1   We thank the reviewers for their comments.

2   **Reply to #1:** While we fully acknowledge that the discussion on the relation to Ghorbani et al. has not made this
3 sufficiently clear, we disagree with the concerns about originality:
4 We present a method which can manipulate an image to obtain an *arbitrary* target explanation. The methods proposed
5 by Ghorbani et al, including their "targeted" method, cannot structurally reproduce a heatmap. It can only increase the
6 accumulated relevance in a certain subsection of the image. As a result, it is not capable to manipulate the heatmap
7 to closely reproduce a target explanation except for very selected and simple cases. We demonstrate this below.
8 Fine-grained control of the heatmap is absolutely essential for attacks on explanations. We also note that their attacks
9 only keep the classification result the same. This leads to significant changes in the network output (see plot d). From
10 their analysis, it is therefore not clear whether the explanation *or* the network is vulnerable (and the heatmap simply
11 reflects the relevance of the perturbation faithfully). Our method keeps the output constant which is crucial for the
12 geometrical interpretation in terms of principal curvatures and all results derived from it, i.e. all of Sec 3+4.
13 The final ms will contain a careful discussion on the relation to Ghorbani et al. and a substantially streamlined Section
14 2. Fig. 3 is moved to the SI. Also we will extend Sec 4 by a discussion of the large scale analysis of $\beta$-smoothing,
15 previously in the SI, and additional pixel flipping results (Samek et al 2017, IEEE) establishing that $\beta$-smoothing
16 performs better than unsmoothed methods (see plot for a preview).
17 *Reply to questions:*
18 • Relevance for relu networks is strongly suggested by relation to SG and indeed confirmed by pixel flipping.
19 • Beta smoothing increases pixel flipping performance.
20 • We did not retrain with softplus. We think it is preferable to modify the explanation method since it is less costly.
21 • In Fig 4, a manipulation of the unsmoothed method is "undone" by smoothing. In all other figures, the smoothed
22 method is attacked directly.
23 • If both networks have softplus non-linearities, we can compare their bound (9). Note that its constant C depends on
24 the weights of the network.
25 *Summary:* Our paper introduces a novel method allowing total control over the heatmap, it explains this manipulability
26 in terms of differential geometry and uses these insights to propose an effective defense with theoretical guarantees.
27 None of these results were contained in Ghorbani et al. We therefore strongly insist on the originality of our paper with
28 respect to Ghorbani et al.

29   **Reply to #3:** We completely agree that our algorithm is for differential explainers (remark will be added). The final ms
30 will contain a detailed comparison to Ghorbani et al (see reply to #1), acknowledging the intuition contained in the
31 figure. (Sanity Checks for Saliency Maps) is already discussed in the text. Bounding the (local) Lipschitz constant
32 of the explanation has the disadvantage that it makes the explanation insensitive to *any* small perturbation, e.g. even
33 those which lead to a substantial change in output. This is clearly undesirable as the heatmap should explain why the
34 perturbation leads to such a drastic "change of mind" of the network. Our method does not have this problem, since
35 it only bounds the same output curvature. The final ms will explain the relation to the nice work by (Alvarez-Melis,
36 Jaakkola) in detail. We fully agree with the indicated relation to PGD. The same 100 images were used for comparability.
37 We implemented all your suggestions.

38   **Reply to #4:** Since our adv attacks move on lines of the same output, they are orthogonal to conventional adv attacks
39 on the classification. They are therefore (locally) independent. We conducted experiments confirming this theoretical
40 prediction (cosine of angle between perturbations averaged over 100 images is $-1.6\times10^{-9}\pm1.6\times10^{-8}$). We added pixel
41 flipping results which show that there is indeed a trade-off between robustifying explanations and their performance.
42 For $\beta \ll 1$, the method is provably more robust but does not perform well. Our choice $\beta \approx 1$ however lies in a sweet
43 spot leading to better explanations *and* robustness. We implemented all your suggestions.
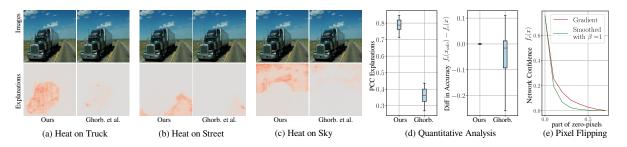


Figure 1: (a-c) Our method can structurally reproduce a heatmap, Ghorbani et al cannot. We use the same image as Ghorbani et al for comparability (d) Similarity to target map (higher is better) and change of winning-class probability for 100 images. (e) $\beta$-smoothing leads to superior pixel flipping performance (smaller is better).