for your constructive feedbacks! Regarding your comments:

4 *Comments 1 & 4:* It turns out that RNP (The baseline in “Rationalizing Neural Predictions” proposed by Lei et. al.) with
5 class-specific generators will converge to a set of degenerated solutions, where, rather than highlighting the informative
6 words for each class, the generators will communicate the class information using trivial symbols (*e.g.* punctuation)
7 or other uninformative features (*e.g.* word position). In one degeneration example, the class-0 generator will always
8 highlight the first word and the class-1 generator will always highlight the last word. As a result, the predictor always
9 has 100% prediction accuracy, which achieves the global optimum of the collaborative target function of RNP. Our
10 empirical experiments verify that RNP with class-specific generators quickly converges to degeneration.

11 In fact, even the original RNP suffers from the degeneration problem. The additional class input merely deteriorates the
12 problem. The problem primarily results from the collaborative nature of the RNP framework. Furthermore, we are able
13 to show that the adversarial game in CAR can fundamentally resolve the degeneration problem. This is because any
14 attempt to communicate class information using uninformative symbols will be easily mimicked by the counterfactual
15 generator, and therefore would not be a good strategy for the factual generator, whose goal is to prevent mimicking.
16 This is another major advantage of CAR, which we did not have enough space to uncover in the paper. Nevertheless,
17 we will add the above discussion to the paper.

18 *Comment 2:* Max-pooling followed by a feed-forward classifier is a commonly-used and well-performed sequence
19 classifier in the sentiment classification community (*e.g.*, Yoon Kim, EMNLP 2014). Therefore we simply follow the
20 usage. Moreover, in our experiments when GRUs are used, max-pooling makes the training converge faster compared
21 to using the final hidden state, which is beneficial for the game training.

22 *Comment 3:* Eq. (10) is slightly different from [18] in the first term as we hope to better control of the lengths for
23 apple-to-apple comparison. We will add the reference in the updated version.

24 *Comment 5:* Directly highlighting top- K positions during the testing phase can be undesirable because during training
25 we used Eq. (10) to constrain on the *average* sparsity level, not the *per-passage* sparsity level. The theoretic guarantee
26 of CAR is also contingent on the average sparsity level. Per-passage sparsity stipulates that each passage has exactly the
27 same proportion of rationale, which is less theoretically and empirically sound, and which also creates a training-test
28 discrepancy if applied to testing only. Using top- K projection for training is not desirable either because the top- K ball
29 is non-convex and does not lead to good convergence.

30 *Comment 6:* RNP and POST-EXP come with their respective strengths and weaknesses. RNP is good at finding the
31 aspects that correlate most closely with the output labels, but lacks class-specific capabilities; whereas POST-EXP
32 is good at finding class-specific words and phrases, but does poorly in finding the right aspect. That is why the two
33 baselines perform differently in Tables 1 and 2. Amazon reviews (Table 1) mostly contain single aspect, but has a lot of
34 mixed reviews, and thus POST-EXP works better. Beer and hotel reviews (Table 2) contain multiple aspects, but within
35 each aspect, the reviews are less mixed, and thus RNP performs better. This can be verified by Tables 4 and 5.

36 **Response to Reviewer 2** We greatly appreciate your acknowledgment of the paper. Regarding your concerns:

37 *Direct comparison to Lei et. al. paper:* The experiment in the Lei et. al. paper is a regression task, where the output
38 scores are real values within $[0, 1]$. In our experiment, we convert the task to a classification task (because class-wise
39 rationales are only well-defined on classification tasks so far) by quantizing the real-valued scores into binary scores.
40 The quantization operation squeezes out information, and deteriorates the performance of the RNP baseline. We
41 explained this difference in the paper (line 213), but we will emphasize it more in the updated version.

42 *Inter-annotator agreement:* Each HIT was originally assigned to one crowd worker. Therefore, we rerun the experiment
43 and obtain a second set of scores. We then compute the agreement between the two sets of scores. Please note that
44 despite that the agreement score can be high or not so high, the new set of accuracy is consistent with the one reported
45 in the paper. For our method on factual rationales, the sentiment agreement is 0.71 on beer reviews, 0.79 on hotel
46 reviews and 0.81 on Amazon reviews. The aspect agreement is 0.58 on beer reviews and 0.64 on hotel reviews. The
47 aspect agreement is lower because for many cases people are making random guesses among the three options, which
48 can be confirmed by the low accuracy in some aspects (beer aroma and hotel cleanliness) in Tables 4 and 5.

49 **Response to Reviewer 4** We greatly appreciate your positive feedback! Please note that although there is only one
50 generator network for each class in our implementation, there are still two generator players for each class. The two
51 players are differentiated by feeding the ground truth label as an additional input to the generator network, and they still
52 have completely distinct target functions in the game. We will make this point clearer in our updated version.