

1 We thank the reviewers for encouraging and insightful comments. Below we respond to reviewers’ major comments.

2 **Question 1 (R1, R4, R5)** *Why study robust attribution regularization? How does it correlate with human perception?*
3 *(e.g., how about robust but totally flawed attributions?)*

4 We appreciate this question about the essential motivation of this work. To begin with, we believe that robust attribution
5 is at least *necessary* for a machine learning model to be trustworthy. Model attributions are *facts* about a model’s
6 behaviors. It is true, as **R1** pointed out, that users may still want to be skeptical when interpreting attributions, even if
7 they are robust. However in the opposite direction, one can never trust a model with *brittle* attributions.

8 Regarding human perception, a very intriguing recent paper [3] from Madry’s group showed that, empirically, adversarial
9 training (or robust prediction training) produces models whose attributions are *much more aligned with human*
10 *perception*, and seems to learn salient features from data. Since robust prediction training is a special case of robust
11 attribution training, their results form a basis for our belief that robust attribution regularization can lead to better
12 correlation with human perception (which is however a somewhat subjective question). Indeed, our visualization results
13 in the draft and appendices also corroborate their findings.

14 Note that, in the worst case, it is easy to construct models that have robust but totally flawed attributions – simply take
15 models that always have “the same” behaviors. However, the situation becomes much more complex if one imposes
16 the constraint that the model should also achieve small training error, which is what we did in this paper. Therefore,
17 robust attributions can be thought of as imposing an inductive bias to encourage learning *invariant* features from data.
18 Since human perception is essentially also related to recognizing invariances, small training error with invariant features
19 intuitively implies alignment with human perception. Next version will incorporate the above points.

20 **Question 2 (R4, R5)** *How useful is evaluating IN and CC metrics?*

21 The use of IN and CC is aligned with previous literature in studying robustness of attributions (in particular, the work of
22 Ghorbani et al. [1]). We agree with reviewers that these two metrics are only one form of brittleness, and are direct
23 consequences of our objectives. We also agree with reviewers that it is hard to quantify how increasing these two
24 metrics improves (directly) human perception. On the other hand, we think that these evaluations are still useful: (1) It
25 answers the scientific questions raised in [1], (2) Perhaps more importantly, it corroborates our analysis that robust
26 prediction training will robustify attributions as well. We will make these points explicit in the next version.

27 **Question 3 (R1)** *Section 3.2 seems redundant.*

28 We apologize for the confusion. In fact Section 3.2 fulfills an important theoretical purpose: Distributional robustness
29 approach forms a different school towards robust prediction training (see [2]). The analysis here shows that generalizing
30 these objectives to robust attributions (a much larger class) essentially still gives very similar objectives in two different
31 robust optimization models, and thus we can “safely” stick to the formulation in Section 3.1, which is reassuring.

32 **Question 4 (R1)** *Is there any functional value of regularizing an intermediate layer?*

33 Proposition 3 proves that if one regularizes by the output layer it gives a natural surrogate loss of Madry et al.’s objective
34 function, which to us makes even more sense as it directly bounds the *absolute difference* between x and x' . We believe
35 that there is more to regularizing intermediate layers and we are actively researching it.

36 **Question 5 (R5)** *More datasets, MNIST is not great.*

37 Besides MNIST we have evaluated Flower and GTSRB (traffic signs). Both are more diverse than MNIST. GTSRB is
38 practically motivated, and Flower is a high-resolution vision dataset well suited for studying attributions. We get similar
39 results in terms of both metrics and visualizations. We are actively working on more datasets and plan to include more
40 in the final version if this paper gets in.

41 **Question 6 (R5)** *More details on optimization difficulty and architectural properties.*

42 We will add more details. Roughly speaking, a main problem is *network depth*, where as depth increases we get very
43 *unstable* trajectories of gradient descent, which seems to be related to the use of second order information during robust
44 attribution optimization (due to summation approximation, we have first order terms in the training objectives).

45 References

- 46 [1] A. Ghorbani, A. Abid, and J. Y. Zou. Interpretation of neural networks is fragile. In *The Thirty-Third AAAI*
47 *Conference on Artificial Intelligence, AAAI 2019*, pages 3681–3688, 2019.
- 48 [2] A. Sinha, H. Namkoong, and J. C. Duchi. Certifying some distributional robustness with principled adversarial
49 training. In *6th International Conference on Learning Representations, ICLR 2018*, 2018.
- 50 [3] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy. In *7th*
51 *International Conference on Learning Representations, ICLR 2019*, 2019.