

1 We thank the reviewers for their effort. We are pleased that everyone appreciated the novelty of the idea and significance  
2 of the problem setting of building locally consistent transparent models. We now address your specific concerns.

3 **R1 and R2 Long rules and interpretability:** Firstly, it is not true that the clauses would involve all the features, as  
4 mentioned in the paper, our Boolean clauses are  $2k$  sparse given  $k$  sparse local explanations i.e. PPs and PNs. So the  
5 sparsity of our clauses depends on the sparsity of the explanations which a contrastive method provides, where our  
6 contribution does not add any additional complexity. For example, this is seen on the Sky Survey dataset in the paper  
7 where the dataset has 17 features but the rules contain only 9 features. Secondly, the issue of many conjunctions and  
8 redundant features etc could arise even with known greedy algorithms for learning such structures. In fact, Rudin  
9 (citation [23] in the paper) argues that even if conjunctive rules are pages long, a person looking for an explanation can  
10 relatively easily "locate by inspection" a sufficiently sparse explanation by navigating the formula. This is cited as the  
11 reason for general desirability of rule/tree based classifiers that may not be succinct.

12 **R1, R2, R3 and R4 Scalability:** If the dataset is high dimensional, as mentioned above, our method being simply a  
13 function of the PPs and PNs which are typically sparse we should be able to scale. If the number of points is large  
14 one can do random sampling or using prototype selection methods to control size of  $F$ . More importantly though,  
15 this size is not a bottleneck since we suggest training simple models on it such as a small decision tree or L1-logistic  
16 regression which scale well even with many features and would end up choosing only a few of these rules. The most  
17 time consuming part is really obtaining PNs from a contrastive method as it is a non-convex optimization problem.  
18 Although, this too could be parallelized across different (batches of) data points.

19 **R2 Correctness of description of Eq 1:** The description is correct since,  $\delta_j$  always lies between  $x_j$  and  $b_j$  no matter the  
20 sign. So  $\delta_j$  can be  $< -|x_j|$  only if  $|b_j| > |x_j|$  and,  $b_j$  and  $x_j$  have opposite signs i.e. one is positive and the other is  
21 negative. Here again  $\delta_j$  would be closer to  $b_j$  than  $x_j$  is to it.

22 **R2 Clarity of binning:** Please see Figure 1 in supplement for a toy example depicting the whole process. We will move  
23 this into the main paper in the final version, since NeurIPS typically allows an extra page.

24 **R2 Reporting fidelity:** Yes, we will add a row in Table 1 showcasing fidelity of the transparent model to the black-box.  
25 Although, the last 5 rows in Table 1 which show Test accuracy hint towards what the fidelities might be.

26 **R3 Regarding the form of the bounds and how they are set:** The bounds  $L_j$  and  $U_j$  in Eq.(1) and Eq.(2) are set manually  
27 and its basically a bound on the domain we know. For example, features like "Age" cannot be negative and will have a  
28 lower bound of 0, while for gray scale images the values would be 0 and 1 for all pixels. There is no guarantee that one  
29 can find a PN vector for every data point using these contrastive methods. However, that is independent of our current  
30 contribution and our method can create rules even if we have no PNs. This can be confirmed by looking at the Magic  
31 and Diabetes datasets in Table 1 where we can still create rules eventhough no PNs were found for those datasets.

32 **R3 Why not Knowledge distillation:** Please look at Table 1 where we perform favorably to knowledge distillation.

33 **R3 Regarding Eq 3:** By the definition of PN, we are aiming to generate a point  $n(x)$  that is of a different class from the  
34 original input  $x$  and thus equation is correct.

35 **R3 Only one PP used in Eq. 5:** Contrastive methods (citations [8] and [31] in the paper) output only a single PP and at  
36 most one PN explanation for each data point. We use this PP and PN in Eq 5 for discretization.

37 **R3 Regarding circular argument:** The local consistency metric just ensures that the predictions as well as the local  
38 explanations (PPs and PNs) for the black-box are consistent with the transparent model. This is independent of how one  
39 might obtain a locally consistent model. When trying to build a such a model one is free to use all the information that is  
40 available to them. So we do not believe there is any circular argument. In fact, the Augmentation baseline we compare  
41 with in Table 1 is more closely related to just directly using PPs and PNs, and we outperform it by a significant margin.

42 **R3 Dimension of PPs and PNs in Eq (5):**  $p(x_j)$  means the  $j$ -th coordinate of the PP vector corresponding to the point  $x$ .  
43 So there is no discrepancy between this and Listing 1. We will clarify.

44 **R4 Selecting between PNs:** Data is typically normalized where PPs and PNs are generated by another method and  
45 given as input to ours. So though you have a valid point that is not the focus of this submission, rather we assume  
46 whichever local contrastive method we use has already done that for us and provided a valid local explanation. Realistic  
47 explanations in previous works were generated using autoencoders etc. (see citation [8] in the paper).

48 **R4 Assumption that  $x$  and the base vector share the same class:** This is not true. There is only one base vector that is  
49 set for the entire dataset. This may be a zero vector or one where the base value for each feature is set based on domain  
50 knowledge (see citation [9] in the paper). Given this, the limitations you mentioned do not apply.

51 **R4 More sensible to look at  $C(x, y) - C(\delta, y)$ :** Interesting suggestion! Although, the issue you mention can be addressed  
52 (indirectly) in Eq. 1 by having a large enough  $\kappa$ , which controls the margin between  $y$  and the next most likely class.